# Introduction to R for Public Health Investigations

Workbook for day 3

# Contents

## Practical Exercise

**Instructions**

Learners are provided with a scenario, questions, and tasks with associated code to perform each task. We recommend:

Novice users (Boatswains): Use this workbook and R script(s) provided on GitHub. Using this workbook as a guide, run the code we've provided piece by piece to understand what each chunk of code does, and what various functions are doing. At this point don't worry too much about being able to write or debug code.

Beginner/Intermediate users (First Mates): R code is provided as a screen capture image in this workbook. You should have sufficient understanding of coding to get a general sense of what the code is doing by reading it (with the assistance of the help documentation, a few Google searches as needed, and comments in the R scripts we've provided on GitHub). It is our intention to have you write the code out from the guide as you progress through the scenario. Cross reference to the R script(s) provided on GitHub if you encounter any tangly problems.

Advanced users (Master Mariners): We encourage you to try writing your own code where you like and contrast it with the code used for the exercise, and to help your peers as questions arise. Cross reference to the R script(s) provided on GitHub if you encounter any tangly problems.

Get as far as you can with this exercise within 4 hours and 10 mins (maximum). Don't worry if you need extra time. The learning curve for R is steep and learners will benefit most from

dedicated time for practicing. Reach out to your course facilitators by Slack or by email if you require assistance with the course material.

## Introduction

In this session you will apply your R and data management skills to prepare and collate data from multiple datasets, apply a case identification algorithm, and conduct a basic descriptive epidemiological analysis. You have been provided with five datasets which are mock extracts from administrative datasets (physician billings, discharges from hospital, and provincial health insurance client roster) and have been simplified for training purposes. These datasets must be combined in order to efficiently apply administrative case definitions (algorithms) so that you can describe the occurrence of health service use for chronic conditions among residents of a small province. You've been provided with the following datasets and case definition:

Physician billings: The msp1213 and msp1314 datasets contain information from fee-for-service physician billings (general practitioner and specialists) in the 2012/13 and 2013/14 fiscal years respectively. Each record equates to a physician visit and a date on which the visit occurred and a single diagnosis code. While persons may have multiple conditions to discuss during their visit with the physician, the physician selects a "most responsible" condition for billing.

Discharges from hospital: The dad1213 and dad1314 datasets contain information from the discharge abstract produced upon discharge from hospital. The data includes acute care hospitalisations only as other hospitalisations, ambulatory or long-term care hospital services are excluded. Multiple diagnoses are included on the discharge abstract, along with the dates on which the individual was admitted to and discharged from the hospital.

Provincial health insurance client roster: The reg1314 dataset contains a listing of every person who is registered for provincial health insurance (e.g. MSP, OHIP, MCP, etc.). Enrolment is automatic upon issue of a birth certificate among residents of the province (in the case of new births) or applied for upon moving to the province. Registration is cancelled upon relocation and registration for provincial health insurance in a new province or death. Publically funded health care services are billed to the province of residence, or federal program providing health insurance coverage, or must be paid for out of pocket or by private insurance by those who are not registered. Regardless, persons who access health care services and who aren't registered with the provincial Medical Services Plan have been excluded from this dataset along with registrants who do not reside in the province.

Table 1: Contents of datasets provided for Introduction to R, Day 3

| Dataset | Fields | Type | Definition |
|---|---|---|---|
| msp1213.csv | uniqueid | String | Unique identifier (person) |

| Dataset | Fields | Type | Definition |
|---|---|---|---|
| | servdat | Date | Date on which physician visit occurred |
| | diag1 | String | Diagnosis (ICD9) code provided on billing |
| msp1314.csv | uniqueid | String | Unique identifier (person) |
| | servdat | Date | Date on which physician visit occurred |
| | diag1 | String | Diagnosis (ICD9) code provided on billing |
| dad1213.csv | uniqueid | String | Unique identifier (person) |
| | admitdat | Date | Date admitted to hospital |
| | sepdat | Date | Date discharged from hospital |
| | diag1 | String | Diagnosis (ICD9) code listed on discharge |
| | diag2 | String | Diagnosis (ICD9) code listed on discharge |
| | diag3 | String | Diagnosis (ICD9) code listed on discharge |
| dad1314.csv | uniqueid | String | Unique identifier (person) |
| | admitdat | Date | Date admitted to hospital |
| | sepdat | Date | Date discharged from hospital |
| | diag1 | String | Diagnosis (ICD9) code listed on discharge |
| | diag2 | String | Diagnosis (ICD9) code listed on discharge |
| | diag3 | String | Diagnosis (ICD9) code listed on discharge |
| reg1314.csv | uniqueid | String | Unique identifier (person) |
| | dob | Date | Date of birth |
| | sex | Numeric | F = 1, M = 2 |

**CCDSS Asthma definition**[1]: The case definition of diagnosed asthma is: an individual aged one year and older having at least two visits to a physician with a diagnosis of asthma in the first diagnostic field in a two-year period, or at least one hospital separation with a diagnosis of asthma ever in any diagnostic field, coded by the International Classification of Diseases (ICD), ninth revision or ICD-9-CM 493 or ICD-10-CA J45-46. Note that the application of this case definition in this exercise is highly simplified.

---

[1] https://www.canada.ca/en/public-health/services/publications/diseases-conditions/asthma-chronic-obstructive-pulmonary-disease-canada-2018.html#a1.2.1

## 1. Set up your workspace

1. Create new folders on your computer to organize the files for today's exercises similar to how they were set up for exercise 1 and 2:
   a. Within the IntoToR folder, create a subfolder for Day 3 called: "Exercise_Day3".
   b. Create new folders within Exercise_Day3 called:
      i. "data";
      ii. "scripts"
      iii. "output"
   c. Move the files for Day 3 from to their corresponding folders.

   **Note: for Day 3 we'll be working from a single .Rmd file in your scripts folder for this exercise, rather than writing several individual scripts.**

   Tips:

   - Use the same RProject from exercise 1 and 2
   - Clean up your environment in RStudio if you have previous work still in it with the following code: rm(list=ls())
   - Remember to save your work regularly

2. Now that you have set up your project folders:

   Open the datasets in excel (msp1213.csv, msp1314.csv, dad1213.csv, dad1314.csv, reg1314.csv).

   a. Is uniqueid really unique? What do you think each record represents?

   b. What field(s) do you think might be useful for linking data across tables?

   c. Refer back to the administrative case definition. What datasets and fields do you think might be needed in identifying cases of asthma and their case dates?

d. Given the administrative case definition and the data you have, what descriptive epidemiological calculations would you like to perform to describe the occurrence of asthma in this population?

e. Considering that reg1314 is a Medical Services Plan client registry for the 2013/14 fiscal year, do you think that this is acceptable for use as a population registry? Why or why not? What limitations might there be if this were to be used as a population registry?

f. What additional information do you need from the other four tables to add to the reg1314 dataset to carry out your descriptive analysis of asthma?

g. How does this information need to be laid out as data fields in the reg1314 dataset so that you can perform your descriptive analysis? What steps do you need to take during data processing to include these data fields in the reg1314 dataset?

*Hint – draw the resulting file. Having your starting point (i.e. datasets provided) and end point (i.e. ultimate flat file) clearly in mind will allow you to identify the specific steps you need to take during data processing.*

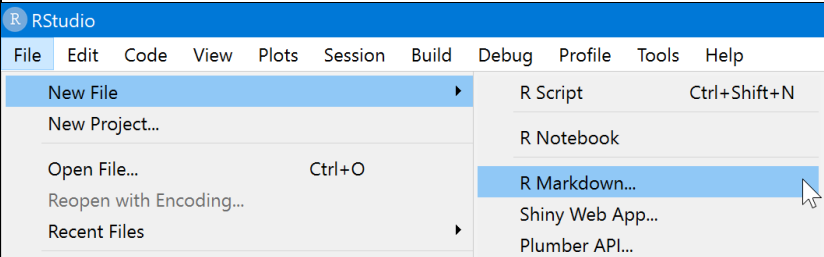In this exercise, the steps below will be taken to create this flat file:

| uniqueid | dob | age | sex | asthma_case | asthma_casedate | popn_1plus |
|----------|-----|-----|-----|-------------|-----------------|------------|
| 1 | 1936-09-21 | 76 | 1 | 0 | NA | 1 |
| 2 | 1916-04-08 | 97 | 2 | 1 | 2012-09-24 | 1 |
| 3 | 1924-10-17 | 88 | 1 | 0 | NA | 1 |
| 4 | 1995-04-01 | 18 | 2 | 0 | NA | 1 |
| 5 | 1997-06-28 | 16 | 1 | 0 | NA | 1 |

Overview of steps to be taken during this exercise:
- Append DAD
- Flag hospital separations with a diagnosis of asthma
  - keep only asthma-related hospitalisations
  - identify duplicates
  - flag and count duplicate records for each uniqueid, keep only earliest sepdat
- Save DAD data as new appended file
- Append MSP
- Flag physician visits with a diagnosis of asthma
  - keep only asthma-related physician visits
  - identify duplicates
  - flag and count duplicate records for each uniqueid, keep only earliest servdat
- Save MSP data as new appended file
- Merge registry to DAD and MSP data
  - Merge DAD and Registry (full join, uniqueid)
  - Merge MSP and Registry (full join, uniqueid)
  - Merge two resultant files (full join, uniqueid dob sex)
- Save as new registry file with chronic disease status information
- Apply case definition: Flag persons who meet the case definition
- Calculate case date (earliest health service coded for asthma)
- Prepare registry for descriptive epi calculations
  - Calculate age at midyear 2013
  - Flag persons who are part of the denominator
  - Apply age group categories
- Analyse the data
  - Sum case definition (numerator) flag and denominator flag, calculate rate
  - Sum case definition (numerator) flag and denominator flag (by agegroup and/or sex), calculate rates
  - Create figures

## 2. Set your RStudio session

Note: Try adding your own comments to the code to keep track of critical things like counts of records and variables (to ensure operations are behaving as expected) and to clarify the code for your understanding (to ensure that you could understand weeks from now what you produced through the course of this exercise). Note that there is a master rmd file provided for this exercise (with comments) to help you should things get tangly. As this exercise is complicated, please do not hesitate to refer to the rmd file provided to run through the exercise, whether you look to develop your understanding by following along with the master rmd file, copying and pasting code into your own rmd file, or typing your own code out from the images provided into your own fresh rmd file. Regardless of your approach to this exercise, it is critical that you pay attention to changes in your files (number of variables and observations for instance) to ensure linkages behave as expected.

| Task | Code |
|---|---|
| Open a new R Markdown file and select HTML as the default output format and call it: "Day3_Exercise_script.Rmd" | ![RStudio File > New File > R Markdown menu screenshot showing options: R Script (Ctrl+Shift+N), R Notebook, R Markdown..., Shiny Web App..., Plumber API...] |
| Note that for the Day 3 exercise, we'll be operating out of a single R-markdown file rather than sourcing several files all at once and treating our R-markdown file as a notebook in HTML format. First, you'll need to | ```<br>---<br>title: "Day3_Exercise_Script"<br>output: html_document<br>---<br>``` |

| | |
|---|---|
| specify this format in the YAML header at the top! | |
| Next, you'll need a code chunk to setup your document for analyses. Here, we'll directly load the libraries we plan to use, initialize the location of the present Rmd script, and set up the location of out 'output' folder pathway for use later. | ```### Setup```<br>```{r setup, include=FALSE}```<br>`require("knitr")`<br>`library(tidyverse)`<br>`library(here)`<br><br>`knitr::opts_chunk$set(echo = TRUE)`<br>`here::i_am("Exercise_Day3/scripts/Exercise_Day3.Rmd")`<br>`output_folder <- here::here("Exercise_Day3", "output")`<br><br>``` ``` ``` ``` |

3. Prepare the data:

| | |
|---|---|
| Load hospitalisation datasets (dad1213.csv and dad1314.csv) and store them as objects in memory. Note the number of observations and variables in each.<br><br>Note: If you need to add additional chunks to your R Markdown file you can either click on Insert  -> R at the top of the Markdown toolbar or | ```{r}```<br>`dad1213 <- read_csv(here("Exercise_Day3", "data", "dad1213.csv"))`<br>`dad1314 <- read_csv(here("Exercise_Day3", "data", "dad1314.csv"))`<br>``` ``` ``` ``` |

Unclassified / Non classifié

**Introduction to R - Workbook for day 3**

| | |
|---|---|
| use the keyboard shortcut Ctrl + Alt + I ( Cmd + Option + I on macOS) | |
| Append hospitalisation datasets (using bind_rows). How many hospitalisations are there in each file?<br><br>Confirm separation dates fall in the correct fiscal year. How many records need to be dropped?<br><br><br>Note: If you only want to run part of the chunk of code, use the Run dropdown at the top of the markdown toolbar. | `##### Append DAD datasets:`<br><br>```{r}<br>dad_appnd <- bind_rows(dad1213, dad1314) %>%<br>  mutate(sepdat = as.Date(sepdat, "%d-%b-%Y")) %>%<br>  filter(sepdat > "2012-04-01" & sepdat < "2014-03-31")<br>```<br><br>Hint: try entering print(unique(dad_appnd$sepdat)) in the command window<br><br>Or, view and sort the file by double clicking on the file in the environment panel and clicking on the set of triangles next to the variable name:<br><br>Filter<br><table><tr><td>uniqueid</td><td>admitdat</td><td>sepdat</td><td>diag1</td><td>diag2</td><td>diag3</td></tr><tr><td>1779</td><td>29-Mar-2012</td><td>2012-04-02</td><td>1011</td><td>9500</td><td>2790</td></tr><tr><td>3016</td><td>02-Apr-2012</td><td>2012-04-02</td><td>1080</td><td>6230</td><td>8850</td></tr><tr><td>1934</td><td>02-Apr-2012</td><td>2012-04-02</td><td>1189</td><td>1100</td><td>NA</td></tr></table> |

11

Flag and keep hospitalisations where a diagnosis of asthma was given. As we are interested in the whole diagnosis code block 493, we will take a substring of that code to filter upon.How many asthma-related hospitalisations were there in 2013/14 and 2012/13?

To practice: how many unique persons (hint: file size at the end of this code)? What is the maximum number of hospitalisations per person during the period of interest (hint: view and sort, or table or count functions)?

```r
##### Process DAD Datasets:
```{r}
#create flag variable.
####Step 1: Extract first three digits of diag code.
dad_appnd <- dad_appnd %>% mutate(diag1_sbstr = str_sub(diag1, end=3)) %>%
  mutate(diag2_sbstr = str_sub(diag2, end=3)) %>%
  mutate(diag3_sbstr = str_sub(diag3, end=3))

####Step 2: create flag where 0 = no asthma, 1= asthma if any diag variables equal to 493.
dad_appnd <- dad_appnd %>% mutate(flag = case_when(diag1_sbstr == "493" | diag2_sbstr =="493" | diag3_sbstr
=="493" ~ 1, TRUE ~ 0))
```

As there are individuals who have been hospitalised more than once in the time period, keep only the earliest separation date. How does the distinct function work to remove duplicates? Do the data need to be sorted before using this function? Should the data be sorted ascending or descending, and on which variable(s)?

Keep and rename only the variables you need for later (i.e., uniqueid, sepdat, count_duplicates)

Tidy the data and work environment.

```r
#keep unique asthma-related hosps
#####Step 1: Keep only asthma related events.
#####Step 2: Group by unique ID (a) and count how many uniqueIDs appear in each group (b)
#####Step 3: Keep only earliest sepdat where  unique persons are duplicated.
#####Step 4: Keep (a) and rename (b) only the variables you need for your analysis.
dad_appnd_asthma <- dad_appnd %>%
  filter(dad_appnd$flag == 1) %>% #Step 1
  group_by(uniqueid) %>% #Step 2 (a)
  add_tally(name = "count_duplicates") %>% #Step 2 (b)
  arrange(uniqueid, sepdat) %>% #Step 3 (a)
  distinct(uniqueid, .keep_all = TRUE) %>% #Step 3 (b)
  select(uniqueid, sepdat, count_duplicates) %>% #Step 4 (a)
  rename(dad_sepdat_asthma = sepdat, dad_count_asthma=count_duplicates) #Step 4 (b)

#tidy workspace
rm(dad_appnd, dad1213,dad1314)

```
```

*Note: The number of records remaining should equal the number of unique individuals obtained above*

| | |
|---|---|
| Append physician visit datasets (msp1213.csv and msp1314.csv). How many physician visits are there in each file? How many in the appended file? | `#### Prepare MSP datasets`<br><br>```r<br>{r}<br>##load files<br>msp1213 <- read_csv(here("Exercise_Day3", "data", "msp1213.csv"))<br>msp1314 <- read_csv(here("Exercise_Day3", "data","msp1314.csv"))<br><br>##append datasets<br>msp_appnd <- bind_rows(msp1213, msp1314)<br>```<br><br>*Tip: If writing from scratch, consider recycling DAD code* |
| Similar to what we did for the DAD datasets, we now need to process the MSP datasets.<br><br><br>For this we want to reformat servdat to a date-class variable to be used for filtering; keep dates between Apr 1, 2012 and Mar 13, 2014; then use the physician codes starting with 493 to create a flag for asthma. | *##### Process MSP Datasets:*<br><br>```r<br>### process MSP.<br>########Step 1 - Reformat servdat to date-class variable.<br>########Step 2 - Filter by servdat: keep if servdat between Apr 1 2012 & Mar 13 2014<br>########Step 3 - Extract first three digits from diag1.<br>########Step 4 - Create new flag variable, 1=asthma, 0=no asthma; code block 493 (incl 4390...4939)<br><br>msp_appnd <- msp_appnd %>%<br>  mutate(servdat = as.Date(servdat, "%d%b%Y")) %>% #Step 1<br>  filter(servdat >= "2012-04-01" & servdat <= "2014-03-31") %>% #Step 2<br>  mutate(diag1_sbstr = as.numeric(str_sub(diag1, end=3))) %>% #Step 3<br>  mutate(flag = case_when(diag1_sbstr == "493" ~ 1, TRUE ~ 0)) #Step 4<br>``` |

| | |
|---|---|
| Drop all non-asthma related physician visits. How many observations remain (hint: file size)? What are the maximum number of physician visits coded for asthma per person (hint: view and sort, or table and count functions)?<br><br>As there are individuals who have seen a physician more than once in the time period, keep only the earliest visit date. How many unique individuals have seen a physician where their visit was coded for asthma (hint: view and sort, or table and count functions)? | ```r<br>#Step 1 - Keep & count phys visits for asthma<br>#Step 2 - Identify and count how many observations per uniqueid<br>#Step 3 - keep only earliest servdat where unique persons are duplicated<br>msp_appnd_asthma <- msp_appnd %>% filter(msp_appnd$flag == 1) %>% #Step 1<br>  group_by(uniqueid) %>% add_tally(name="count_duplicates") %>% #Step 2<br>  arrange(uniqueid, servdat) %>% distinct(uniqueid, .keep_all = TRUE) #Step 3<br>``` |
| Tidy the data and work environment<br><br>Keep and rename only the variables you need for later (i.e. uniqueid, servdat, count_duplicates) | ```r<br>#Tidy the data.<br>msp_appnd_asthma <- msp_appnd_asthma %>%<br>  select(uniqueid, servdat, count_duplicates) %>%<br>  rename(msp_servdat_asthma = servdat, msp_count_asthma = count_duplicates)<br><br>#clear work environment.<br>rm(msp_appnd, msp1213, msp1314)<br>```<br>``` |

### 3. Clean, process, and assemble the data

| Task | Code |
|------|------|
| Load reg1314. How many registrants are there? | ```###### Registry Dataset```<br>` ```{r} `<br>`#####Load registry file`<br>`reg1314 <- read_csv(here("Exercise_Day3", "data", "reg1314.csv"))` |
| In order to assemble MSP, DAD, and registry data, what type of merge(s) will you need? How will you merge all three files (i.e., what key variable or variables) and in what steps? Provide your justifications.<br><br>Tip: Review the join types in your participant guide | |

| | |
|---|---|
| Merge hospitalisation data with the registry. How many records are there in the merged datafile? How many registrants in 2013/14 were hospitalised where a diagnosis of asthma was given in the previous two years? | ````##### Merge the dad dataset with the registry:`<br>```` ```{r} ````<br>`#--- Merge dad to reg`<br>`asthma_dad_reg1314 <- full_join(reg1314, dad_appnd_asthma, by=c("uniqueid"="uniqueid"))`<br><br><br>`#--- Frequency of each result in dad_count_asthma`<br>`asthma_dad_reg1314 %>%`<br>`  count(dad_count_asthma)`<br><br>` ``` `<br><br><br>*Hint: You can find the count of observations in the resulting merged file in the environment window. View the merged file and see how many individuals have data pertaining to asthma hospitalisations. The number of persons should equal the number of persons in the dad_appnd_asthma file. The merged file also contains records pertaining to everyone else in the registry.* |
| Merge physician visit data with the registry. How many records are there in the merged data file? How many registrants in 2013/14 visited a physician where the billing was coded for asthma in the previous two years? | ````##### Merge the MSP dataset with the registry:`<br>```` ```{r} ````<br>`#--- Merge msp to reg`<br>`asthma_msp_reg1314 <- full_join(reg1314, msp_appnd_asthma, by=c("uniqueid"="uniqueid"))`<br><br>`asthma_msp_reg1314 %>%`<br>`  count(msp_count_asthma)`<br><br>` ``` `<br><br><br>*Hint: You can find the count of observations in the resulting merged file in the environment window. View the merged file and see how many individuals have data pertaining to asthma visits. The number of persons should equal the number of persons in the msp_appnd_asthma file. The merged file also contains records pertaining to everyone else in the registry.* |

| | |
|---|---|
| Combine asthma datasets. How many observations are there in the resulting file? How many people are found in both datasets? | ```r<br>##### Merge Asthma-Registry datasets<br>```{r}<br>#create combined asthma datasets<br>asthma_reg1314 <- full_join(asthma_dad_reg1314, asthma_msp_reg1314, by=c("uniqueid"="uniqueid", "sex"="sex", "dob"="dob"))<br><br>table(asthma_reg1314$msp_count_asthma, asthma_dad_reg1314$dad_count_asthma)<br>count(asthma_reg1314, msp_count_asthma)<br>count(asthma_dad_reg1314, dad_count_asthma)<br><br>```<br><br>*Note: As sex and dob come from same file, mismatches would indicate bad things have happened - don't want duplicated columns.*<br><br>*Hint:*<br><br>- *You can find the count of observations in the resulting merged file in the environment window.*<br>- *Using the table command, we can see a cross tabulation of count of hospitalisations by physician visits. This table will illustrate persons who have shown up in both files, and what their pattern of health service utlisation for asthma was in the period of interest.*<br>- *Using the count command we can see the number of persons in the dataset broken down by count of physician visits for asthma or hospitalisations for asthma.* |

| | |
|---|---|
| How many individuals meet the hospitalisation and physician component of the case definition for asthma? What are the other components of this definition? | ````` ##### Process Asthma-Registry dataset ```{r Flag individuals who meet hospitalisation and physician vist portion of the case definition}  # 1H Ever or 2P in 2Y asthma_reg1314 <- asthma_reg1314 %>%   mutate(PH_flag = case_when(dad_count_asthma >= 1 ~ 1,                             msp_count_asthma >= 2 ~ 1,                             TRUE~0))  asthma_reg1314 %>%   count(PH_flag) # table(asthma_reg1314$PH_flag) ``` ` |
| Calculate case age in 2013/14. What reference date would you use for this calculation? | ```` ```{r Calculate age}  asthma_reg1314 <- asthma_reg1314 %>%   mutate(dob = as.Date(dob, "%d%b%Y")) %>%  # format as date   mutate(age = as.numeric(difftime(time1 = as.Date("01jul2013", "%d%b%Y"), time2 = dob, units="days"))/365.25) ` |
| How will you identify who should be included in the case definition based on age? Remember you will need to obtain denominator data from the registry for your descriptive epi calculations. | ``` #--- Identify asthma cases based on age and case definition #143 cases asthma_reg1314 <- asthma_reg1314 %>%   mutate(popn_1plus = case_when(age >= 1 ~ 1,                                 TRUE ~ 0)) %>%   mutate(asthma_case = case_when(popn_1plus == 1 & PH_flag == 1 ~ 1,                                  TRUE~0)) ``` |
| How many individuals met the case definition for asthma in 2013/14? | ``` asthma_reg1314 %>%   count(asthma_case) ``` |

Specify the case date, or the date of the earliest health service coded for asthma. We will do this by creating a series of flags based on hospital separation date and physician visit date. We can use these flags to identify and subsequently recode to specify which date to use as the case date, or the earliest point a person met the case definition.

Check your logic statements through appropriate tabulations and cross tabulations (i.e., see count functions to the right and in the rmd file for this exercise).

```r
```{r Process when case met definition}

#identify case date, or the earliest point at which a person met the case definition
asthma_reg1314 <- asthma_reg1314 %>%
  mutate(date_flag = case_when(
    asthma_case ==1 & dad_sepdat_asthma < msp_servdat_asthma ~ 1,
    asthma_case ==1 & dad_sepdat_asthma > msp_servdat_asthma ~ 2,
    asthma_case ==1 & dad_sepdat_asthma == msp_servdat_asthma ~ 3,
    asthma_case ==1 & is.na(dad_sepdat_asthma) == TRUE ~ 4,
    asthma_case ==1 & is.na(msp_servdat_asthma) == TRUE ~ 5,
    TRUE~0))


asthma_reg1314 <- asthma_reg1314 %>%
  mutate(asthma_casedate = case_when(
    date_flag == 1 ~ dad_sepdat_asthma,
    date_flag == 2 ~ msp_servdat_asthma,
    date_flag == 3 ~ msp_servdat_asthma,
    date_flag == 4 ~ msp_servdat_asthma,
    date_flag == 5 ~ dad_sepdat_asthma))


#count asthma case flags
asthma_reg1314 %>%
  count(asthma_case)

#count date flags
asthma_reg1314 %>%
  count(date_flag)

#compare case and date flags
asthma_reg1314 %>%
  count(asthma_case, date_flag)

#compare casedate and case flags
asthma_reg1314 %>%
  count(asthma_casedate, date_flag)
```

Does the count of asthma cases count reflect incident cases? Why or why not?

| | |
|---|---|
| Tidy the final product and workspace as needed.<br><br>Keep only the variables you need for later: uniqueid, dob, age, sex, asthma_case, asthma_casedate, popn_1plus | <pre>#---Tidy data<br>str(asthma_reg1314) #review structure of the data<br><br>#select variables to keep<br>asthma_reg1314 <- asthma_reg1314 %>%<br>  select(uniqueid, dob, age, sex, asthma_case, asthma_casedate, popn_1plus)<br><br>#clear out work environment<br>rm(asthma_dad_reg1314, asthma_msp_reg1314, dad_appnd_asthma, msp_appnd_asthma, reg1314)<br>```</pre> |

### 4. Analyze and visualize the data

| Task | Code |
|------|------|
| What was the crude asthma prevalence rate in 2013/14?<br><br>Note – code duplicates in this block for demonstrative purposes related to chaining statements together. | ```#### Data Analysis```<br><br>```` ```{r Calculate crude prevalence} ````<br><br>```#crude prevalence:```<br>```#add up n cases and n population```<br>```asthma_overall <- asthma_reg1314 %>%```<br>```  summarize(cases = sum(asthma_case), popn = sum(popn_1plus))```<br><br>```asthma_overall <- asthma_reg1314 %>%```<br>```  summarize(cases = sum(asthma_case), popn = sum(popn_1plus)) %>%```<br>```  mutate(crude_rate=(cases/popn)*100000)```<br><br>```print(asthma_overall)```<br><br>```` ``` ```` |

Calculate age-specific counts of prevalent asthma cases by first creating age categorisations and tabulating cases by those categories. What age groups will you use?

We will calculate rates and create a figure in the next step.

```r
```{r Calculate age-specific prevalence}
#age-specific prevalence
asthma_reg1314 <- asthma_reg1314 %>%
  mutate(agegrp = case_when(
    age >=1 & age <=9 ~ "<10",
    age >= 10 & age <= 19 ~ "10-19",
    age >= 20 & age <= 29 ~ "20-29",
    age >= 30 & age <= 39 ~ "30-39",
    age >= 40 & age <= 49 ~ "40-49",
    age >= 50 & age <= 59 ~ "50-59",
    age >= 60 & age <= 69 ~ "60-69",
    age >= 70 & age <= 79 ~ "70-79",
    age >= 80 & age <= 89 ~ "80-89",
    age >= 90 ~ "90+"))

#tabulate age groupings
asthma_reg1314 %>%
  count(agegrp)

#tabulate cases per age grouping
asthma_reg1314 %>%
  count(agegrp,asthma_case) %>%
  pivot_wider(names_from = asthma_case, values_from = n)
```
```

| | |
|---|---|
| Calculate and plot age-specific asthma rates using the same 10 year age groups you created in the previous step.. | ```r {r Age-specific asthma case rate}\n\n# Summarize cases and population per age range\nasthma_agespecific <- asthma_reg1314 %>%\n  group_by(agegrp) %>%\n  summarize(cases = sum(asthma_case), popn = sum(popn_1plus), .groups = "drop") %>%\n  filter(!is.na(agegrp))\n\n# Calculate rates\nasthma_agespecific <- asthma_agespecific %>%\n  mutate(agespecific_rate=(cases/popn)*100000)\n\n\n# Basic Plot\nagespec_barchart <- ggplot(data=asthma_agespecific, aes(x=agegrp, y=agespecific_rate)) +\n geom_bar(stat="identity")\n\n# Add simple options to the plot\nagespec_barchart <- agespec_barchart +\nlabs(title="Age-Specific Asthma Rate", x="Age Group (10 Year)", y = "Rate (per 100,000)")+\ntheme_classic()\n\n#Call plot\nagespec_barchart\n``` |

Age-Specific Asthma Rate

| | |
|---|---|
| Calculate and plot age- and sex-specific asthma rates. | ```` ```{r Age and sex-specific asthma rates} # Age and sex-specific prevalence  #add sex variable to the dataset asthma_reg1314 <- asthma_reg1314 %>%   mutate(sexchar = case_when(asthma_reg1314$sex == 1 ~ "Female",                             asthma_reg1314$sex == 2 ~ "Male",                             TRUE~"Unknown"))   #summarize based on age and sex asthma_as_specific <- asthma_reg1314 %>%   group_by(agegrp, sexchar) %>%   summarize(cases = sum(asthma_case), popn = sum(popn_1plus), .groups = "drop") %>%   filter(!is.na(agegrp))  #calculate rates asthma_as_specific <- asthma_as_specific %>%   mutate(as_specific_rate=(cases/popn)*100000)  # Plot age and sex-specific results agesexspec_barchart <- ggplot(data=asthma_as_specific, aes(fill=sexchar, x=agegrp, y=as_specific_rate)) + geom_bar(position="dodge", stat="identity") + labs(title="Age- and Sex-Specific Asthma Rate",     x="Age Group (10 Year)",     y = "Rate (per 100,000)",     fill="Sex") + theme_classic()  #Call plot agesexspec_barchart ``` ```` |

| | |
|---|---|
| | Age- and Sex-Specific Asthma Rate<br><br>Rate (per 100,000) axis: 900, 600, 300, 0<br><br>Sex<br>Female<br>Male<br><br>Age Group (10 year): <10, 10-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80-89, 90+ |
| Ensure pertinent details are captured in the script heading. | |
| Knit the document<br><br>Note: The document will be saved in your scripts folder or the folder you specified at the beginning of the markdown. | Knit |

## 5. Example output

This is an example of what your output will look like when you open it in your web-browser:

# R for Public Health Investigations: Day 3

Script Information:

*Purpose:* Exercise in advanced data management training for epidemiologists in R by using mock admin datasets to identify cases of Asthma

*Author:* J. Stares

*Date:* 26-Oct-2020

*Last modified by:* B. Hetman

*Date last modified:* 2020-11-15

Notes:

- Mock admin SAS datasets processed in R and Stata in 2017
- Training module location: Dropbox
- Case definition - Asthma: 1 year and older - 1H Ever or 2P in 2Y https://www.canada.ca/en/public-health/services/publications/canadian-chronic-disease-surveillance-system-factsheet.html

Setup:

```
require("knitr")
```

```
## Loading required package: knitr
```

```
## setting working directory
opts_knit$set(root.dir = "C:/IntroToR/Exercise_Day3/data/")
opts_chunk$set(fig.path = "C:/IntroToR/Exercise_Day3/data/")


#Load libraries
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 4.0.3
```

```
## -- Attaching packages ---------------------------------------------------------- tidyverse 1.3.0 --
```

```
## v ggplot2 3.3.2     v purrr   0.3.4
## v tibble  3.0.4     v dplyr   1.0.2
```

## 6. Critical thinking

1. Why did we build the final product as a registry table that flags numerators and denominators for asthma? Is there a better way?

2. What are possible limitations to taking the approach of using administrative health records to identify cases of disease? What are the strengths with this approach?

3. What are your thoughts on how this exercise was structured in comparison with exercises from Day 1 and Day 2 (i.e., several short scripts vs. one long one, chunks of code vs. code broken into lines, and working directory structure)?

## 7. Bonus task 1!

Bonus advanced content! Demonstration of user defined functions:

| Task | Code |
|------|------|
| Review any of the following material on Functions in R:<br><br>https://www.tutorialspoint.com/r/r_functions.htm<br><br>https://towardsdatascience.com/your-first-user-defined-function-in-r-1eedc634ead4<br><br>https://swcarpentry.github.io/r-novice-inflammation/02-func-R/<br><br>https://datasciencebeginners.com/2018/11/02/10-user-defined-functions-in-r/<br><br><br><br>Examine and run the demo code provided here to your rmd file.<br><br>What are the basic components of a user defined function? Which part of the code defines the name of the function? Are there any arguments input to this function? What does this function do? | <pre>##### So many plotting options!<br>```{r Demonstrating a custom function}<br># Create your own function to format plot styles!<br><br>theme_awesome <- function(){<br>  theme_classic() +<br>    theme(<br>      line = element_line(color="black"),<br>      panel.background = element_rect(fill = "white", colour = "darkgrey"),<br>      panel.grid.major.y =  element_line(colour = "grey", linetype = 3, size = 0.5),<br>      plot.background = element_rect(fill="white"),<br>      strip.background = element_rect(fill = "white"),<br>      legend.background = element_rect(fill = "white"),<br>      legend.key = element_rect(fill = "white"),<br>      legend.position = "bottom",<br>      plot.title = element_text(size = 14, family = "Helvetica Neue", face = "bold"),<br>      plot.subtitle = element_text(size = 9, family = "Helvetica Neue", face = "italic"),<br>      axis.title.x = element_text(),<br>      text = element_text(size = 11),<br><br>  )<br>}</pre> |

| | |
|---|---|
| Apply your custom function to the agesexspec_barchart object created earlier. | ```
agesexspec_barchart +
  theme_awesome() +
  scale_fill_manual(values = c("#4D648D", "#D0E1F9")) +
  scale_color_manual(values = "darkgrey") +
  labs(subtitle = "Great job for getting this far. Enjoy some cake!")
```<br><br>**Age- and Sex-Specific Asthma Rate**<br>*Great job for getting this far. Enjoy some cake!* |
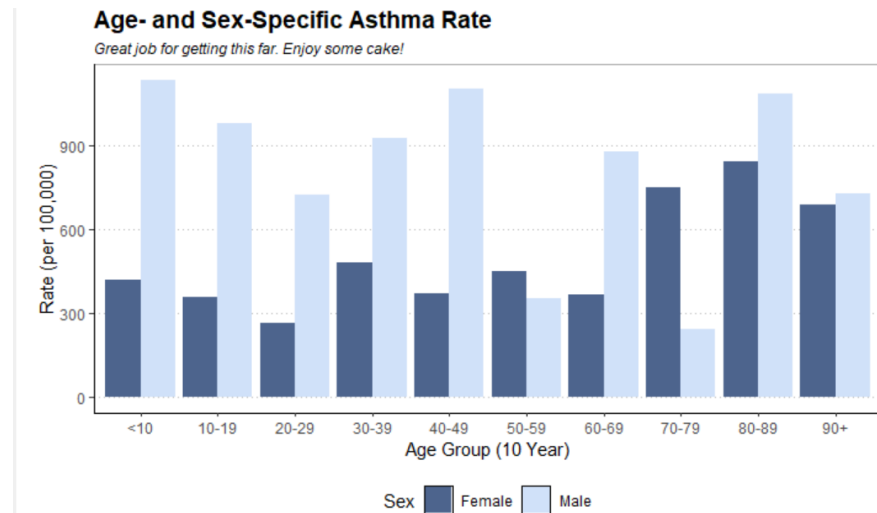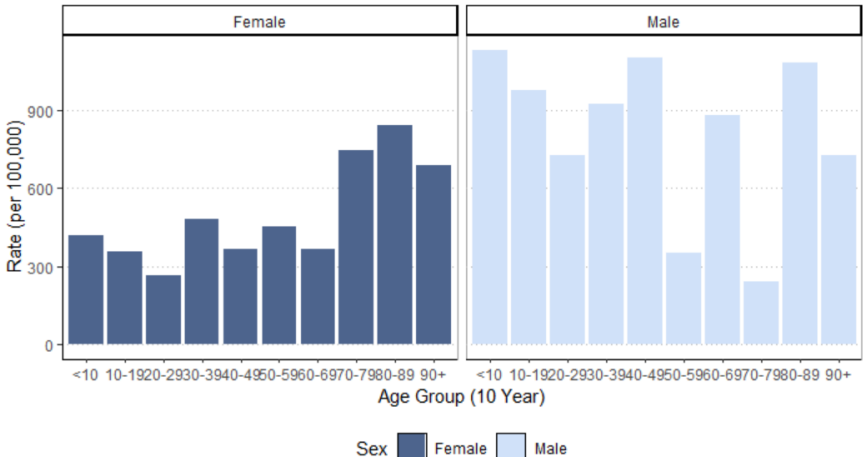| Knit the document | Knit |

## 7. Bonus task 2!

Bonus advanced content! Better living through modern technology - demonstration of facets and loops:

| Task | Code |
|------|------|
| Review the following tutorial on faceting for ggplot in R.<br><br>http://www.sthda.com/english/wiki/ggplot2-facet-split-a-plot-into-a-matrix-of-panels<br><br>Examine and add the demo code provided here to your rmd file.<br><br>What is a facet? What does this facet do specifically? How are facets applicable to your work? | ```\n##### Demonstrating facets and a simple loop:\n```{r Facetting on sex}\n# What if we wanted to compare sexes by themselves?\nagesexspec_barchart +\n  theme_awesome() +\n  scale_fill_manual(values = c("#4D648D", "#D0E1F9")) +\n  scale_color_manual(values = "darkgrey") +\n  labs(subtitle = "Great job for getting this far. Enjoy some cake!")+\n  facet_wrap(~sexchar)\n```<br><br><br>**Age- and Sex-Specific Asthma Rate**<br>Great job for getting this far. Enjoy some cake! |

| | |
|---|---|
| Review the following tutorial on For Loops in R.<br><br>https://www.r-bloggers.com/2015/12/how-to-write-the-first-for-loop-in-r/<br><br><br>Examine and run the demo code provided here.<br><br>What are the basic components of a for loop? What is the purpose of the "i in unique" statement? What does the code block do specifically? What additional functionality does the loop provide? How are loops applicable to your work? | <pre>```{r A simple loop}<br># What if we wanted a separate chart for each age group?<br># The following code will iterate through each age group created earlier,<br># and produce a unique chart for that age category only; then save that chart<br># in your output directory.<br><br># Check where your working directory is currently set to:<br>#getwd()<br><br>for (i in unique(asthma_as_specific$agegrp)){<br># Create titles for the individual plots based on the age group<br>title <- paste0("Asthma rate among age group ", i)<br># Special note: R will fail if the file name includes non-alpha-numeric characters<br># We use str_replace_all to remove these<br>title <- str_replace_all(title, "[^[:alnum:]]", " ")<br><br>#subset the data we need to plot<br>loop_plot <- asthma_as_specific %>% filter(agegrp==i)<br><br># Now we plot and save the individual files to our working directory<br>ggplot(data=loop_plot, aes(fill=sexchar, x=agegrp, y=as_specific_rate))+<br>  geom_bar(position="dodge", stat="identity")+<br>  labs(title="Age- and Sex-Specific Asthma Rate",<br>      x=paste0("Age (years)"),<br>      y="Rate (per 100,000)",<br>      fill="Sex") +<br>  theme_awesome() +<br>  scale_fill_manual(values = c("#4D648D", "#D0E1F9")) +<br>  scale_color_manual(values = "darkgrey") +<br>  labs(title = title,<br>      subtitle = "Great job for getting this far. Enjoy some cake!")<br><br>ggsave(path = output_folder, filename = paste0(title, ".png"), device = "png")<br>}<br>#Check your output directory for the exported plots!<br>```</pre> |
| Knit the document | ![Knit] |