

# export\_data

Group P

2022-09-27

This is a team assignment; each team should complete and turn in a PDF created from an Rmd via Github. Please include code and output for the following components:

Description of data set 1. What is the data source? (1-2 sentences on where the data is coming from, dates included, etc.)

*Data on monkey pox cases in the EU/EEA is accessible at: <https://www.ecdc.europa.eu/en/publications-data/data-monkeypox-cases-eueea>*

2. How does the data set relate to the group problem statement and question?

*The data set is going to help us understand how monkey pox case rates may differ by region and various demographic factors, additionally, the data set will allow us to determine if there is a relationship between certain demographics and monkey pox case rates.*

Import statement NOTE: Please use data sets available in the PHW251 Project Data github repo Links to an external site. (this is important to make sure everyone is using the same data sets)(done) Use appropriate import function and package based on the type of file(done) Utilize function arguments to control relevant components (i.e. change column types, column names, missing values, etc.)(working on it) Document the import process(working on it)

Loading the library to be used in this assignment

```
library(tidyverse)
```

```
## Warning in system("timedatectl", intern = TRUE): running command 'timedatectl'
## had status 1
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr  0.3.4
## v tibble  3.1.6      v dplyr  1.0.8
## v tidyr   1.2.0      v stringr 1.4.0
## v readr   2.1.2      v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(readr)
library(janitor)
```

```
##
## Attaching package: 'janitor'
```

```
## The following objects are masked from 'package:stats':
##
##   chisq.test, fisher.test
```

Importing data from git repositories

```
file_path1<-"https://raw.githubusercontent.com/PHW290/phw251_projectdata/main/euro_mpx_cases.csv"
monkey_pox <-read_csv(file_path1,na = c("", "NA", "*", "n/a"))%>% clean_names()
```

```
## Rows: 2987 Columns: 5
## -- Column specification -----
## Delimiter: ","
## chr  (3): CountryExp, CountryCode, Source
## dbl  (1): ConfCases
## date (1): DateRep
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```

file_path2 <- "https://raw.githubusercontent.com/PHW290/phw251_projectdata/main/euro_pop_denominators.csv"
pop_denominator<- read.csv(file_path2,na = c("", "NA", "*", "n/a"))>% clean_names()

file_path3 <- "https://raw.githubusercontent.com/PHW290/phw251_projectdata/main/euro_census_stats.csv"
census_stats <- read.csv(file_path3,na = c("", "NA", "*", "n/a"))>% clean_names()

file_path4<- "https://raw.githubusercontent.com/PHW290/phw251_projectdata/main/world_country_regions.csv"
world_country_region <- read.csv(file_path4,na = c("", "NA", "*", "n/a"))>% clean_names()

# checking if there is something unique about the variables in different dataset
x<-unique(monkey_pox$country_code)
y<-unique(pop_denominator$geo)

```

## Data manipulation and cleaning

```
#aggregating the monkey pox data set to the desired time period (week or month)& renaming the country c  
monkey_pox
```

```
## # A tibble: 2,987 x 5  
##   date_rep  country_exp country_code source conf_cases  
##   <date>    <chr>      <chr>      <chr>      <dbl>  
## 1 2022-05-09 Austria    AT        TESSy        0  
## 2 2022-05-09 Belgium   BE        TESSy        0  
## 3 2022-05-09 Bulgaria  BG        TESSy        0  
## 4 2022-05-09 Croatia   HR        TESSy        0  
## 5 2022-05-09 Cyprus    CY        TESSy        0  
## 6 2022-05-09 Czechia   CZ        TESSy        0  
## 7 2022-05-09 Denmark   DK        EI           0  
## 8 2022-05-09 Estonia   EE        EI           0  
## 9 2022-05-09 Finland   FI        EI           0  
## 10 2022-05-09 France    FR        EI           0  
## # ... with 2,977 more rows
```

```
library(lubridate)
```

```
##  
## Attaching package: 'lubridate'  
  
## The following objects are masked from 'package:base':  
##  
##   date, intersect, setdiff, union
```

```
monkey_pox <- monkey_pox%>% mutate(months_date= month(date_rep))%>%  
  rename(geo=country_code)  
monkey_pox
```

```
## # A tibble: 2,987 x 6  
##   date_rep  country_exp geo  source conf_cases months_date  
##   <date>    <chr>      <chr> <chr>      <dbl>      <dbl>  
## 1 2022-05-09 Austria    AT    TESSy        0          5  
## 2 2022-05-09 Belgium   BE    TESSy        0          5  
## 3 2022-05-09 Bulgaria  BG    TESSy        0          5  
## 4 2022-05-09 Croatia   HR    TESSy        0          5  
## 5 2022-05-09 Cyprus    CY    TESSy        0          5  
## 6 2022-05-09 Czechia   CZ    TESSy        0          5  
## 7 2022-05-09 Denmark   DK    EI           0          5  
## 8 2022-05-09 Estonia   EE    EI           0          5  
## 9 2022-05-09 Finland   FI    EI           0          5  
## 10 2022-05-09 France    FR    EI           0          5  
## # ... with 2,977 more rows
```

Joining monkey pox data set with pop denominator

```
merged_data <- merge(monkey_pox, pop_denominator, by="geo")  
# choosing the most recent time period
```

Identify data types for 5+ data elements/columns/variables Identify 5+ data elements required for your specified scenario. If <5 elements are required to complete the analysis, please choose additional variables of interest in the data set to explore in this milestone. Utilize functions or resources in RStudio to determine the types of each data element (i.e. character, numeric, factor) Identify the desired type/format for each variable—will you need to convert any columns to numeric or another type?

Provide a basic description of the 5+ data elements Numeric: mean, median, range Character: unique values/categories Or any other descriptives that will be useful to the analysis