

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC & KỸ THUẬT MÁY TÍNH



CẤU TRÚC RỜI RẠC CHO KHMT (CO1007)

Thông kê khảo sát kết quả Covid-19
môn Cấu trúc rời rạc

GVHD: Huỳnh Tường Nguyên
Nguyễn Ngọc Lễ
SV thực hiện: Đỗ Văn Băng – 2110813
Võ Văn Dũng – 2110102
Trần Hoàng Đại Sơn – 2110509
Nguyễn Tấn Dũng – 2110098
Hà Đình Nghĩa – 2114174
Trần Minh Khoa – 2110278

Tp. Hồ Chí Minh, Tháng 03/2022

Mục lục

1	Động cơ nghiên cứu	2
2	Mục tiêu	2
3	Kiến thức và kết quả chuẩn bị	2
3.1	Thống kê và các đặc trưng của mẫu số liệu	2
3.1.1	Thống kê là gì	2
3.1.2	Các đặc trưng của mẫu số liệu	2
3.2	Tìm hiểu ngôn ngữ lập trình R	5
3.2.1	Giới thiệu về R	5
3.2.2	Các điểm nổi bật của R	5
3.2.3	Các kiểu và cấu trúc dữ liệu cơ bản trong R	5
3.2.4	Các cấu trúc điều khiển trong R	6
3.2.5	Function trong R	7
3.2.6	Package trong R	7
4	Mô tả dữ liệu	7
5	Nhiệm vụ	7
	Nhóm câu hỏi liên quan đến tổng quát dữ liệu	8
	Nhóm câu hỏi liên quan đến mô tả thống kê cơ bản dữ liệu	15
	Nhóm câu hỏi liên quan đến dữ liệu thu thập	20
	Nhóm câu hỏi liên quan đến trực quan dữ liệu	24
	Nhóm câu hỏi liên quan đến trực quan dữ liệu theo thời gian là tháng	29
	Nhóm câu hỏi liên quan đến trực quan dữ liệu theo trung bình 7 ngày gần nhất	41
	Nhóm câu hỏi liên quan đến tất cả quốc gia theo thời gian là tháng	68
	Nhóm câu hỏi liên quan đến tất cả quốc gia theo trung bình 7 ngày gần nhất	76
	Nhóm câu hỏi liên quan đến sự tương quan giữa nhiễm bệnh và tử vong	91
	Nhóm câu hỏi riêng	122
	Tài liệu tham khảo	129

1 Động cơ nghiên cứu

Bệnh Corona do virus gây ra còn gọi là COVID-19 đã tạo ra những tác động tiêu cực đến nền đời sống của cư dân trên thế giới. Các đợt bùng phát của COVID-19 hay những biến thể virus đã mang đến những thách thức chưa từng có và được dự báo sẽ có tác động đáng kể đến sự phát triển kinh tế. Nhiều thông tin, tin tức về tình hình dịch bệnh cũng như dữ liệu về COVID-19 được phổ biến rộng rãi trong đời sống hay trên internet để giúp cho mọi người quan sát, phân tích, nghiên cứu được cập nhật hàng ngày.

Phân tích & thống kê dữ liệu về COVID-19 giúp cho ta thấy được số ca nhiễm bệnh, tử vong của một quốc gia, so sánh tình trạng của các quốc gia trong khu vực hay diễn biến dịch trên thế giới. Từ số liệu được báo cáo mọi chúng ta muốn biết các ca nhiễm bệnh có xu hướng tăng lên hay giảm xuống quy mô các đợt bùng phát ở mỗi quốc gia. Dữ liệu dùng cho bài tập lớn có tham khảo từ nguồn có thể xử lý trước với một vài thống kê cơ bản trước khi nó được truyền đi để khai thác dữ liệu thông minh sâu hơn.

2 Mục tiêu

Trong bài tập lớn này, các sinh viên sẽ bắt đầu với các bài toán thống kê đơn giản từ những dữ liệu được cung cấp. Qua đó, các em sẽ tìm ra những con số thú vị, có ý nghĩa đối với các dữ liệu thực tế từ tình hình dịch corona. Những kết quả mà các em tìm ra sẽ là bước khởi đầu cho việc khai phá nguồn dữ liệu của hệ thống sau này, nhằm đạt tới mục tiêu nâng cao kỹ năng lập trình, kỹ năng giải quyết vấn đề cho người học, kỹ năng làm việc nhóm cũng như hướng tới mục tiêu cao hơn là đam mê trong làm việc, học tập và nghiên cứu.

3 Kiến thức và kết quả chuẩn bị

3.1 Thống kê và các đặc trưng của mẫu số liệu

3.1.1 Thống kê là gì

Thống kê là khoa học về các phương pháp thu thập, tổ chức, trình bày, phân tích số liệu. Thống kê giúp ta phân tích số liệu một cách khách quan và rút ra các tri thức, thông tin chứa đựng trong các số liệu đó. Trên cơ sở này, chúng ta mới có thể đưa ra những dự báo và quyết định đúng đắn. Thống kê cần thiết cho mọi lực lượng lao động, đặc biệt rất cần cho các ngành quản lí, hoạch định chính sách.

3.1.2 Các đặc trưng của mẫu số liệu

- **Mẫu** là một tập con hữu hạn các đơn vị điều tra.
- **Mẫu số liệu** là các giá trị của dấu hiệu thu được trên một mẫu.
- **Dữ liệu ngoại lệ (outliers)** là một điểm dữ liệu có sự khác biệt đáng kể so với các quan sát khác. Dữ liệu ngoại lệ có thể xuất hiện do sự thay đổi thang đo hoặc do lỗi từ dữ liệu thu thập (thông thường dữ liệu ngoại lệ dạng này sẽ bị loại khỏi tập dữ liệu). Một giá trị ngoại lệ có thể gây ra vấn đề nghiêm trọng trong quá trình phân tích dữ liệu.
- **Kích thước mẫu N** là số phần tử của một mẫu.
- **Tần số n_i** của một giá trị x_i là số lần xuất hiện của giá trị đó trong mẫu.
- **Mốt** của dấu hiệu là số liệu có tần số lớn nhất trong mẫu số liệu, kí hiệu M_o .
- **Tần suất f_i** của giá trị x_i là tỉ số giữa tần số n_i và kích thước mẫu N , thường được biểu diễn dưới dạng phần trăm.

$$f_i = \frac{n_i}{N} \quad (1)$$

- **Số trung bình** của một mẫu số liệu có vai trò đại diện cho các số liệu của mẫu. Giá trị số trung bình của một mẫu kích thước N gồm $m(m \leq N)$ phần tử x_1, x_2, \dots, x_m với tần số tương ứng là n_1, n_2, \dots, n_m được tính bởi công thức:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^m n_i x_i \quad (2)$$

- **Số trung vị** M_e của một mẫu số liệu kích thước N được sắp xếp theo thứ tự không giảm là số liệu đứng ở vị trí thứ $\frac{N+1}{2}$ đối với trường hợp N là số lẻ. Nếu N là số chẵn thì số trung vị là trung bình cộng của hai số liệu ở vị trí thứ $\frac{N}{2}$ và $\frac{N}{2} + 1$.
- **Bách phân vị (Percentile)** là đại lượng dùng để ước tính tỷ lệ dữ liệu trong một tập số liệu rơi vào vùng cao hơn hoặc thấp hơn so với một giá trị cho trước. Bách phân vị chia dữ liệu có thứ tự theo hàng trăm. Để xác định giá trị v_p của phân vị thứ p trong một tập dữ liệu, ta thực hiện các bước sau:

1. Sắp xếp dữ liệu theo thứ tự từ bé đến lớn.
2. Tính chỉ số i :

$$i = \frac{p(n+1)}{100}$$

Trong đó:

- i là vị trí của giá trị dữ liệu tại phân vị p .
- p là phân vị thứ p .
- n là tổng số quan sát.

3. Xác định vị trí v_p :

- Nếu i là số nguyên thì phân vị thứ p là giá trị dữ liệu ở vị trí thứ i trong tập dữ liệu.
- Nếu i không phải là số nguyên thì làm tròn i lên và làm tròn i xuống số nguyên gần nhất, sau đó tính trung bình hai giá trị dữ liệu ở hai vị trí này trong tập dữ liệu.

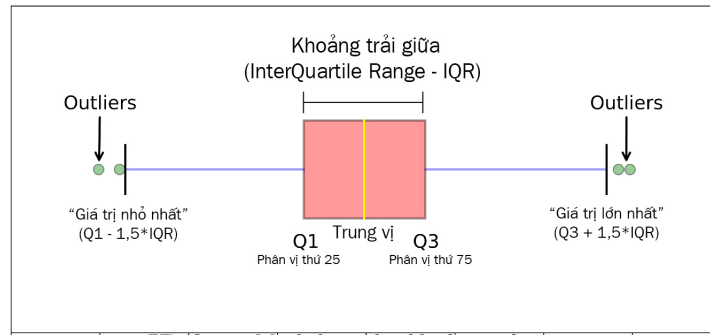
- **Tứ phân vị (Quartile)** là một trường hợp đặc biệt của bách phân vị. Tứ phân vị có 3 giá trị, đó là tứ phân vị thứ nhất, thứ nhì, và thứ ba. Ba giá trị này chia một tập hợp dữ liệu đã sắp xếp theo thứ tự thành 4 phần có số lượng quan sát đều nhau.

- Tứ phân vị thứ nhất Q_1 bằng trung vị phần dưới, tương đương bách phân vị thứ 25.
- Tứ phân vị thứ hai Q_2 chính bằng giá trị trung vị, tương đương bách phân vị thứ 50.
- Tứ phân vị thứ ba Q_3 bằng giá trị trung vị phần trên, tương đương bách phân vị thứ 75.

- **Khoảng trải giữa (Interquartile Range-IQR)** hay còn gọi là khoảng tứ phân vị của tập dữ liệu. Khoảng trải giữa là một con số cho biết mức độ lan truyền của nửa giữa hoặc 50% phần giữa của tập dữ liệu. IQR thường được sử dụng thay cho khoảng biến thiên (Range) vì nó loại trừ hầu hết giá trị bất thường hay giá trị ngoại lệ (Outliers) của dữ liệu. Công thức tính IQR:

$$IQR = Q_3 - Q_1$$

IQR có thể giúp xác định các giá trị ngoại lệ. Một giá trị bị nghi ngờ là một giá trị ngoại lệ nếu nó nhỏ hơn $1,5 * IQR$ dưới phần tư đầu tiên ($Q_1 - 1,5 * IQR$) hoặc lớn hơn $(1,5 * IQR)$ trên phần tư thứ ba ($Q_3 + 1,5 * IQR$) (Xem hình dưới). Các giá trị ngoại lệ luôn yêu cầu việc rà soát, kiểm tra lại dữ liệu. Những điểm dữ liệu đặc biệt này có thể do lỗi hoặc do sự bất thường trong dữ liệu nhưng cũng có thể là chìa khóa để hiểu dữ liệu.



Hình 1: Minh họa cho khoảng trải giữa

- **Phương sai (variance) và độ lệch chuẩn (standard deviation)** đo mức độ phân tán của các số liệu trong mẫu quanh số trung bình. Phương sai và độ lệch chuẩn càng lớn thì mức độ phân tán của số liệu càng lớn. Giả sử ta có mẫu số liệu kích thước N là x_1, x_2, \dots, x_n . Phương sai của mẫu số liệu này, kí hiệu s^2 được tính bởi công thức:

$$s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (3)$$

Căn bậc hai của phương sai được gọi là độ lệch chuẩn, kí hiệu s :

$$s = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \quad (4)$$

- **Độ méo lệch (skewness)** là một đại lượng đo lường mức độ bất đối xứng của phân phối xác suất của một biến ngẫu nhiên.

$$\gamma_1 = \frac{\sum_{i=1}^N (x_i - \bar{X})^3 / N}{s^3} \quad (5)$$

Trong đó:

- γ_1 là độ méo lệch.
- s là độ lệch chuẩn.
- N là số phần tử của mẫu.
- \bar{X} là giá trị trung bình của mẫu.
- x_i là giá trị thứ i của mẫu.
- **Độ nhọn (kurtosis)** là một đại lượng đo mức độ tập trung của phân phối xác suất của một biến ngẫu nhiên, cụ thể là mức độ tập trung của các quan sát quanh trung tâm của phân phối trong mối quan hệ với hai đuôi.

$$\gamma_2 = \frac{\sum_{i=1}^N (x_i - \bar{X})^4 / N}{s^4} \quad (6)$$

Trong đó:

- γ_2 là độ méo lệch.
- s là độ lệch chuẩn.
- N là số phần tử của mẫu.
- \bar{X} là giá trị trung bình của mẫu.
- x_i là giá trị thứ i của mẫu.

- **Hệ số tương quan Pearson (Pearson correlation coefficient)**: là số liệu thống kê kiểm tra đo lường mối quan hệ thống kê hoặc liên kết giữa các biến phụ thuộc với các biến liên tục. Hệ số tương quan sẽ trả lời cho các câu hỏi chẳng hạn như: Có mối quan hệ tương quan giữa nhiệt độ và doanh thu bán kem? Có mối quan hệ tương quan giữa sự hài lòng công việc, năng suất và thu nhập? Hay hai biến nào có mối liên hệ chặt chẽ nhất giữa tuổi, chiều cao, cân nặng, quy mô gia đình và thu nhập gia đình?

Hệ số tương quan Pearson (r) có giá trị giao động trong khoảng liên tục từ -1 đến $+1$:

- $r = 0$: Hai biến không có tương quan tuyến tính.
- $r = 1$ hoặc $r = -1$: Hai biến có mối tương quan tuyến tính tuyệt đối.
- $r < 0$: Hệ số tương quan âm. Nghĩa là giá trị biến x tăng thì giá trị biến y giảm và ngược lại, giá trị biến y tăng thì giá trị biến x giảm.
- $r > 0$: Hệ số tương quan dương. Nghĩa là giá trị biến x tăng thì giá trị biến y tăng và ngược lại, giá trị biến y tăng thì giá trị biến x cũng tăng.

3.2 Tìm hiểu ngôn ngữ lập trình R

3.2.1 Giới thiệu về R

R là một ngôn ngữ lập trình hàm cấp cao vừa là một môi trường dành cho tính toán thống kê. R hỗ trợ rất nhiều công cụ cho phân tích dữ liệu, khám phá tri thức và khai mở dữ liệu nhờ đó có thể phát triển nhanh các ứng dụng tính toán xác suất thống kê, phân tích dữ liệu.

Động lực ra đời của R là vào khoảng năm 1993, hai nhà thống kê học Ross Ihaka và Robert Gentleman ở University of Auckland (New Zealand) nhận thấy rằng các phần mềm thống kê thương mại sử dụng cho các tính toán thống kê vào thời điểm ấy còn đắt đỏ, không phù hợp và linh hoạt cho mục đích giảng dạy cũng như công việc. Sau đó hai ông đã quyết định lựa chọn ngôn ngữ S được phát triển bởi Bell Laboratories với nỗ lực viết một phần mềm thống kê mới. Cuối cùng Ross Ihaka và Robert Gentleman đã cho ra đời R và miễn phí cho tất cả người sử dụng vào năm 1995.



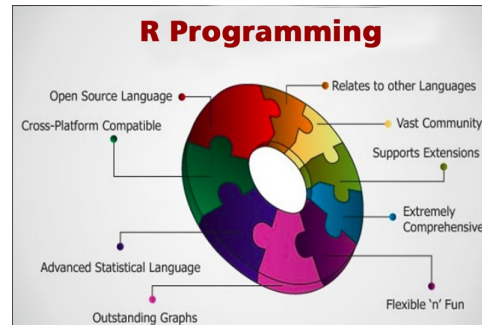
Hình 2: Logo của R

3.2.2 Các điểm nổi bật của R

- R có tính mã nguồn mở (open - source), nên liên tục được nâng cấp, cập nhật bởi cộng đồng người sử dụng trên toàn thế giới.
- R tương thích đa nền tảng, do đó có thể chạy trên bất kỳ hệ điều hành nào
- R tương thích với nhiều ngôn ngữ lập trình khác như C, C++, và FORTRAN. Những ngôn ngữ lập trình khác như .NET, Java, Python cũng có thể thao tác trực tiếp trên đối tượng
- R chứa nhiều package cho phép tương tác với nhiều cơ sở dữ liệu. Một số trong đó là Roracle, Rmysql, ...
- R có một cộng đồng lớn và tích cực. Mọi người từ khắp nơi trên thế giới có thể giúp đỡ và hỗ trợ. Nhiều ý tưởng và công nghệ mới nhất xuất hiện trong cộng đồng R..

3.2.3 Các kiểu và cấu trúc dữ liệu cơ bản trong R

- **Character/String**: Là kiểu dữ liệu dùng để lưu kí tự.
- **Numeric**: Là kiểu dữ liệu dùng để lưu trữ các số thực trong hệ thập phân.
- **Integer**: Là kiểu dữ liệu dùng để lưu trữ số nguyên.
- **Complex**: Là kiểu dữ liệu dùng để lưu trữ số phức.

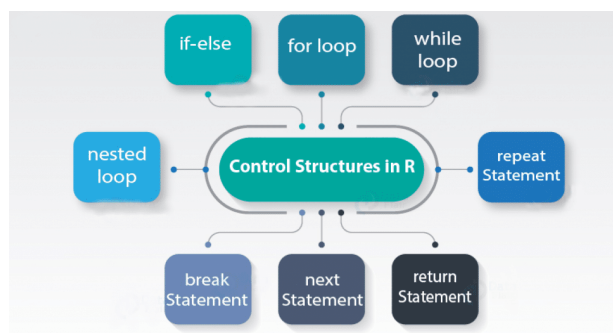


Hình 3: Ưu điểm nổi bật của R

- **Logical:** Là kiểu dữ liệu dùng để lưu trữ giá trị kiểu Boolean (0 và 1).
- **Vector:** Là cấu trúc dữ liệu đơn giản nhất trong R. Một vectơ là một tập hợp các phần tử có cùng kiểu dữ liệu. Một vectơ hỗ trợ các kiểu dữ liệu như logical, integer, double, character, complex, hoặc raw.
- **Matrix:** Là một tập hợp các vectơ có cùng độ dài, hay còn gọi là mảng hai chiều hoặc ma trận.
- **List:** Là một tập hợp các phần tử có kiểu dữ liệu mà kiểu dữ liệu của chúng có thể khác nhau. Một phần tử của List có thể là 1 kiểu dữ liệu cơ bản, hoặc có thể là 1 vectơ, 1 matrix, 1 list,...
- **Data Frame:** Là một kiểu dữ liệu cao cấp hơn matrix, được tích hợp từ nhiều cột hoặc nhiều hàng, tạo thành 1 bảng. Data Frames được tích hợp thêm rất nhiều hàm có ích cho việc thống kê, xử lý số liệu.
- **Array:** Là mảng nhiều chiều. Số chiều của array có thể lớn hơn 2. Người dùng có thể định nghĩa số chiều của Array.
- **Factor:** Là mảng 1 chiều giống vectơ nhưng lưu trữ có phân loại. Ta có thể thống kê các giá trị và tần số của các phần tử trong factor.

3.2.4 Các cấu trúc điều khiển trong R

- **Cấu trúc tuần tự:** Đây là cấu trúc có cơ bản trong mọi ngôn ngữ. Trừ khi có rẽ nhánh hoặc vòng lặp, thứ tự thực hiện câu lệnh được là từ trên xuống dưới, một cách tuần tự.
- **Cấu trúc rẽ nhánh có điều kiện:** Tương tự như các ngôn ngữ khác, câu lệnh *if* được dùng để điều khiển việc rẽ nhánh có điều kiện. Các câu lệnh *if* có thể lồng nhau tạo ra *Nested if* và *if* còn có thể đi kèm với *else* hoặc *else if*, giúp cho việc điều khiển rẽ nhánh có điều kiện trở nên hiệu quả hơn. Ngoài ra, R còn có câu lệnh *switch* có chức năng tương tự như *switch case* trong ngôn ngữ C.
- **Cấu trúc điều khiển vòng lặp:** Tương tự như các ngôn ngữ lập trình phổ biến, ngôn ngữ cũng có câu lệnh *for* và *while*. Các câu lệnh điều khiển việc kết thúc vòng lặp và rẽ nhánh không điều kiện trong R là *break* và *next* (giống continue của C). Ngoài ra R còn có *repeat* giúp tạo ra vòng lặp không cần điều kiện, chỉ kết thúc khi gặp câu lệnh *next*.



Hình 4: Cấu trúc điều khiển của R

3.2.5 Function trong R

- Built-in function: Đây là những hàm có thể gọi trực tiếp trong chương trình bởi người dùng như: `sq()`, `mean()`, `max()`, `min()`,...
- User-defined function: Đây là những hàm do người dùng tự thiết kế.

– **Cú pháp:**

```
function_name = function(arguments) {  
    statements  
    ...  
}
```

3.2.6 Package trong R

- **R-Package** là tập hợp các function và dữ liệu mẫu được lưu trữ trong một thư mục gọi là "library" trong môi trường làm việc của **R**.
- Các package này thường được cài đặt tự động trong quá trình cài đặt **R**. Chúng ta cũng có thể cài đặt thêm một số package khác cho mục đích cụ thể.
- Các package của R được chứa trong các repository(kho chứa), trong đó repository phổ biến và đầy đủ nhất là **CRAN** (Comprehensive R Archive Network).
- Đưa package vào library: để sử dụng package nào đó thì chúng ta phải nó lên môi trường làm việc hiện tại của R bằng câu lệnh `library("package Name", lib.loc = "path to library")`.

4 Mô tả dữ liệu

Dữ liệu gồm các thuộc tính chính “**iso_code**, **continent**, **location**, **date**, **new_cases**, **new_deaths**” được lưu trong file **csv**.

1. **iso_code**: Định danh đất nước
2. **continent**: Tên châu lục
3. **location**: Tên quốc gia
4. **date**: Ngày quan sát với định dạng Month-Day-Year
5. **new_cases**: Số trường hợp COVID-19 mới được xác nhận
6. **new_deaths**: Số tử vong mới do COVID-19

5 Nhiệm vụ

Gọi **MD** là mã đề riêng cho mỗi nhóm (gồm 4 ký số) không trùng nhau, nhóm sinh viên sẽ thực hiện các yêu cầu dưới đây với các giá trị xác định như sau:

- Mỗi nhóm sẽ dùng R để thao tác trên số file dữ liệu khác nhau được chọn theo cột “STT” theo cách tính $kq = (kytu1 + kytu2 + kytu3 + kytu4) \% 6$:
 - Nếu $kq = 0$ thì làm các stt là 1,2,3
 - Nếu $kq = 1$ thì làm các stt là 4,5,6
 - Nếu $kq = 2$ thì làm các stt là 7,8,9
 - Nếu $kq = 3$ thì làm các stt là 10,11,12
 - Nếu $kq = 4$ thì làm các stt là 13,14,15
 - Nếu $kq = 5$ thì làm các stt là 16,17,18

STT	đất nước	STT	đất nước
1	Kenya	10	Canada
2	Lesotho	11	Greenland
3	Morocco	12	United States
4	Indonesia	13	Australia
5	Japan	14	New Caledonia
6	Vietnam	15	New Zealand
7	Andorra	16	Brazil
8	Slovenia	17	Chile
9	United Kingdom	18	Venezuela

Hoàn thành các bài tập:

- Đối với các bài tập chung gồm các phần $\{i, \dots, xiii\}$, tất cả các nhóm đều phải làm
- Mỗi nhóm sẽ thực hiện 4 câu riêng của mình trong phần ix bằng các lấy 4 ký số trong mã đề của mình là chỉ số câu hỏi tương ứng. Nếu ký số là 0 thì làm câu 10.

Do MD của nhóm là 6320, nên nhóm sẽ giải quyết vấn đề của:

- Ba quốc gia là Brazil, Chile và Venezuela.
- Đối với phần x nhóm sẽ giải quyết 4 câu là 2, 3, 6 và 10.

i) Nhóm câu hỏi liên quan đến tổng quát dữ liệu

Dùng tập dữ liệu để trả lời các câu hỏi và trình bày theo định dạng

1) Tập mẫu thu thập dữ liệu vào các năm nào

- Hiện thực trong R

```
1 library("tidyverse")
2 Data <- read.csv("C:/Users/Asus/Documents/covidData.csv")
3 years = substring(Data$date, nchar(Data$date)-3, nchar(Data$date))
4 fac = factor(years)
5 print("Dữ liệu được thu thập vào những năm: ")
6 print(as.integer(levels(fac)))
```

- Kết quả

```
[1] "Dữ liệu được thu thập vào những năm: "
> print(as.integer(levels(fac)))
[1] 2020 2021 2022
```

2) Số lượng đất nước và định danh của mỗi đất nước (hiển thị 10 đất nước đầu tiên).

- Hiện thực trong R

```
1 library("tidyverse")
2 library("gt")
3 data <- read.csv("owid-covid-data.csv", TRUE, ",")
4 x <- distinct(data.frame(factor(data$iso_code)))
5 y <- distinct(data.frame(factor(data$location)))
6 g <- head(cbind(x,y), n = 10)
7 colnames(g) <- c("iso_code:", "Country")
8 g$iso_code: <- as.character(g$iso_code:)
9 g$Country <- as.character(g$Country)
10 g[nrow(g)+1,] <- c("Count", nrow(x))
11 gt(g) %>%
12 cols_align(
13   align = c("center"),
14   columns = "iso_code:"
15 )
```

- Kết quả

iso_code:	Country
AFG	Afghanistan
OWID_AFR	Africa
ALB	Albania
DZA	Algeria
AND	Andorra
AGO	Angola
AIA	Anguilla
ATG	Antigua and Barbuda
ARG	Argentina
ARM	Armenia
Count	238

3) Số lượng châu lục trong tập mẫu

- Hiện thực trong R

```
1 library("tidyverse")
2 library("gt")
3 data <- read.csv("owid-covid-data.csv",TRUE, ",")
4 x <- distinct(data.frame(factor(data$continent)))
5 x <- data.frame(x[-c(2),])
6 colnames(x) <- c("Continent:")
7 y <- data.frame(c("Chau A", "Chau Au", "Chau Phi", "Chau Bac My", "Chau Nam My", "Chau Dai Duong"
8   ))
9 colnames(y) <- c("So chau luc:")
10 gt(cbind(x,y)) %>%
11   cols_align(
12     align = c("center"),
13     columns = "Continent:"
14   )
```

- Kết quả

Continent:	So chau luc:
Asia	Chau A
Europe	Chau Au
Africa	Chau Phi
North America	Chau Bac My
South America	Chau Nam My
Oceania	Chau Dai Duong

4) Số lượng dữ liệu thu thập được trong từng từng châu lục và tổng số

- Hiện thực trong R

```
1 library("gt")
2 data <- read.csv("owid-covid-data.csv",TRUE, ",")
3 w <- data.frame(table(data$continent))
4 colnames(w) <- c("Continent:", "Observations")
5 w <- w[-c(1),]
6 w$"Continent:" <- as.character(w$"Continent:")
7 w[nrow(w)+1,] <- c("Tong:", sum(w$Observations))
8 w$"Continent:" <- as.factor(w$"Continent:")
```

```

9 gt(w) %>%
10   cols_align(
11     align = c("center"),
12     columns = "Continent:"
13   )

```

- Kết quả

Continent:	Observations
Africa	38647
Asia	35528
Europe	36375
North America	24438
Oceania	8993
South America	9335
Tong:	153316

5) Số lượng dữ liệu thu thập được trong từng từng đất nước (hiển thị 10 đất nước cuối cùng) và tổng số

- Hiện thực trong R

```

1 library("gt")
2 data <- read.csv("owid-covid-data.csv", TRUE, ",")
3 w <- data.frame(table(data$iso_code))
4 y = tail(w, n = 10)
5 colnames(y) <- c("iso_code", "Observations")
6 y$iso_code <- as.character(y$iso_code)
7 y[nrow(y)+1,] <- c("Tong:", sum(y$Observations))
8 y$iso_code <- as.factor(y$iso_code)
9 gt(y) %>%
10   cols_align(
11     align = c("center"),
12     columns = "iso_code"
13   )

```

- Kết quả

iso_code	Observations
VEN	708
VGB	694
VNM	759
VUT	467
WLF	489
WSM	459
YEM	681
ZAF	744
ZMB	704
ZWE	702
Tong:	6407

6) Cho biết các châu lục nào có lượng dữ liệu thu thập nhỏ nhất và giá trị nhỏ nhất đó?

- Hiện thực trong R

```
1 data <- read.csv("owid-covid-data.csv",TRUE,"")
2 w <- data.frame(table(data$continent))
3 colnames(w) <- c("Continent", "MinimumValue")
4 print(subset(w,w$MinimumValue == min(w$MinimumValue)))
5 min(w$Observations)
```

- Kết quả

	Continent	MinimumValue
6	Oceania	8993

7) Cho biết các châu lục nào có lượng dữ liệu thu thập lớn nhất và giá trị lớn nhất đó?

- Hiện thực trong R

```
1 data <- read.csv("owid-covid-data.csv",TRUE,"")
2 w <- data.frame(table(data$continent))
3 colnames(w) <- c("Continent", "MaximumValue")
4 print(subset(w,w$MaximumValue == max(w$MaximumValue)))
```

- Kết quả

	Continent	MaximumValue
2	Africa	38647

8) Cho biết các nước nào có lượng dữ liệu thu thập nhỏ nhất và giá trị nhỏ nhất đó?

- Hiện thực trong R

```
1 data <- read.csv("owid-covid-data.csv",TRUE,"")
2 w <- data.frame(table(data$location))
3 colnames(w) <- c("Country", "MinimumValue")
4 print(subset(w,w$MinimumValue == min(w$MinimumValue)))
```

- Kết quả

	Country	MinimumValue
171	Pitcairn	85

9) Cho biết các nước nào có lượng dữ liệu thu thập lớn nhất và giá trị lớn nhất đó?

- Hiện thực trong R

```
1 data <- read.csv("owid-covid-data.csv",TRUE,"")
2 w <- data.frame(table(data$location))
3 colnames(w) <- c("Country", "MaximumValue")
4 print(subset(w,w$MaximumValue == max(w$MaximumValue)))
```

- Kết quả

	Country	MaximumValue
9	Argentina	781
137	Mexico	781

10) Cho biết các date nào có lượng dữ liệu thu thập nhỏ nhất và giá trị nhỏ nhất đó?

- Hiện thực trong R

```
1 data <- read.csv("owid-covid-data.csv",TRUE,"")
2 w <- data.frame(table(data$date))
3 colnames(w) <- c("Date", "MinimumValue")
4 print(subset(w,w$MinimumValue == min(w$MinimumValue)))
```

- Kết quả

	Date	MinimumValue
1	1/1/2020	2

34	1/2/2020	2
67	1/3/2020	2

11) Cho biết các date nào có lượng dữ liệu thu thập lớn nhất và giá trị lớn nhất đó?

- Hiện thực trong R

```
1 data <- read.csv("owid-covid-data.csv", TRUE, ",")
2 w <- data.frame(table(data$date))
3 colnames(w) <- c("Date", "MaximumValue")
4 print(subset(w, w$MaximumValue == max(w$MaximumValue)))
```

- Kết quả

	Date	MaximumValue
689	8/22/2021	238
691	8/23/2021	238
693	8/24/2021	238
695	8/25/2021	238
697	8/26/2021	238
699	8/27/2021	238
701	8/28/2021	238
703	8/29/2021	238

12) Cho biết số lượng dữ liệu thu thập được theo date và châu lục.

- Hiện thực trong R

```
1 library("tidyverse")
2 Data <- read.csv("C:/Users/Asus/Documents/covidData.csv")
3 continents <- Data$continent[Data$continent != ""]
4 continents <- c(levels(factor(continents)))
5 dates <- Data$date[Data$date != ""]
6 dates <- c(levels(factor(dates)))
7 size_c = NROW(continents)
8 size_d = NROW(dates)
9 size_res = size_c*size_d
10 res <- data.frame(matrix(NA, nrow = size_res, ncol = 3))
11 colnames(res) <- c("Continent", "Date", "Observations")
12 j = 1
13 k = 1
14 for (i in 1:size_res) {
15   res$Continent[i] <- continents[j]
16   res$Date[i] <- dates[k]
17   k <- k + 1
18   if (i == size_d) j <- j + 1
19   if (k == size_d) k <- 1
20 }
21 for (i in 1:size_res) {
22   temp <- subset(Data, Data$continent == res$Continent[i] & Data$date == res$Date[i])
23   res$Observations[i] <- nrow(temp)
24 }
25 print(res)
```

- Kết quả sau khi thực hiện chương trình trên ("...]" biểu thị còn nhiều giá trị do không đủ không gian nên nhóm không nhập vào):

	Continent	Date	Observations
1	Africa	1/1/2020	0
2	Africa	1/1/2021	55
3	Africa	1/1/2022	55
4	Africa	1/10/2020	0
[...]			
329	Africa	2/27/2021	55
330	Africa	2/28/2020	6
331	Africa	2/28/2021	55

```
332 Africa 2/29/2020 6
333 Africa 2/3/2020 0
[ reached 'max' / getOption("max.print") -- omitted 4353 rows
```

13) Cho biết số lượng dữ liệu thu thập được là lớn nhất theo date và châu lục.

- Dựa vào kết quả câu 12, ta tái sử dụng code của câu 12 và dùng câu lệnh dưới đây để thực hiện yêu cầu đề bài:

```
1 result <- subset(res, res$Observations == max(res$Observations))
2 row.names(result) <- c(seq(1:NROW(result)))
3 print(result)
```

- Kết quả sau khi thực hiện như sau:

```
Continent Date Observations
1 Africa 1/1/2021 55
2 Africa 1/1/2022 55
3 Africa 1/10/2021 55
4 Africa 1/10/2022 55
[...]
329 Africa 4/13/2021 55
330 Africa 4/14/2021 55
331 Africa 4/15/2021 55
332 Africa 4/16/2021 55
333 Africa 4/17/2021 55
[ reached 'max' / getOption("max.print") -- omitted 197 rows ]
```

14) Cho biết số lượng dữ liệu thu thập được là nhỏ nhất theo date và châu lục.

- Tương tự câu trên, ta giải quyết yêu cầu bài toán trên như sau:

```
1 result <- subset(res, res$Observations == min(res$Observations))
2 row.names(result) <- c(seq(1:NROW(result)))
3 print(result)
```

- Kết quả

```
Continent Date Observations
1 Africa 1/1/2020 0
2 Africa 1/10/2020 0
3 Africa 1/11/2020 0
4 Africa 1/12/2020 0
[...]
48 Asia 1/2/2020 0
49 Asia 1/3/2020 0
50 Asia 1/1/2020 0
51 Asia 1/2/2020 0
52 Asia 1/3/2020 0
53 Asia 1/1/2020 0
```

15) Với một date là k và châu lục t cho trước, hãy cho biết số lượng dữ liệu thu thập được.

- Tương tự ta sử dụng đoạn code ở câu 12, để hiện thực bài toán ta cần thêm các dòng lệnh sau:

```
1 k = readline("Nhập date bạn muốn tra cứu: ");
2 t = readline("Nhập continent bạn muốn tra cứu: ");
3 result <- res$Observations[res$Continent == t & res$Date == k]
4 print("Kết quả là: ")
5 print(result)
```

- Kết quả trả về màn hình với trường hợp ví dụ k = 11/1/2020, t = Africa:

```
> k = readline("Nhap date ban muon tra cuu: ");
Nhap date ban muon tra cuu: 11/1/2020
> t = readline("Nhap continent ban muon tra cuu: ");
Nhap continent ban muon tra cuu: Africa
> result <- res$Observations[res$Continent == t & res$Date == k]
> print("Ket qua la: ")
[1] "Ket qua la: "
> print(result)
[1] 55
```

16) Có đất nước nào mà số lượng dữ liệu thu thập được là bằng nhau không? Hãy cho biết các iso_code của đất nước đó.

- Hiện thực trong R

```
1 library("tidyverse")
2 Data <- read.csv("C:/Users/Asus/Documents/covidData.csv")
3 countrys <- Data$location[is.na(Data$iso_code) == 0]
4 countrys <- c(levels(factor(countrys)))
5 num_of_data <- c()
6 for (i in 1:NROW(countrys)) {
7   num_of_data[i] <- NROW(subset(Data, Data$location == countrys[i]))
8 }
9 isoCode <- c()
10 for (i in 1:NROW(countrys)) {
11   isoCode[i] <- Data$iso_code[Data$location == countrys[i]][i]
12 }
13 res <- cbind(isoCode, countrys, num_of_data)
14 colnames(res) <- c("iso_code", "Country", "Num_of_data")
15 size_res <- nrow(res)
16 res <- data.frame(res)
17 res <- subset(res, is.na(res$iso_code) == 0)
18 row.names(res) <- c(seq(1:NROW(res)))
19 eq_data <- c(levels(factor(res$Num_of_data)))
20 for (i in 1:NROW(eq_data)) {
21   temp <- c()
22   k <- 1
23   t <- FALSE
24   for (j in 1:NROW(res)) {
25     if (res$Num_of_data[j] == eq_data[i]) {
26       temp[k] <- res$iso_code[j]
27       k <- k + 1
28     }
29     if (k > 2) t <- TRUE
30   }
31   if (t == TRUE) {
32     temp <- data.frame(temp)
33     temp <- t(temp)
34     colnames(temp) <- c(seq(1:NCOL(temp)))
35     rownames(temp) <- c(paste(eq_data[i], ":"))
36     print(temp)
37   }
38 }
```

- Kết quả được thể hiện theo định dạng: với mỗi hàng (trừ hàng đánh thứ tự), con số đầu tiên biểu thị số lượng dữ liệu thu thập được ở các quốc gia, các chuỗi kí tự tiếp theo là các iso_code của các quốc gia có số lượng dữ liệu thu thập được bằng với con số đầu tiên mỗi hàng. Khi đó ta thấy có nhiều trường hợp số lượng dữ liệu thu thập được ở các quốc gia bằng nhau:

```
1      2
686 : "SPM" "SSD"
1      2
691 : "BDI" "SLE"
1      2      3
694 : "AIA" "VGB" "TCA"
1      2      3
```

```
697 : "GNB" "MLI" "KNA"
      1     2     3     4     5
700 : "DMA" "GRD" "MOZ" "SYR" "TLS"
      1     2
[... ]
      1     2
781 : "ARG" "MEX"
```

17) Liệt kê iso_code, tên đất nước mà chiều dài iso_code lớn hơn 3.

- Hiện thực trong R

```
1 library("tidyverse")
2 Data <- read.csv("C:/Users/Asus/Documents/covidData.csv")
3 iso_code <- Data$iso_code[is.na(Data$iso_code) == 0]
4 iso_code <- c(levels(factor(iso_code)))
5 countrys <- c()
6 for (i in 1:(NROW(iso_code))) {
7   countrys[i] <- Data$location[Data$iso_code == iso_code[i]]
8 }
9 res <- cbind(iso_code, countrys)
10 res <- data.frame(res)
11 colnames(res) <- c("iso_code", "Country")
12 result <- subset(res, nchar(res$iso_code) > 3)
13 row.names(result) <- c(seq(1:NROW(result)))
14 print(result)
```

- Kết quả

	iso_code	Country
1	OWID_AFR	Africa
2	OWID_ASI	Asia
3	OWID_CYN	Northern Cyprus
4	OWID_EUN	European Union
5	OWID_EUR	Europe
6	OWID_HIC	High income
7	OWID_INT	International
8	OWID_KOS	Kosovo
9	OWID_LIC	Low income
10	OWID_LMC	Lower middle income
11	OWID_NAM	North America
12	OWID_OCE	Oceania
13	OWID_SAM	South America
14	OWID_UMC	Upper middle income
15	OWID_WRL	World

ii) Nhóm câu hỏi liên quan đến mô tả thống kê cơ bản dữ liệu

Với mỗi quốc gia cần tính số liệu thống kê lần lượt cho nhiễm và tử vong do coronavirus được báo cáo mới:

- Chuẩn bị dữ liệu cho toàn bộ phần ii (phần code dùng chung cho cả phần ii)

```
1 library(dataset)
2 library(tidyverse)
3 library(dplyr)
4 covid = read.csv("C:/Users/Admin/Downloads/covidData.csv")
5 covid = covid %>% filter(!continent == '') %>% mutate(new_cases = abs(new_cases), new_deaths = abs(
  new_deaths))
6 brazil = subset(covid, iso_code == "BRA")
7 chile = subset(covid, iso_code == "CHL")
8 venezuela = subset(covid, iso_code == "VEN")
9 brazil[is.na(brazil)] = 0
10 chile[is.na(chile)] = 0
11 venezuela[is.na(venezuela)] = 0
```



```
12 Countries = c("Brazil", "Chile", "Venezuela")
```

1) Tính giá trị nhỏ nhất, lớn nhất

- Hướng giải quyết: Tìm giá trị nhỏ nhất/ lớn nhất bằng min/max.
- Hiện thực trong R

```
1 #for covid new_cases
2 infections = data.frame(Countries)
3 infections[["Min"]] <- NA
4 infections[1, "Min"] = min(brazil$new_cases)
5 infections[2, "Min"] = min(chile$new_cases)
6 infections[3, "Min"] = min(venezuela$new_cases)
7 infections[1, "Max"] = max(brazil$new_cases)
8 infections[2, "Max"] = max(chile$new_cases)
9 infections[3, "Max"] = max(venezuela$new_cases)
10
11
12 # for covid new_deaths
13 deaths = data.frame(Countries)
14 deaths[["Min"]] <- NA
15 deaths[1, "Min"] = min(brazil$ new_deaths)
16 deaths[2, "Min"] = min(chile$ new_deaths)
17 deaths[3, "Min"] = min(venezuela$ new_deaths)
18 deaths[1, "Max"] = max(brazil$ new_deaths)
19 deaths[2, "Max"] = max(chile$ new_deaths)
20 deaths[3, "Max"] = max(venezuela$ new_deaths)
```

- Kết quả đối với các ca nhiễm

Countries	Min	Max
Brazil	0	287149
Chile	0	41651
Venezuela	0	4418

- Kết quả đối với các ca tử vong

Countries	Min	Max
Brazil	0	4148
Chile	0	1057
Venezuela	0	35

2) Tính tứ phân vị thứ nhất(Q1), thứ hai(Q2), thứ ba(Q3)

- Hướng giải quyết: Dùng quantile để tìm tứ phân vị.
- Hiện thực trong R

```
1 infections = data.frame(Countries)
2 infections[["Q1"]] = NA
3 infections[["Q2"]] = NA
4 infections[["Q3"]] = NA
5 for (i in c(2,3,4)) {
6   infections[[i]][1] = quantile(brazil$new_cases)[[i]]
7 }
8 for (i in c(2,3,4)) {
9   infections[[i]][2] = quantile(chile$new_cases)[[i]]
10 }
11 for (i in c(2,3,4)) {
12   infections[[i]][3] = quantile(venezuela$new_cases)[[i]]
13 }
14
15 deaths = data.frame(Countries)
16 deaths[["Q1"]] = NA
17 deaths[["Q2"]] = NA
18 deaths[["Q3"]] = NA
19 for (i in c(2,3,4)) {
20   deaths[[i]][1] = quantile(brazil$new_cases)[[i]]
21 }
```

```

22 for (i in c(2,3,4)) {
23   deaths[[i]][2] = quantile(chile$new_cases)[[i]]
24 }
25 for (i in c(2,3,4)) {
26   deaths[[i]][3] = quantile(venezuela$new_cases)[[i]]
27 }

```

- Kết quả đối với các ca nhiễm

	Countries	Q1	Q2	Q3
1	Brazil	13276.00	31149.0	53989.00
2	Chile	1179.75	1997.5	4315.25
3	Venezuela	245.00	590.0	1100.25

- Kết quả đối với các ca tử vong

	Countries	Q1	Q2	Q3
1	Brazil	13276.00	31149.0	53989.00
2	Chile	1179.75	1997.5	4315.25
3	Venezuela	245.00	590.0	1100.25

3) Tính giá trị trung bình (Avg)

- Hướng giải quyết: Dùng mean để tìm giá trị trung bình.
- Hiện thực trong R

```

1 infections = data.frame(Countries)
2 infections[["Avg"]] =NA
3 infections[1, "Avg"] = mean(brazil$ new_cases)
4 infections[2, "Avg"] = mean(chile$ new_cases)
5 infections[3, "Avg"] = mean(venezuela$ new_cases)
6
7 deaths = data.frame(Countries)
8 deaths[["Avg"]] =NA
9 deaths[1, "Avg"] = mean(brazil$ new_deaths)
10 deaths[2, "Avg"] = mean(chile$ new_deaths)
11 deaths[3, "Avg"] = mean(venezuela$ new_deaths)

```

- Kết quả đối với các ca nhiễm

	Countries	Avg
1	Brazil	38747.6979
2	Chile	3872.5907
3	Venezuela	721.1201

- Kết quả đối với các ca tử vong

	Countries	Avg
1	Brazil	888.544828
2	Chile	57.173077
3	Venezuela	7.929379

4) Tính giá trị độ lệch chuẩn (Std)

- Hướng giải quyết: Dùng sd để tìm giá trị độ lệch chuẩn.
- Hiện thực trong R

```

1 infections = data.frame(Countries)
2 infections[["Std"]] =NA
3 infections[[2]][1] = sd(brazil$new_cases)
4 infections[[2]][2] = sd(chile$new_cases)
5 infections[[2]][3] = sd(venezuela$new_cases)
6
7 deaths = data.frame(Countries)
8 deaths[["Std"]] =NA
9 deaths[[2]][1] = sd(brazil$new_deaths)
10 deaths[[2]][2] = sd(chile$new_deaths)
11 deaths[[2]][3] = sd(venezuela$new_deaths)

```

- Kết quả đối với các ca nhiễm

	Countries	Std
1	Brazil	37015.4040
2	Chile	6093.3562
3	Venezuela	631.3701

- Kết quả đối với các ca tử vong

	Countries	Std
1	Brazil	771.595312
2	Chile	67.164298
3	Venezuela	7.214482

- 5) Đếm xem có bao nhiêu outliers, một quan sát mà giá trị của nó nằm trong khoảng sau:

$$IQR = Q3 - Q1$$

$$outliers < Q1 - 1.5 * IQR \text{ hoặc } outliers > Q3 + 1.5 * IQR$$

- Hướng giải quyết: Dùng `boxplot.stats(...)` để tìm các outliers rồi dùng `length(...)` để đếm xem có bao nhiêu outliers.
- Hiện thực trong R

```
1 infections = data.frame(Countries)
2 infections[["Outlier"]] = NA
3 infections[[2]][[1]] = length(boxplot.stats(brazil$new_cases)$out)
4 infections[[2]][[2]] = length(boxplot.stats(chile$new_cases)$out)
5 infections[[2]][[3]] = length(boxplot.stats(venezuela$new_cases)$out)
6
7 deaths = data.frame(Countries)
8 deaths[["Outlier"]] = NA
9 deaths[[2]][[1]] = length(boxplot.stats(brazil$new_deaths)$out)
10 deaths[[2]][[2]] = length(boxplot.stats(chile$new_deaths)$out)
11 deaths[[2]][[3]] = length(boxplot.stats(venezuela$new_deaths)$out)
```

- Kết quả đối với các ca nhiễm

	Countries	Outlier
1	Brazil	26
2	Chile	38
3	Venezuela	15

- Kết quả đối với các ca tử vong

	Countries	Outlier
1	Brazil	40
2	Chile	28
3	Venezuela	21

- 6) Lập bảng mô tả số liệu thống kê các đất nước cho từng thể loại:

`infections[deaths]:`

Countries	Min	Q1	Q2	Q3	Max	Avg	Std	Outlier
ctr_i	?	?	?	?	?	?	?	?

- Hướng giải quyết: Kết hợp thông tin từ các câu trên.
- Hiện thực trong R

```
1 infections = data.frame(Countries)
2 infections[["Min"]] <- NA
3 infections[1, "Min"] = min(brazil$new_cases)
4 infections[2, "Min"] = min(chile$new_cases)
5 infections[3, "Min"] = min(venezuela$new_cases)
6 infections[["Q1"]] = NA
7 infections[["Q2"]] = NA
8 infections[["Q3"]] = NA
9 infections[["Max"]] = NA
10 infections[["Avg"]] = NA
11 infections[["Std"]] = NA
12 infections[["Outlier"]] = NA
13
14 infections[1, "Max"] = max(brazil$new_cases)
```

```
15 infections[2, "Max"] = max(chile$new_cases)
16 infections[3, "Max"] = max(venezuela$new_cases)
17 infections[1, "Avg"] = mean(brazil$new_cases)
18 infections[2, "Avg"] = mean(chile$new_cases)
19 infections[3, "Avg"] = mean(venezuela$new_cases)
20 for (i in c(3,4,5)) {
21   infections[[i]][[1]] = quantile(brazil$new_cases)[[i-1]]
22 }
23 for (i in c(3,4,5)) {
24   infections[[i]][[2]] = quantile(chile$new_cases)[[i-1]]
25 }
26 for (i in c(3,4,5)) {
27   infections[[i]][[3]] = quantile(venezuela$new_cases)[[i-1]]
28 }
29 infections[[9]][[1]] = length(boxplot.stats(brazil$new_cases)$out)
30 infections[[8]][[1]] = sd(brazil$new_cases)
31
32
33 infections[[9]][[2]] = length(boxplot.stats(chile$new_cases)$out)
34 infections[[8]][[2]] = sd(chile$new_cases)
35
36
37 infections[[9]][[3]] = length(boxplot.stats(venezuela$new_cases)$out)
38 infections[[8]][[3]] = sd(venezuela$new_cases)
39 # for covid new_deaths
40 deaths = data.frame(Countries)
41 deaths[["Min"]] <- NA
42 deaths[1, "Min"] = min(brazil$ new_deaths)
43 deaths[2, "Min"] = min(chile$ new_deaths)
44 deaths[3, "Min"] = min(venezuela$ new_deaths)
45 deaths[["Q1"]] = NA
46 deaths[["Q2"]] = NA
47 deaths[["Q3"]] = NA
48 deaths[["Max"]] = NA
49 deaths[["Avg"]] = NA
50 deaths[["Std"]] = NA
51 deaths[["Outlier"]] = NA
52
53 deaths[1, "Max"] = max(brazil$ new_deaths)
54 deaths[2, "Max"] = max(chile$ new_deaths)
55 deaths[3, "Max"] = max(venezuela$ new_deaths)
56 deaths[1, "Avg"] = mean(brazil$ new_deaths)
57 deaths[2, "Avg"] = mean(chile$ new_deaths)
58 deaths[3, "Avg"] = mean(venezuela$ new_deaths)
59 for (i in c(3,4,5)) {
60   deaths[[i]][[1]] = quantile(brazil$ new_deaths)[[i-1]]
61 }
62 for (i in c(3,4,5)) {
63   deaths[[i]][[2]] = quantile(chile$ new_deaths)[[i-1]]
64 }
65 for (i in c(3,4,5)) {
66   deaths[[i]][[3]] = quantile(venezuela$ new_deaths)[[i-1]]
67 }
68
69 deaths[[9]][[1]] = length(boxplot.stats(brazil$new_deaths)$out)
70 deaths[[8]][[1]] = sd(brazil$new_deaths)
71
72
73 deaths[[9]][[2]] = length(boxplot.stats(chile$new_deaths)$out)
74 deaths[[8]][[2]] = sd(chile$new_deaths)
75
76
77 deaths[[9]][[3]] = length(boxplot.stats(venezuela$new_deaths)$out)
78 deaths[[8]][[3]] = sd(venezuela$new_deaths)
```

- Thống kê cho các ca nhiễm

Countries	Min	Q1	Q2	Q3	Max	Avg	Std	Outlier
Brazil	0	13276.00	31149.0	53989.00	287149	38747.6979	37015.4040	26
Chile	0	1179.75	1997.5	4315.25	41651	3872.5907	6093.3562	38
Venezuela	0	245.00	590.0	1100.25	4418	721.1201	631.3701	15

- Thống kê cho các ca tử vong

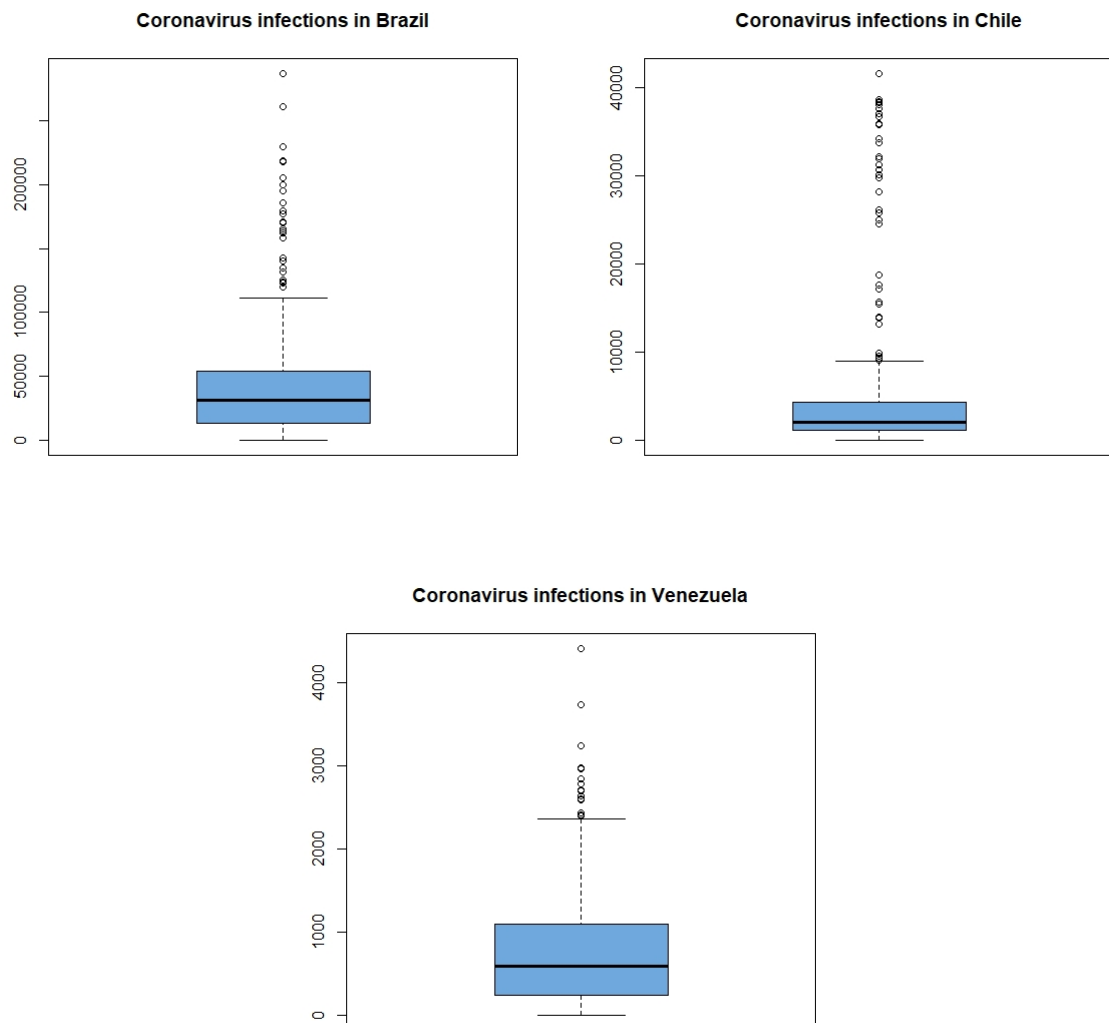
Countries	Min	Q1	Q2	Q3	Max	Avg	Std	Outlier
Brazil	0	317	721	1190	4148	888.544828	771.595312	40
Chile	0	17	40	82	1057	57.173077	67.164298	28
Venezuela	0	3	6	12	35	7.929379	7.214482	21

7) Vẽ biểu đồ boxplot hay còn được gọi là box-and-whisker cho nhiễm coronavirus

- Hiện thực trong R

```
1 boxplot(brazil$new_cases, main = "Coronavirus infections in Brazil", col = "#6fa8dc")
2 boxplot(chile$new_cases, main = "Coronavirus infections in Chile", col = "#6fa8dc")
3 boxplot(venezuela$new_cases, main = "Coronavirus infections in Venezuela", col = "#6fa8dc" )
```

- Biểu đồ boxplot cho các ca nhiễm:



iii) Nhóm câu hỏi liên quan đến dữ liệu thu thập

Với mỗi quốc gia cần tính số liệu thống kê lần lượt cho nhiễm và tử vong do coronavirus:

Chuẩn bị dữ liệu cho toàn bộ phần iii (phần code dùng chung cho cả phần iii)

```
1 library(tidyverse)
2 library(dplyr)
3 covid = read.csv("C:/Users/Admin/Downloads/covidData.csv")
4 covid = covid %>% filter(!continent == '') %>% mutate(new_cases = abs(new_cases), new_deaths = abs(
  new_deaths))
5 brazil = subset(covid, iso_code == "BRA")
6 chile = subset(covid, iso_code == "CHL")
7 venezuela = subset(covid, iso_code == "VEN")
8 brazil[is.na(brazil)]=0
9 chile[is.na(chile)]=0
10 venezuela[is.na(venezuela)]=0
11 Countries = c("Brazil", "Chile", "Venezuela")
```

1) Có bao nhiêu ngày có số lần dữ liệu không được báo cáo mới.

Note: Dữ liệu không được báo cáo mới là missing value (NA) hoặc số liệu bằng không

- Hướng giải quyết: Đếm các hàng có dữ liệu không được báo cáo mới.
- Hiện thực trong R

```
1 none_report = data.frame(Countries)
2 none_report[["New cases"]] = NA
3 none_report[["New deaths"]] = NA
4
5 none_report[[2]][[1]] = nrow(brazil %>% filter(new_cases ==0))
6 none_report[[2]][[2]] = nrow(chile %>% filter(new_cases ==0))
7 none_report[[2]][[3]] = nrow(venezuela %>% filter(new_cases ==0))
8
9 none_report[[3]][[1]] = nrow(brazil %>% filter(new_deaths ==0))
10 none_report[[3]][[2]] = nrow(chile %>% filter(new_deaths ==0))
11 none_report[[3]][[3]] = nrow(venezuela %>% filter(new_deaths ==0))
```

- Kết quả

	Countries	New cases	New deaths
1	Brazil	9	22
2	Chile	14	33
3	Venezuela	63	129

2) Có bao nhiêu ngày có số lần thu thập dữ liệu là thấp nhất được báo cáo mới.

- Hướng giải quyết: Đếm các hàng có dữ liệu bằng với giá trị nhỏ nhất của số ca nhiễm/ ca tử vong.
- Hiện thực trong R

```
1 iii2 = data.frame(Countries)
2 brazil = brazil %>% filter(!new_cases ==0, !new_deaths==0)
3 chile = chile %>% filter(!new_cases ==0, !new_deaths==0)
4 venezuela = venezuela %>% filter(!new_cases ==0, !new_deaths==0)
5 iii2[["New cases"]]=NA
6 iii2[["New deaths"]]=NA
7 iii2[[2]][[1]]=nrow(brazil %>% filter( new_cases == min(brazil$new_cases)))
8 iii2[[2]][[2]]=nrow(chile %>% filter( new_cases == min(chile$new_cases)))
9 iii2[[2]][[3]]=nrow(venezuela %>% filter( new_cases == min(venezuela$new_cases)))
10 iii2[[3]][[1]]=nrow(brazil %>% filter( new_deaths == min(brazil$new_deaths)))
11 iii2[[3]][[2]]=nrow(chile %>% filter( new_deaths == min(chile$new_deaths)))
12 iii2[[3]][[3]]=nrow(venezuela %>% filter( new_deaths == min(venezuela$new_deaths)))
```

- Kết quả

	Countries	New cases	New deaths
1	Brazil	1	1
2	Chile	1	11
3	Venezuela	3	14

3) Có bao nhiêu ngày có số lần thu thập dữ liệu là cao nhất được báo cáo mới

Tương tự câu 2

- Hiện thực trong R

```
1 iii3 = data.frame(Countries)
2 brazil = brazil %>% filter(!new_cases ==0, !new_deaths==0)
3 chile = chile %>% filter(!new_cases ==0, !new_deaths==0)
4 venezuela = venezuela %>% filter(!new_cases ==0, !new_deaths==0)
5 iii3[["New cases"]]=NA
6 iii3[["New deaths"]]=NA
7 iii3[[2]][[1]]=nrow(brazil %>% filter( new_cases == max(brazil$new_cases)))
8 iii3[[2]][[2]]=nrow(chile %>% filter( new_cases == max(chile$new_cases)))
9 iii3[[2]][[3]]=nrow(venezuela %>% filter( new_cases == max(venezuela$new_cases)))
10 iii3[[3]][[1]]=nrow(brazil %>% filter( new_deaths == max(brazil$new_deaths)))
11 iii3[[3]][[2]]=nrow(chile %>% filter( new_deaths == max(chile$new_deaths)))
12 iii3[[3]][[3]]=nrow(venezuela %>% filter( new_deaths == max(venezuela$new_deaths)))
```

- Kết quả

	Countries	New cases	New deaths
1	Brazil	1	1
2	Chile	1	1
3	Venezuela	1	3

- 4) Thể hiện bảng số liệu như sau:

Không được báo cáo mới:

Countries	Infections	Deaths
ctr_i	value	value

Báo cáo mới:

Countries	Infections	Deaths
ctr_i	value	value

- Hướng giải quyết: Đếm số ngày không có dữ liệu/ có dữ liệu về các ca nhiễm mới/ ca tử vong mới.
- Hiện thực trong R

```
1 unreport = data.frame(Countries)
2 brazil = subset(covid, iso_code == "BRA")
3 chile = subset(covid, iso_code == "CHL")
4 venezuela = subset(covid, iso_code == "VEN")
5 brazil[is.na(brazil)]=0
6 chile[is.na(chile)]=0
7 venezuela[is.na(venezuela)]=0
8 unreport[["Infections value"]] = NA
9 unreport[["Deaths value"]] = NA
10 unreport[[2]][[1]] = nrow(brazil %>% filter(new_cases ==0))
11 unreport[[2]][[2]] = nrow(chile %>% filter(new_cases ==0))
12 unreport[[2]][[3]] = nrow(venezuela %>% filter(new_cases ==0))
13 unreport[[3]][[1]] = nrow(brazil %>% filter(new_deaths ==0))
14 unreport[[3]][[2]] = nrow(chile %>% filter(new_deaths ==0))
15 unreport[[3]][[3]] = nrow(venezuela %>% filter(new_deaths ==0))
16
17 report = data.frame(Countries)
18 brazil = subset(covid, iso_code == "BRA")
19 chile = subset(covid, iso_code == "CHL")
20 venezuela = subset(covid, iso_code == "VEN")
21 brazil[is.na(brazil)]=0
22 chile[is.na(chile)]=0
23 venezuela[is.na(venezuela)]=0
24 report[["Infections value"]] = NA
25 report[["Deaths value"]] = NA
26 report[[2]][[1]] = nrow(brazil %>% filter(!new_cases ==0))
27 report[[2]][[2]] = nrow(chile %>% filter(!new_cases ==0))
28 report[[2]][[3]] = nrow(venezuela %>% filter(!new_cases ==0))
29 report[[3]][[1]] = nrow(brazil %>% filter(!new_deaths ==0))
30 report[[3]][[2]] = nrow(chile %>% filter(!new_deaths ==0))
31 report[[3]][[3]] = nrow(venezuela %>% filter(!new_deaths ==0))
```

- Không được báo cáo mới

Countries	Infections value	Deaths value
Brazil	9	22
Chile	14	33
Venezuela	63	129

- Báo cáo mới

Countries	Infections value	Deaths value
Brazil	716	703
Chile	714	695
Venezuela	645	579

5) Cho biết số ngày ngắn nhất liên tiếp mà không có dữ liệu được báo cáo

- Hướng giải quyết: Đánh dấu các dữ liệu khác không (tức là dữ liệu được báo cáo) bằng 1, giữ nguyên các dữ liệu bằng không, đếm số dòng có cùng giá trị (bằng 0 hoặc bằng 1) liên tiếp.
- Hiện thực trong R

```

1 #Nhưng ngày có dữ liệu thì danh dau =1, không có dữ liệu danh dau =0
2 bra = brazil
3 for (i in 1:nrow(bra)) {
4   if(bra[[5]][i] != 0 ){
5     bra[[5]][i] =1
6   }
7   if(bra[[6]][i] != 0 ){
8     bra[[6]][i] =1
9   }
10 }
11 chi = chile
12 for (i in 1:nrow(chi)) {
13   if(chi[[5]][i] != 0 ){
14     chi[[5]][i] =1
15   }
16   if(chi[[6]][i] != 0 ){
17     chi[[6]][i] =1
18   }
19 }
20 ven = venezuela
21 for (i in 1:nrow(ven)) {
22   if(ven[[5]][i] != 0 ){
23     ven[[5]][i] =1
24   }
25   if(ven[[6]][i] != 0 ){
26     ven[[6]][i] =1
27   }
28 }
29 iii5 = data.frame(Countries)
30 iii5[['New cases']]=NA
31 iii5[['New deaths']]=NA
32 iii5[[2]][1]= tapply(rle(bra$new_cases)$length, rle(bra$new_cases)$value, min)[1][1]
33 iii5[[2]][2]= tapply(rle(chi$new_cases)$length, rle(chi$new_cases)$value, min)[1][1]
34 iii5[[2]][3]= tapply(rle(ven$new_cases)$length, rle(ven$new_cases)$value, min)[1][1]
35 iii5[[3]][1]= tapply(rle(bra$new_deaths)$length, rle(bra$new_deaths)$value, min)[1][1]
36 iii5[[3]][2]= tapply(rle(chi$new_deaths)$length, rle(chi$new_deaths)$value, min)[1][1]
37 iii5[[3]][3]= tapply(rle(ven$new_deaths)$length, rle(ven$new_deaths)$value, min)[1][1]

```

- Kết quả: Số ngày ngắn nhất liên tiếp mà không có dữ liệu được báo cáo cho các ca nhiễm và các ca tử vong

	Countries	New cases	New deaths
1	Brazil	1	2
2	Chile	1	1
3	Venezuela	1	1

6) Cho biết số ngày dài nhất liên tiếp mà không có dữ liệu được báo cáo

Tương tự câu 5, thay đổi min thành max ở dòng cuối

```

1 iii6 = data.frame(Countries)
2 iii6[['New cases']]=NA

```



```

3 iii6[['New deaths']]=NA
4 iii6[[2]][[1]]= apply(rle(bra$new_cases)$length, rle(bra$new_cases)$value, max)[[1]][[1]]
5 iii6[[2]][[2]]= apply(rle(chi$new_cases)$length, rle(chi$new_cases)$value, max)[[1]][[1]]
6 iii6[[2]][[3]]= apply(rle(ven$new_cases)$length, rle(ven$new_cases)$value, max)[[1]][[1]]
7 iii6[[3]][[1]]= apply(rle(bra$new_deaths)$length, rle(bra$new_deaths)$value, max)[[1]][[1]]
8 iii6[[3]][[2]]= apply(rle(chi$new_deaths)$length, rle(chi$new_deaths)$value, max)[[1]][[1]]
9 iii6[[3]][[3]]= apply(rle(ven$new_deaths)$length, rle(ven$new_deaths)$value, max)[[1]][[1]]

```

- Kết quả: Số ngày dài nhất liên tiếp mà không có dữ liệu được báo cáo cho các ca nhiễm và các ca tử vong

	Countries	New cases	New deaths
1	Brazil	3	20
2	Chile	6	28
3	Venezuela	1	24

- 7) Cho biết số ngày ngắn nhất liên tiếp mà không có người nhiễm bệnh mới
Tương tự hai câu trên

- Hiện thực trong R

```

1 iii7 = data.frame(Countries)
2 iii7[['New cases']]=NA
3 iii7[[2]][[1]]= apply(rle(bra$new_cases)$length, rle(bra$new_cases)$value, min)[[1]][[1]]
4 iii7[[2]][[2]]= apply(rle(chi$new_cases)$length, rle(chi$new_cases)$value, min)[[1]][[1]]
5 iii7[[2]][[3]]= apply(rle(ven$new_cases)$length, rle(ven$new_cases)$value, min)[[1]][[1]]

```

- Kết quả: Số ngày ngắn nhất liên tiếp mà không có người nhiễm bệnh mới

	Countries	New cases
1	Brazil	1
2	Chile	1
3	Venezuela	1

- 8) Cho biết số ngày dài nhất liên tiếp mà không có người nhiễm bệnh mới
Tương tự câu trên (thay min thành max)

- Hiện thực trong R

```

1 iii8 = data.frame(Countries)
2 iii8[['New cases']]=NA
3 iii8[[2]][[1]]= apply(rle(bra$new_cases)$length, rle(bra$new_cases)$value, max)[[1]][[1]]
4 iii8[[2]][[2]]= apply(rle(chi$new_cases)$length, rle(chi$new_cases)$value, max)[[1]][[1]]
5 iii8[[2]][[3]]= apply(rle(ven$new_cases)$length, rle(ven$new_cases)$value, max)[[1]][[1]]

```

- Kết quả: Số ngày dài nhất liên tiếp mà không có người nhiễm bệnh mới

	Countries	New cases
1	Brazil	3
2	Chile	6
3	Venezuela	1

iv) Nhóm câu hỏi liên quan đến trực quan dữ liệu

- 1) Vẽ biểu đồ tần số tích lũy quốc gia cho các châu lục

- Tần số tích lũy: Trong một tập dữ liệu, tần số tích lũy đối với một giá trị x là tổng số điểm mà nhỏ hơn hoặc bằng x .
- Biểu đồ tần số tích lũy: biểu thị những thông tin dạng tích lũy. Nó thể hiện số lượng hay tỉ lệ những quan sát nhỏ hơn hoặc bằng một giá trị cụ thể.
- Hiện thực trong R

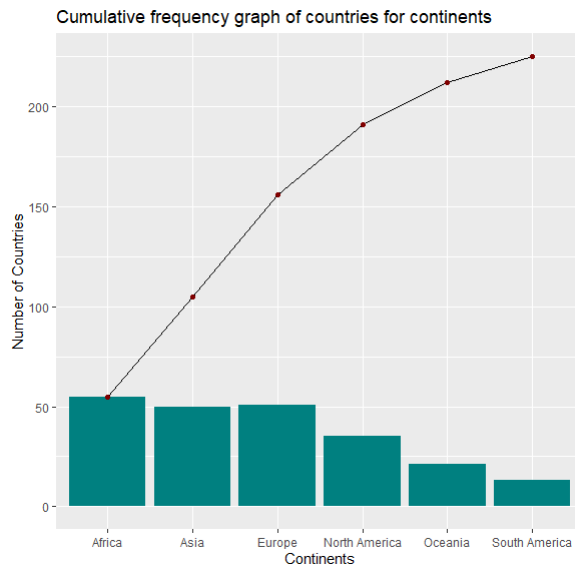
```

1 continents = covid %>% filter(!continent == '') %>% select(iso_code:location) %>% unique() %>%
  group_by(continent) %>% summarize(n = n())
2 ggplot(data, aes(x = continent, y = n, group = 1)) +
3   xlab("Continents") +
4   ylab("Number of Countries") +

```

```
5 geom_line(aes(y = cumsum(n))) + geom_bar(stat="identity",fill = "#008080" ) +
6 geom_point(aes(x = continent, y = cumsum(n)), color = "#800000") +
7 labs(title = "Cumulative frequency graph of countries for continents" , substitute= "Tan so
8 tích lũy")
```

- Kết quả:

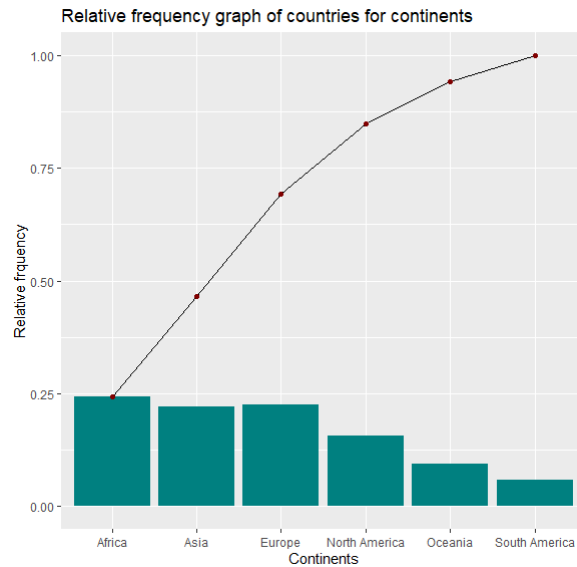


2) Vẽ biểu đồ tần số tương đối quốc gia cho các châu lục

- Biểu đồ tần số tương đối: Việc đếm tần số có thể được thể hiện dưới dạng số tuyệt đối hoặc tương đối (ví dụ phân số hoặc tỉ lệ phần trăm). Biểu đồ dưới đây thể hiện nội dung y nguyên như biểu đồ tần số tích lũy ở trên, chỉ khác một điều đơn vị tính của các cột là tỉ lệ phần trăm.
- Hiện thực trong R

```
1 continents =continents %>% mutate("m" = n/sum(n))
2 ggplot(data = continents, aes(x = continent, y = m, group =1)) +
3   geom_line(aes(y = cumsum(m))) + geom_bar(stat="identity",fill = "#008080" )+
4   xlab("Continents") +
5   ylab("Relative frequency") +
6   geom_point(aes(x = continent, y = cumsum(m)), color = "#800000") +
7   labs(title = "Relative frequency graph of countries for continents" , substitute= "Tan so
   tương doi" )
```

- Kết quả:

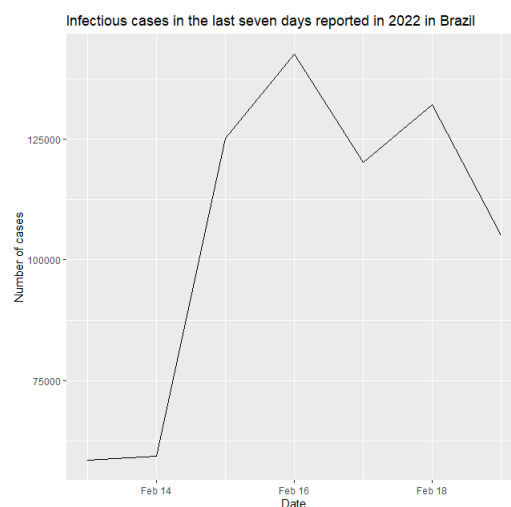


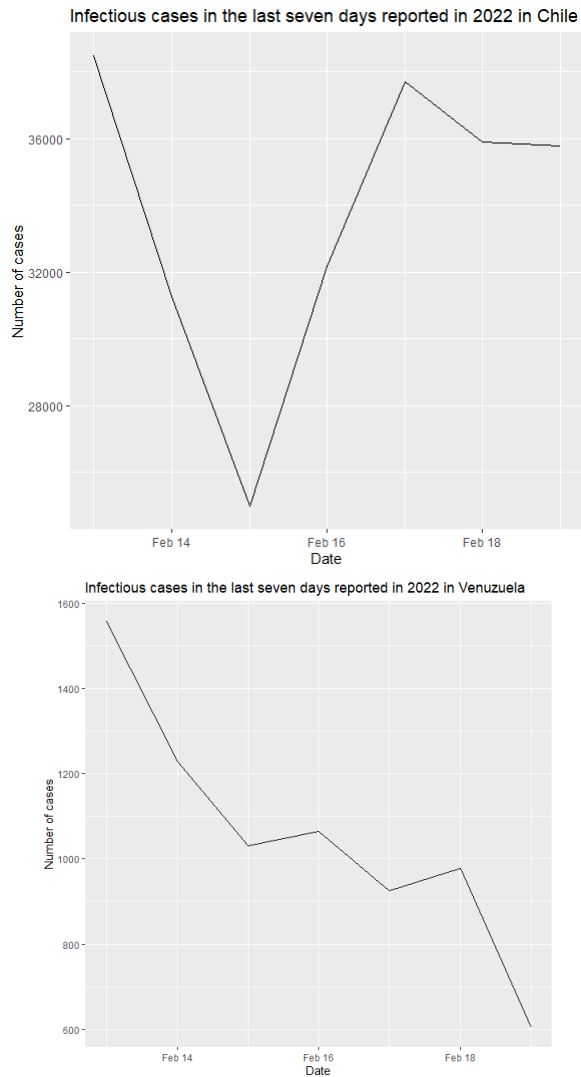
3) Vẽ biểu đồ thể hiện nhiễm bệnh đã báo cáo của các quốc gia trong 7 ngày cuối của năm cuối cùng

- Hướng giải quyết: Lọc dữ liệu của 7 ngày cuối được báo cáo của năm cuối. Vẽ biểu đồ dựa trên dữ liệu đó

```
1 covid$date = as.Date(covid$date, format = "%m/%d/%Y")
2 bra = brazil[order(brazil$date, decreasing = TRUE),]
3 bra = bra[1:7,] #select data from the last 7 days
4 ggplot(data= bra, aes(x= date, y =new_cases)) + geom_line()+ labs(title = "Infectious cases in
the last seven days reported in 2022 in Brazil", x="Date", y="Number of cases")
5
6 chi = chile[order(chile$date, decreasing = TRUE),]
7 chi = chi[1:7,] #select the last 7 days
8 ggplot(data= chi, aes(x= date, y =new_cases)) + geom_line()+ labs(title = "Infectious cases in
the last seven days reported in 2022 in Chile", x="Date", y="Number of cases")
9
10 ven = venezuela[order(venezuela$date, decreasing = TRUE),]
11 ven = ven[1:7,] #select the last 7 days
12 ggplot(data= ven, aes(x= date, y =new_cases)) + geom_line()+ labs(title = "Infectious cases in
the last seven days reported in 2022 in Venuzuela", x="Date", y="Number of cases")
```

- Kết quả:

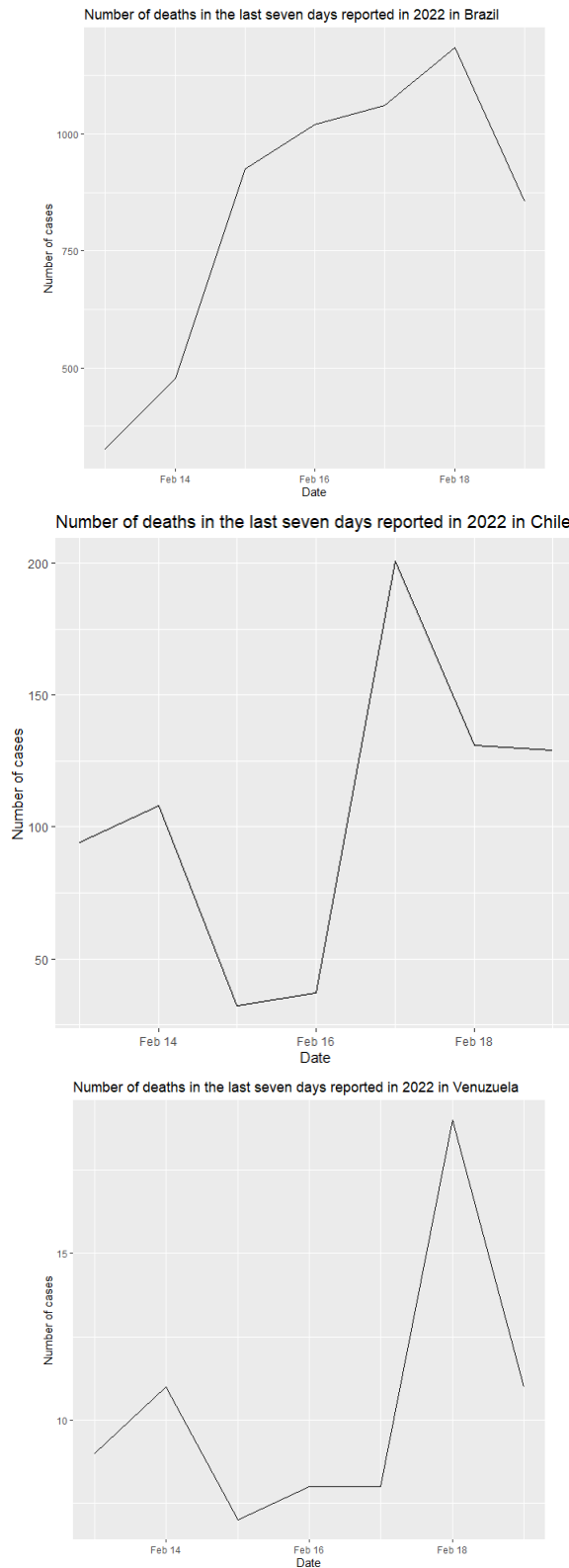




4) Vẽ biểu đồ thể hiện tử vong đã báo cáo của các quốc gia trong 7 ngày cuối của năm cuối cùng
*Tương tự câu trên*Hiện thực trong R

- `ggplot(data= bra, aes(x= date, y =new_deaths)) + geom_line()+ labs(title = "Number of deaths in the last seven days reported in 2022 in Brazil", x="Date", y="Number of cases")`
- 2 `ggplot(data= chi, aes(x= date, y =new_deaths)) + geom_line()+ labs(title = "Number of deaths in the last seven days reported in 2022 in Chile", x="Date", y="Number of cases")`
- 3 `ggplot(data= ven, aes(x= date, y =new_deaths)) + geom_line()+ labs(title = "Number of deaths in the last seven days reported in 2022 in Venezuela", x="Date", y="Number of cases")`

- Kết quả:



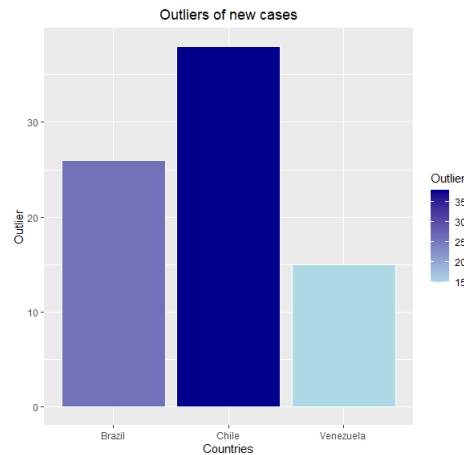
5) Vẽ biểu đồ phổ đất nước xuất hiện outliers cho nhiễm bệnh

- Dùng thông tin về Outliers ở phần ii) để vẽ biểu đồ
- Hiện thực trong R

```
1 ggplot(data = df_cases, aes(x= Countries, y = Outlier)) +
2   geom_histogram(stat='identity',col="white", aes(fill = Outlier), binwidth = 7) +
3   scale_fill_gradient("Outlier",low= "lightblue", high = "darkblue") +ggtitle("Outliers of
```

```
new cases") +  
4 theme(plot.title = element_text(hjust = 0.5))
```

- Kết quả:

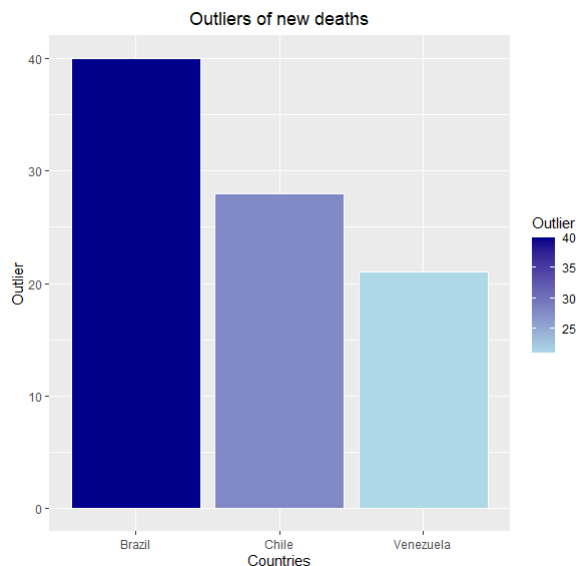


6) Vẽ biểu đồ phổ đất nước xuất hiện outliers cho tử vong
Tương tự câu trên

- Hiện thực trong R

```
1 ggplot(data = df_death, aes(x= Countries, y = Outlier)) +  
2 geom_histogram(stat='identity', col="white", aes(fill = Outlier), binwidth = 7) +  
3 scale_fill_gradient("Outlier", low= "lightblue", high = "darkblue") + ggtitle("Outliers of  
new deaths") +  
4 theme(plot.title = element_text(hjust = 0.5))
```

- Kết quả:



v) **Nhóm câu hỏi liên quan đến trực quan dữ liệu theo thời gian là tháng**

Trên từng năm hãy vẽ biểu đồ thể hiện trục Ox là thời gian, trục Oy là nhiễm bệnh/tử vong. Hãy dùng 4 ký số của mã đề để vẽ 4 tháng tương ứng theo ký số đó. Nếu ký số là 0 thì lấy tháng là 10.

- Chuẩn bị dữ liệu cho toàn bộ phần v (phần code dùng chung cho cả phần v)

```
1 library(tidyverse)  
2 library(lubridate)  
3 library(dplyr)  
4 library(xlsx)  
5 setwd("E:/BTL_CTRR")
```

```
6 data_covid = read.csv("owid-covid-data.csv", header = TRUE)
7 library("ggplot2")
8 data <- data_covid %>%
9   separate(date, into = c("month", "day", "year"), sep = "/")
10 data = mutate(data, new_cases = abs(new_cases))
11 data = mutate(data, new_deaths = abs(new_deaths))
12 data <- subset(data, data$continent != "")
13 data[is.na(data)] <- 0
14 data <- transform(data, month=as.numeric(month))
15
16 year_2020_brazil <- data[data$location == "Brazil" & data$year=="2020",]
17 year_2020_chile <- data[data$location == "Chile" & data$year=="2020",]
18 year_2020_ven <- data[data$location == "Venezuela" & data$year=="2020",]
19
20 year_2021_brazil <- data[data$location == "Brazil"&data$year=="2021",]
21 year_2021_chile <- data[data$location == "Chile"&data$year=="2021",]
22 year_2021_ven <- data[data$location == "Venezuela"&data$year=="2021",]
23
24 year_2022_brazil <- data[data$location == "Brazil"&data$year=="2022",]
25 year_2022_chile <- data[data$location == "Chile"&data$year=="2022",]
26 year_2022_ven <- data[data$location == "Venezuela"&data$year=="2022",]
27
28 bra_2020 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2020_brazil, FUN = sum)
29 chi_2020 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2020_chile, FUN = sum)
30 ven_2020 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2020_ven, FUN = sum)
31
32 bra_2021 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2021_brazil, FUN = sum)
33 chi_2021 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2021_chile, FUN = sum)
34 ven_2021 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2021_ven, FUN = sum)
35
36 bra_2022 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2022_brazil, FUN = sum)
37 chi_2022 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2022_chile, FUN = sum)
38 ven_2022 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2022_ven, FUN = sum)
39
40 bra_2020_1<-subset(bra_2020,bra_2020$month==2 | bra_2020$month==10 |
41                   bra_2020$month==3 | bra_2020$month==6)
42 bra_2021_1<-subset(bra_2021,bra_2021$month==2 | bra_2021$month==10 |
43                   bra_2021$month==3 | bra_2021$month==6)
44 bra_2022_1<-subset(bra_2022,bra_2022$month==2 | bra_2022$month==10 |
45                   bra_2022$month==3 | bra_2022$month==6)
46
47 chi_2020_1<-subset(chi_2020,chi_2020$month==2 | chi_2020$month==10 |
48                   chi_2020$month==3 | chi_2020$month==6)
49 chi_2021_1<-subset(chi_2021,chi_2021$month==2 | chi_2021$month==10 |
50                   chi_2021$month==3 | chi_2021$month==6)
51 chi_2022_1<-subset(chi_2022,chi_2022$month==2 | chi_2022$month==10 |
52                   chi_2022$month==3 | chi_2022$month==6)
53
54 ven_2020_1<-subset(ven_2020,ven_2020$month==2 | ven_2020$month==10 |
55                   ven_2020$month==3 | ven_2020$month==6)
56 ven_2021_1<-subset(ven_2021,ven_2021$month==2 | ven_2021$month==10 |
57                   ven_2021$month==3 | ven_2021$month==6)
58 ven_2022_1<-subset(ven_2022,ven_2022$month==2 | ven_2022$month==10 |
59                   ven_2022$month==3 | ven_2022$month==6)
60
61 bra_2020_2<-subset(bra_2020,bra_2020$month==11 | bra_2020$month==12)
62 bra_2021_2<-subset(bra_2021,bra_2021$month==11 | bra_2021$month==12)
63
64 chi_2020_2<-subset(chi_2020,chi_2020$month==11 | chi_2020$month==12)
65 chi_2021_2<-subset(chi_2021,chi_2021$month==11 | chi_2021$month==12)
66
67 ven_2020_2<-subset(ven_2020,ven_2020$month==11 | ven_2020$month==12)
68 ven_2021_2<-subset(ven_2021,ven_2021$month==11 | ven_2021$month==12)
69 new_row=c(2,0,0)
70 ven_2020_1 <- rbind(ven_2020_1,new_row)
71 temp=list(bra_2020_1,chi_2020_1,ven_2020_1)
72 data_2020=temp %>% reduce(full_join, by='month')
73 colnames(data_2020) <- c("month","brazil_cases","brazil_deaths","chile_cases","chile_deaths","
74   vene_cases","vene_deaths")
75 temp=list(bra_2021_1,chi_2021_1,ven_2021_1)
```

```

76 data_2021=temp %>% reduce(full_join, by='month')
77 colnames(data_2021) <- c("month", "brazil_cases", "brazil_deaths", "chile_cases", "chile_deaths", "
  vene_cases", "vene_deaths")
78
79 temp=list(bra_2022_1,chi_2022_1,ven_2022_1)
80 data_2022=temp %>% reduce(full_join, by='month')
81 colnames(data_2022) <- c("month", "brazil_cases", "brazil_deaths", "chile_cases", "chile_deaths", "
  vene_cases", "vene_deaths")
82
83 temp=list(bra_2020_2,chi_2020_2,ven_2020_2)
84 data_2020_1=temp %>% reduce(full_join, by='month')
85 colnames(data_2020_1) <- c("month", "brazil_cases", "brazil_deaths", "chile_cases", "chile_deaths", "
  vene_cases", "vene_deaths")
86
87 temp=list(bra_2021_2,chi_2021_2,ven_2021_2)
88 data_2021_1=temp %>% reduce(full_join, by='month')
89 colnames(data_2021_1) <- c("month", "brazil_cases", "brazil_deaths", "chile_cases", "chile_deaths", "
  vene_cases", "vene_deaths")

```

1) Biểu đồ thu thập nhiễm bệnh cho từng tháng

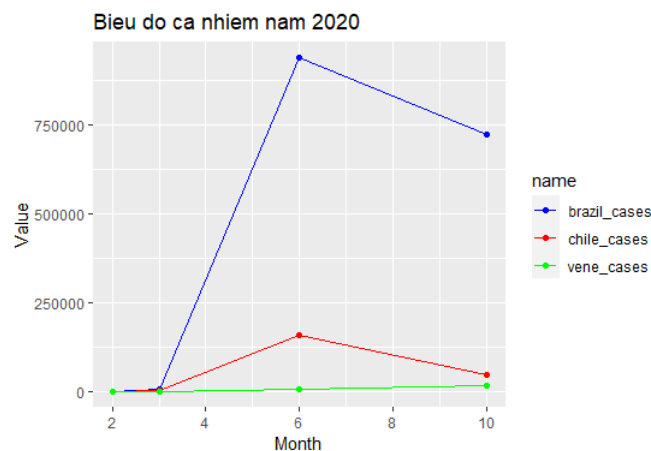
- Hiện thực trong R

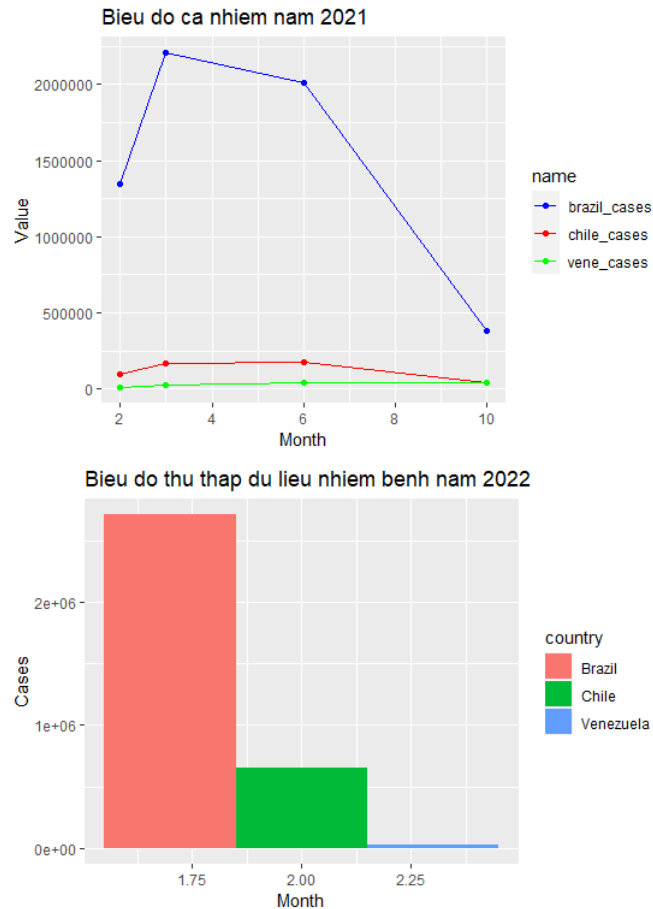
```

1 temp <- data_2020%>%
2   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
3 ggplot(data=temp,aes(x=month, y=value,color=name))+
4   geom_line()+
5   geom_point()+
6   scale_color_manual(values =c("blue","red","green"))+
7   labs(title="Biểu đồ ca nhiễm nam 2020",x="Month",y="Value")
8
9 temp <- data_2021%>%
10  tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
11 ggplot(data=temp,aes(x=month, y=value,color=name))+
12  geom_line()+
13  geom_point()+
14  scale_color_manual(values =c("blue","red","green"))+
15  labs(title="Biểu đồ ca nhiễm nam 2021",x="Month",y="Value")
16
17 V_1_2022<- data.frame(month=c(2,2,2),
18   country=c('Brazil','Chile','Venezuela'),
19   cases=c(bra_2022_1$new_cases,chi_2022_1$new_cases,ven_2022_1$new_cases))
20 ggplot(V_1_2022, aes(fill=country, y=cases, x=month)) +
21   geom_bar(position='dodge', stat='identity') +
22   labs(x='Month', y='Cases', title='Biểu đồ thu thập dữ liệu từ vòng năm 2022')

```

- Kết quả:



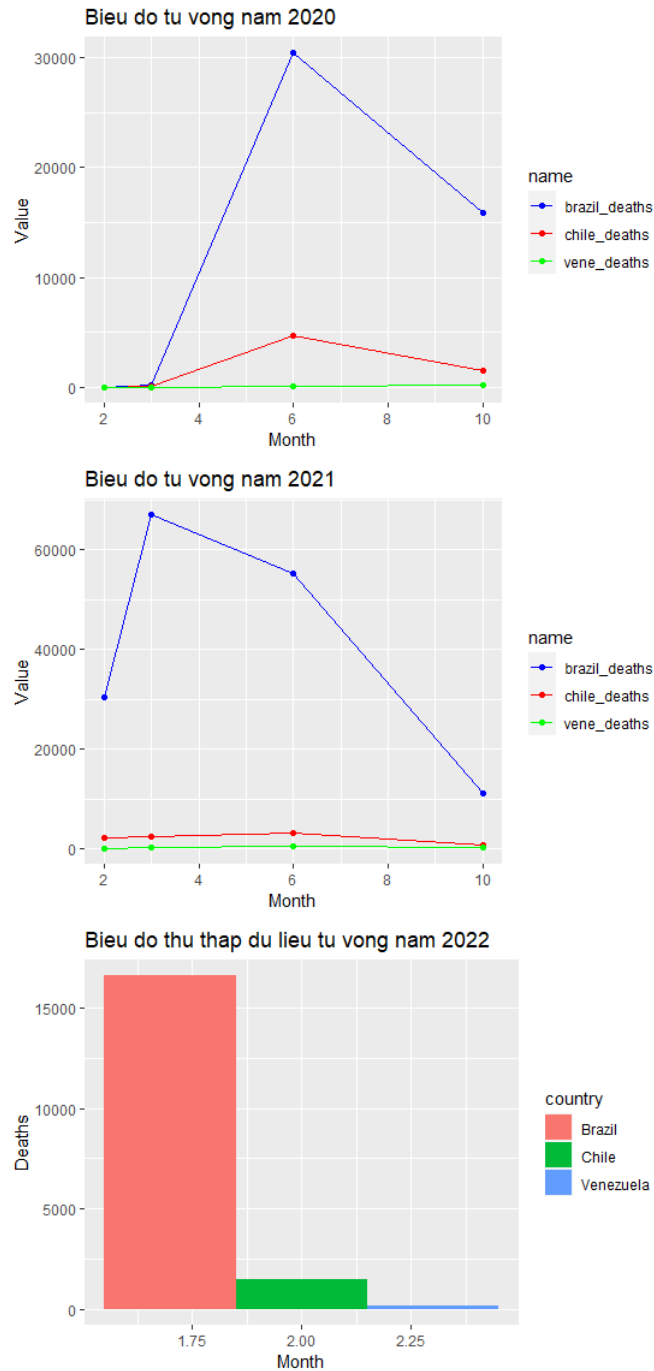


2) Biểu đồ thu thập tử vong cho từng tháng

- Hiện thực trong R

```
1 temp <- data_2020%>%
2   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
3   ggplot(data=temp,aes(x=month, y=value,color=name))+
4     geom_line()+
5     geom_point()+
6     scale_color_manual(values =c("blue","red","green"))+
7     labs(title="Bieu do tu vong nam 2020",x="Month",y="Value")
8
9 temp <- data_2021%>%
10   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
11   ggplot(data=temp,aes(x=month, y=value,color=name))+
12     geom_line()+
13     geom_point()+
14     scale_color_manual(values =c("blue","red","green"))+
15     labs(title="Bieu do tu vong nam 2021",x="Month",y="Value")
16
17
18 V_2_2022<-data.frame(month=c(2,2,2),
19   country=c('Brazil','Chile','Venezuela'),
20   deaths=c(bra_2022_1$new_deaths,chi_2022_1$new_deaths,ven_2022_1$
21     new_deaths))
22   ggplot(V_2_2022, aes(fill=country, y=deaths, x=month)) +
23     geom_bar(position='dodge', stat='identity') +
24     labs(x='Month', y='Deaths', title='Bieu do thu thap du lieu tu vong nam 2022')
```

- Kết quả:



3) Biểu đồ thu thập gồm nhiễm bệnh và tử vong cho từng tháng

- Hiện thực trong R

```

1 temp<- data_2020%>%
2   tidyr::pivot_longer(cols = c(
3     brazil_cases,brazil_deaths,chile_cases,chile_deaths,vene_cases,vene_deaths))
4   ggplot(data=temp,aes(x=month, y=value,color=name))+
5     geom_line()+
6     geom_point()+
7     scale_color_manual(values =c("blue","red","green","black","violet","orange"))+
8     labs(title="Bieu do ca nhien va tu vong nam 2020",x="Month",y="Value")
9
10 temp<- data_2021%>%
11   tidyr::pivot_longer(cols = c(
12     brazil_cases,brazil_deaths,chile_cases,chile_deaths,vene_cases,vene_deaths))
13   ggplot(data=temp,aes(x=month, y=value,color=name))+

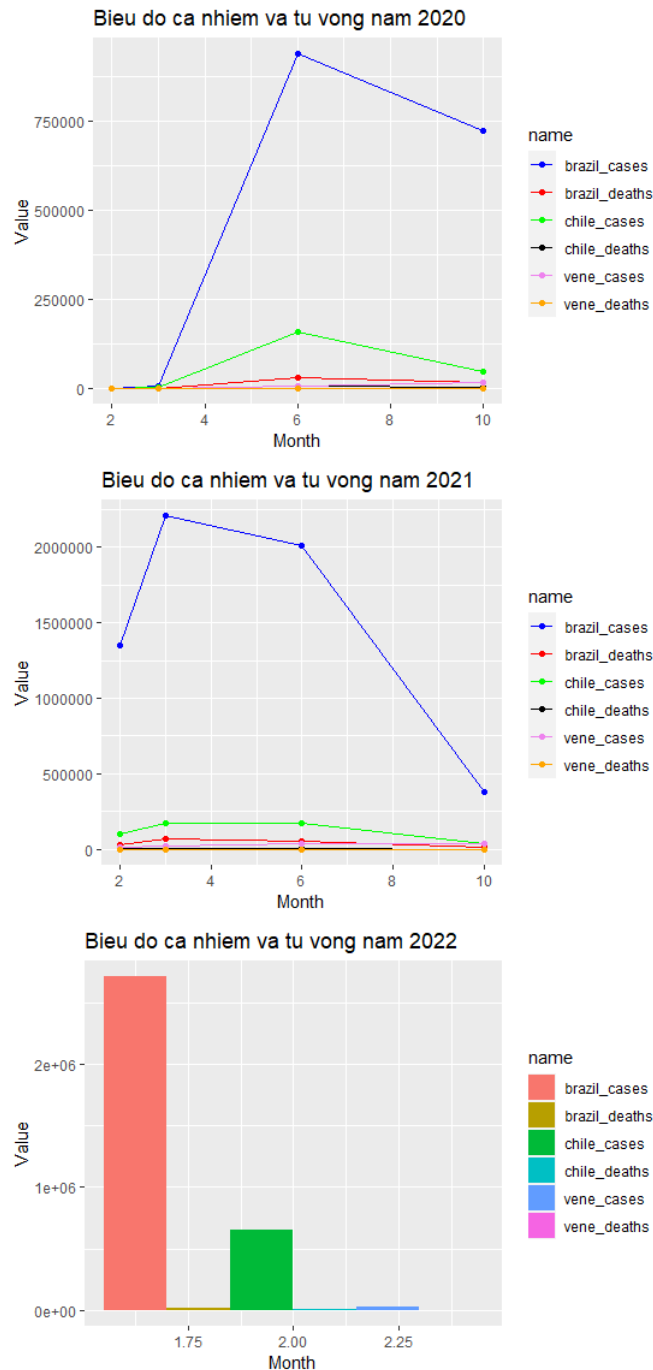
```

```

12 geom_line()+
13 geom_point()+
14 scale_color_manual(values =c("blue","red","green","black","violet","orange"))+
15 labs(title="Bieu do ca nhien va tu vong nam 2021",x="Month",y="Value")
16
17 temp<- data_2022%>%
18 tidyr::pivot_longer(cols = c(
19   brazil_cases,brazil_deaths,chile_cases,chile_deaths,vene_cases,vene_deaths))
20 ggplot(temp, aes(fill=name, y=value, x=month)) +
21   geom_bar(position='dodge', stat='identity') +
22   labs(x='Month', y='Value', title='Bieu do ca nhien va tu vong nam 2022')

```

- Kết quả:



4) Biểu đồ thu thập nhiễm bệnh gồm 2 tháng cuối của năm

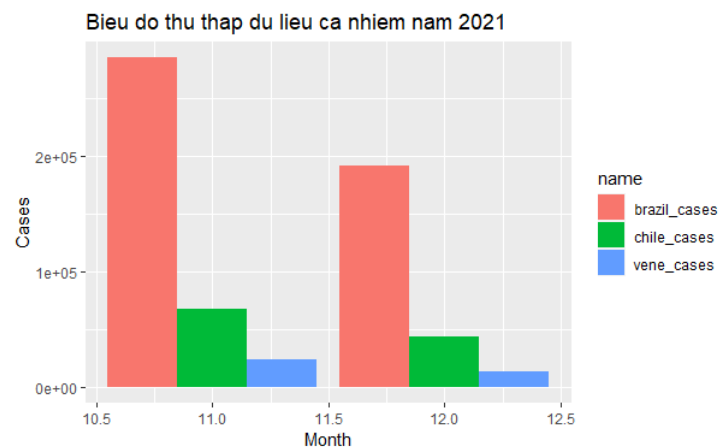
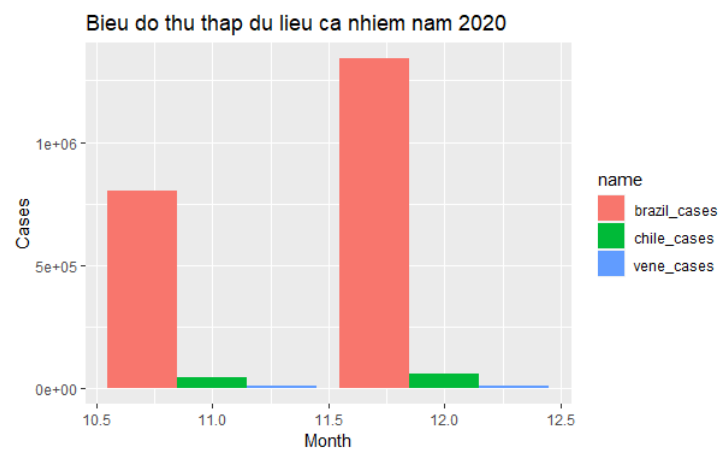
- Hiện thực trong R

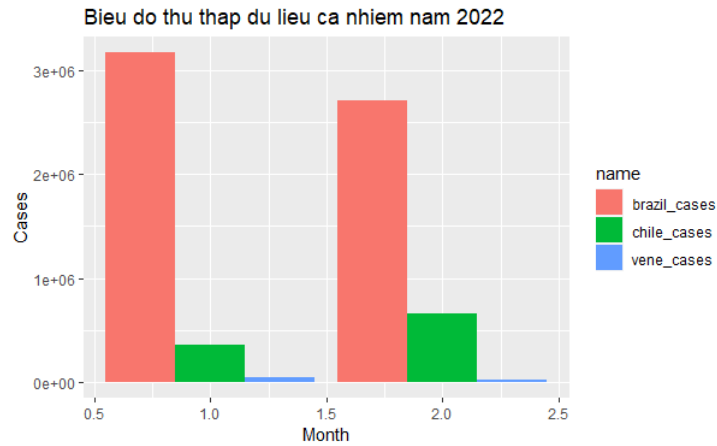
```

1 temp<- data_2020_1%>%
2   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
3   ggplot(temp, aes(fill=name, y=value, x=month)) +
4     geom_bar(position='dodge', stat='identity') +
5     labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem nam 2020')
6
7 temp<- data_2021_1%>%
8   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
9   ggplot(temp, aes(fill=name, y=value, x=month)) +
10    geom_bar(position='dodge', stat='identity') +
11    labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem nam 2021')
12
13 temp<- data_2022_1%>%
14   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
15   ggplot(temp, aes(fill=name, y=value, x=month)) +
16     geom_bar(position='dodge', stat='identity') +
17     labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem nam 2022')

```

- Kết quả:



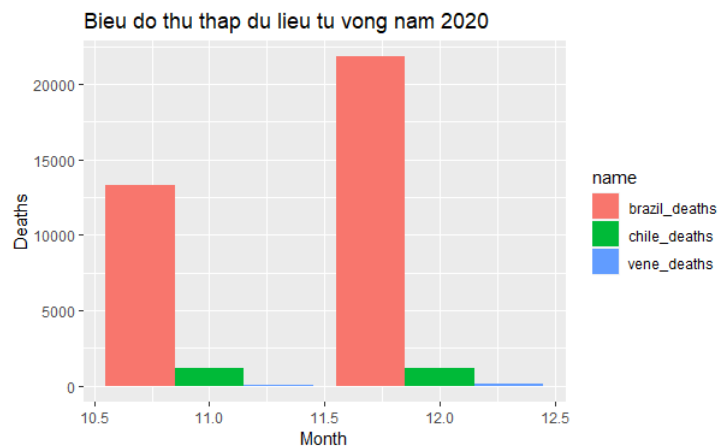


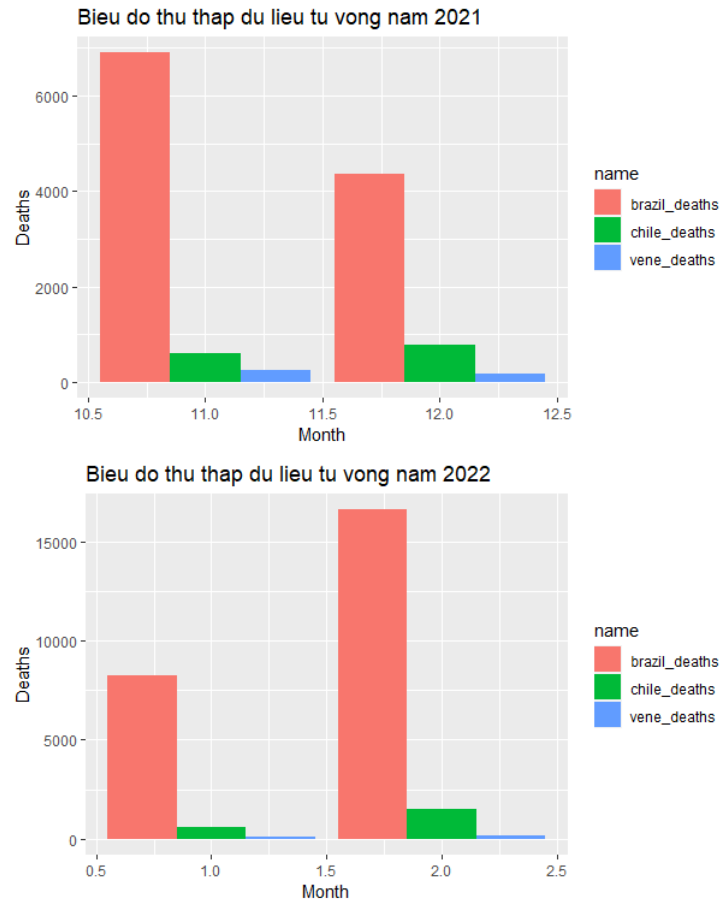
5) Biểu đồ thu thập tử vong gồm 2 tháng cuối của năm

- Hiện thực trong R

```
1 temp <- data_2020_1%>%
2   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
3   ggplot(temp, aes(fill=name, y=value, x=month)) +
4     geom_bar(position='dodge', stat='identity') +
5     labs(x='Month', y='Deaths', title='Bieu do thu thap du lieu tu vong nam 2020')
6
7 temp <- data_2021_1%>%
8   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
9   ggplot(temp, aes(fill=name, y=value, x=month)) +
10    geom_bar(position='dodge', stat='identity') +
11    labs(x='Month', y='Deaths', title='Bieu do thu thap du lieu tu vong nam 2021')
12
13 temp <- data_2022_1%>%
14   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
15   ggplot(temp, aes(fill=name, y=value, x=month)) +
16     geom_bar(position='dodge', stat='identity') +
17     labs(x='Month', y='Deaths', title='Bieu do thu thap du lieu tu vong nam 2022')
```

- Kết quả:



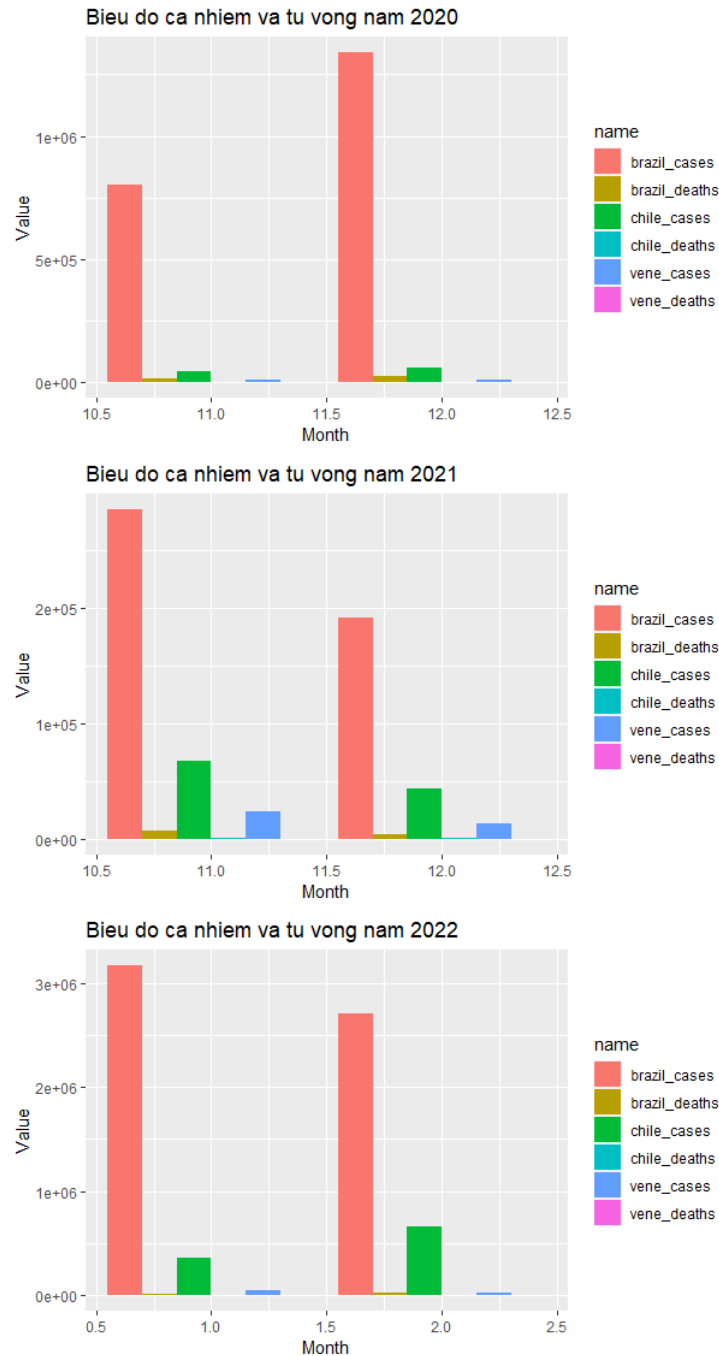


6) Biểu đồ thu thập gồm nhiễm bệnh và tử vong gồm 2 tháng cuối của năm

- Hiện thực trong R

```
1 temp<- data_2020_1%>%
2   tidyr::pivot_longer(cols = c(
3     brazil_cases,brazil_deaths,chile_cases,chile_deaths,vene_cases,vene_deaths))
4   ggplot(temp, aes(fill=name, y=value, x=month)) +
5     geom_bar(position='dodge', stat='identity') +
6     labs(x='Month', y='Value', title='Bieu do ca nhiem va tu vong nam 2020')
7
8 temp<- data_2021_1%>%
9   tidyr::pivot_longer(cols = c(
10    brazil_cases,brazil_deaths,chile_cases,chile_deaths,vene_cases,vene_deaths))
11   ggplot(temp, aes(fill=name, y=value, x=month)) +
12     geom_bar(position='dodge', stat='identity') +
13     labs(x='Month', y='Value', title='Bieu do ca nhiem va tu vong nam 2021')
14
15 temp<- data_2022_1%>%
16   tidyr::pivot_longer(cols = c(
17    brazil_cases,brazil_deaths,chile_cases,chile_deaths,vene_cases,vene_deaths))
18   ggplot(temp, aes(fill=name, y=value, x=month)) +
19     geom_bar(position='dodge', stat='identity') +
20     labs(x='Month', y='Value', title='Bieu do ca nhiem va tu vong nam 2022')
```

- Kết quả:



7) Biểu đồ thu thập nhiễm bệnh tích lũy cho từng tháng

- Hiện thực trong R

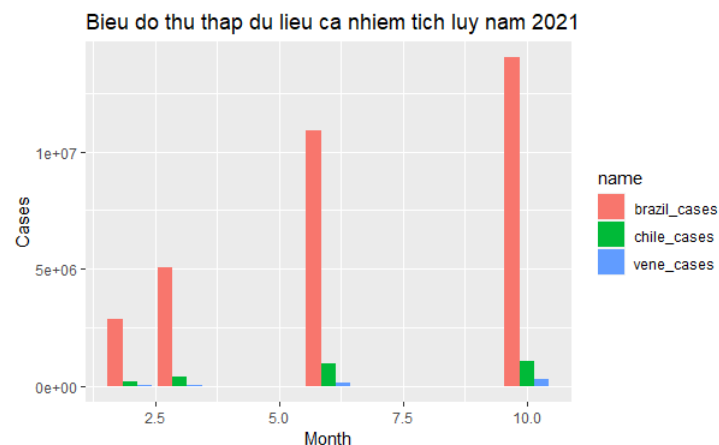
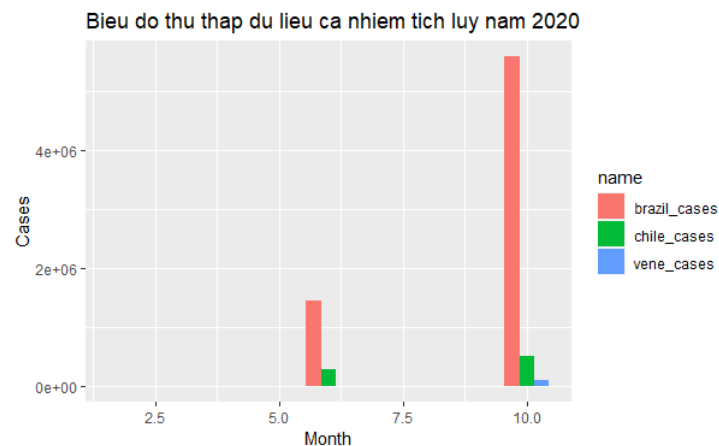
```
1 temp=list(bra_2020,chi_2020,ven_2020)
2 sum_2020=temp %>% reduce(full_join, by='month')
3 colnames(sum_2020) <- c("month","brazil_cases","brazil_deaths","chile_cases","chile_deaths","
4 vene_cases","vene_deaths")
5 sum_2020[is.na(sum_2020)] <- 0
6
7 temp=list(bra_2021,chi_2021,ven_2021)
8 sum_2021=temp %>% reduce(full_join, by='month')
9 colnames(sum_2021) <- c("month","brazil_cases","brazil_deaths","chile_cases","chile_deaths","
10 vene_cases","vene_deaths")
11
12 temp=list(bra_2022,chi_2022,ven_2022)
13 sum_2022=temp %>% reduce(full_join, by='month')
```

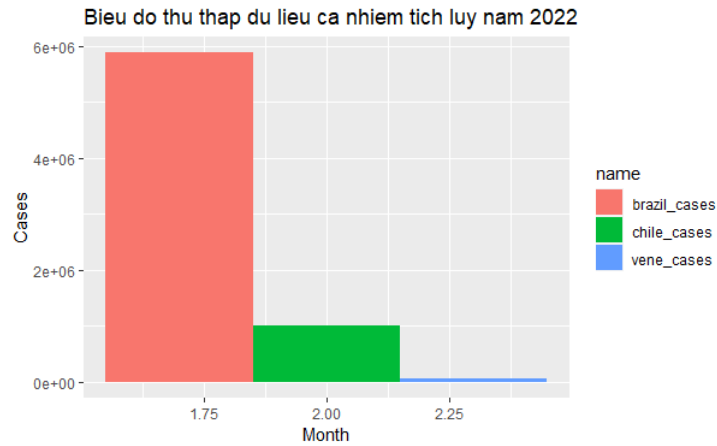
```

12 colnames(sum_2022) <- c("month","brazil_cases","brazil_deaths","chile_cases","chile_deaths","
    vene_cases","vene_deaths")
13
14 cumsum_2020=as.data.frame(lapply(sum_2020,cumsum))
15 cumsum_2020$month = c(2:12)
16 cumsum_2021=as.data.frame(lapply(sum_2021,cumsum))
17 cumsum_2021$month =c(1:12)
18 cumsum_2022=as.data.frame(lapply(sum_2022,cumsum))
19 cumsum_2022$month = c(1,2)
20
21 temp<- cumsum_2020%>%
22   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
23 temp <- subset(temp,temp$month==2 | temp$month==10 |
24   temp$month==3 | temp$month==6)
25 ggplot(temp, aes(fill=name, y=value, x=month)) +
26   geom_bar(position='dodge', stat='identity') +
27   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2020')
28
29 temp<- cumsum_2021%>%
30   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
31 temp <- subset(temp,temp$month==2 | temp$month==10 |
32   temp$month==3 | temp$month==6)
33 ggplot(temp, aes(fill=name, y=value, x=month)) +
34   geom_bar(position='dodge', stat='identity') +
35   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2021')
36
37 temp<- cumsum_2022%>%
38   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,vene_cases))
39 temp <- subset(temp,temp$month==2 | temp$month==10 |
40   temp$month==3 | temp$month==6)
41 ggplot(temp, aes(fill=name, y=value, x=month)) +
42   geom_bar(position='dodge', stat='identity') +
43   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2022')

```

● Kết quả:



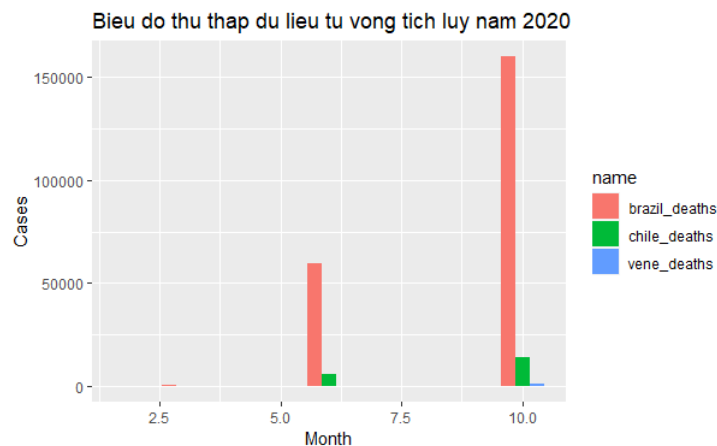


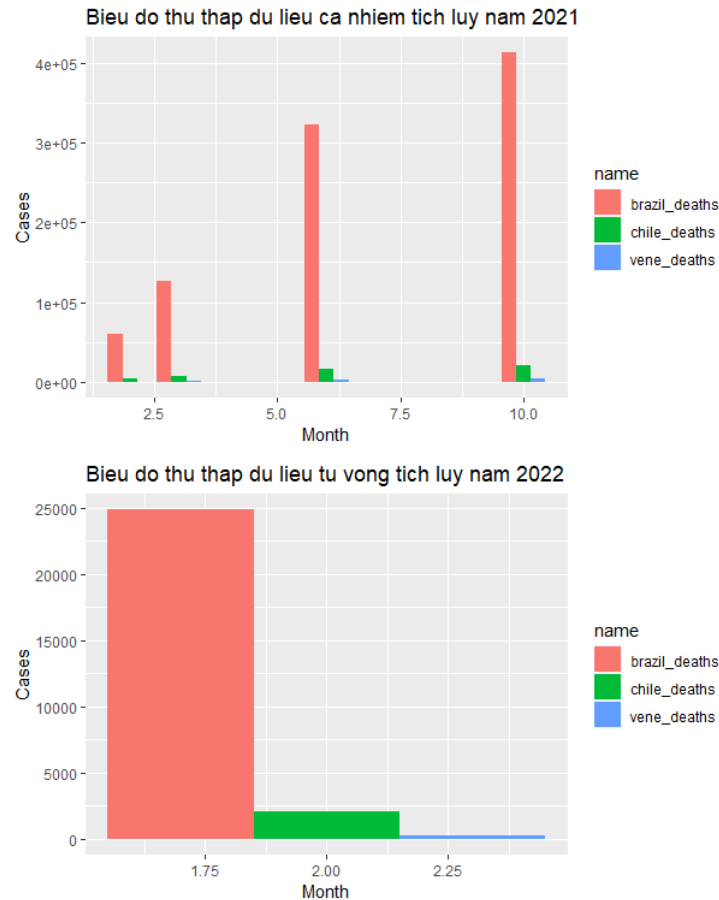
8) Biểu đồ thu thập tử vong tích lũy cho từng tháng

- Hiện thực trong R

```
1 temp<- cumsum_2020%>%
2   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
3 temp <- subset(temp,temp$month==2 | temp$month==10 |
4   temp$month==3 | temp$month==6)
5 ggplot(temp, aes(fill=name, y=value, x=month)) +
6   geom_bar(position='dodge', stat='identity') +
7   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu tu vong tích luy nam 2020')
8 temp<- cumsum_2021%>%
9   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
10 temp <- subset(temp,temp$month==2 | temp$month==10 |
11   temp$month==3 | temp$month==6)
12 ggplot(temp, aes(fill=name, y=value, x=month)) +
13   geom_bar(position='dodge', stat='identity') +
14   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2021')
15 temp<- cumsum_2022%>%
16   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,vene_deaths))
17 temp <- subset(temp,temp$month==2 | temp$month==10 |
18   temp$month==3 | temp$month==6)
19 ggplot(temp, aes(fill=name, y=value, x=month)) +
20   geom_bar(position='dodge', stat='identity') +
21   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu tu vong tích luy nam 2022')
```

- Kết quả:





vi) **Nhóm câu hỏi liên quan đến trực quan dữ liệu theo trung bình 7 ngày gần nhất:**

- Trên từng năm hãy vẽ biểu đồ thể hiện trục Ox là thời gian, trục Oy là nhiễm bệnh/tử vong. Hãy dùng 4 ký số của mã đề để vẽ 4 tháng tương ứng theo ký số đó. Nếu ký số là 0 thì lấy tháng là 10.
- Dùng trung bình của các ca nhiễm bệnh và tử vong được báo cáo trong 7 ngày gần nhất để loại trừ một số báo cáo không thường xuyên và đưa chúng ta đến gần hơn với con số hàng ngày.

- Chuẩn bị dữ liệu cho toàn bộ phần vi (phần code dùng chung cho cả phần vi)

```
1 library(tidyverse)
2 library(datasets)
3 library(dplyr)
4 library("ggplot2")
5 library(lubridate)
6 setwd("E:/BTL_CTRR")
7 covid = read.csv("owid-covid-data.csv", header = TRUE)
8 covid$date = as.Date(covid$date, format = "%m/%d/%Y")
9 covid = covid %>% filter(!continent == '') %>% mutate(new_cases = abs(new_cases), new_deaths = abs(
10   new_deaths))
11 covid[is.na(covid)] = 0
12 brazil = subset(covid, iso_code == "BRA")
13 chile = subset(covid, iso_code == "CHL")
14 venezuela = subset(covid, iso_code == "VEN")
15 brazil = brazil %>% mutate(day = as.numeric(format(as.Date(brazil$date), "%d"))) %>%
16   mutate(month = as.numeric(format(as.Date(brazil$date), "%m"))) %>%
17   mutate(year = as.numeric(format(as.Date(brazil$date), "%Y")))
18 brazil[["Newcases"]] = NA
19 brazil[["Newdeaths"]] = NA
20 for(i in 1:nrow(brazil)) {
21   if(i <= 7) {
22     brazil$Newcases[[i]] = mean(brazil[[5]][1:i])
23     brazil$Newdeaths[[i]] = mean(brazil[[6]][1:i])
24   }
```

```
25 }
26 else{
27   brazil$Newcases[[i]] = sum( brazil[[5]][(i-6):i] )/7
28   brazil$Newdeaths[[i]] = sum( brazil[[6]][(i-6):i] )/7
29 }
30 }
31 brazil = brazil %>% filter(!Newcases==0) %>% filter(!Newdeaths==0)
32
33 chile = chile %>% mutate(day =as.numeric(format(as.Date(chile$date),"%d"))) %>%
34   mutate(month = as.numeric(format(as.Date(chile$date),"%m"))) %>%
35   mutate(year = as.numeric(format(as.Date(chile$date),"%Y")))
36 # thay du lieu new cases, deaths bang gia tri trung binh 7 ngay gan nhat
37 chile[["Newcases"]]=NA
38 chile[["Newdeaths"]]=NA
39 for(i in 1:nrow(chile)) {
40
41   if(i <=7) {
42     chile$Newcases[[i]] = mean(chile[[5]][1:i])
43     chile$Newdeaths[[i]] = mean(chile[[6]][1:i])
44   }
45   else{
46     chile$Newcases[[i]] = sum( chile[[5]][(i-6):i] )/7
47     chile$Newdeaths[[i]] = sum( chile[[6]][(i-6):i] )/7
48   }
49 }
50 chile = chile %>% filter(!Newcases==0) %>% filter(!Newdeaths==0)
51
52 #Dung cot Newcases thay cho new_cases v Newdeaths thay cho new_deaths
53
54 venezuela = venezuela %>% mutate(day =as.numeric(format(as.Date(venezuela$date),"%d"))) %>%
55   mutate(month = as.numeric(format(as.Date(venezuela$date),"%m"))) %>%
56   mutate(year = as.numeric(format(as.Date(venezuela$date),"%Y")))
57 # thay du lieu new cases, deaths bang gia tri trung binh 7 ngay gan nhat
58 venezuela[["Newcases"]]=NA
59 venezuela[["Newdeaths"]]=NA
60 for(i in 1:nrow(venezuela)) {
61
62   if(i <=7) {
63     venezuela$Newcases[[i]] = mean(venezuela[[5]][1:i])
64     venezuela$Newdeaths[[i]] = mean(venezuela[[6]][1:i])
65   }
66   else{
67     venezuela$Newcases[[i]] = sum( venezuela[[5]][(i-6):i] )/7
68     venezuela$Newdeaths[[i]] = sum( venezuela[[6]][(i-6):i] )/7
69   }
70 }
71 venezuela = venezuela %>% filter(!Newcases==0) %>% filter(!Newdeaths==0)
72
73 #Dung cot Newcases thay cho new_cases v Newdeaths thay cho new_deaths
74
75 brazil_2020 <- subset(brazil,brazil$year == "2020")
76 brazil_2021 <- subset(brazil,brazil$year == "2021")
77 brazil_2022 <- subset(brazil,brazil$year == "2022")
78
79 chile_2020 <- subset(chile,chile$year == "2020")
80 chile_2021 <- subset(chile,chile$year == "2021")
81 chile_2022 <- subset(chile,chile$year == "2022")
82 vene_2020 <- subset(venezuela,venezuela$year == "2020")
83 vene_2021 <- subset(venezuela,venezuela$year == "2021")
84 vene_2022 <- subset(venezuela,venezuela$year == "2022")
85
86 temp=list(brazil_2020,chile_2020,vene_2020)
87 data_2020=temp %>% reduce(full_join, by='date')
88 data_2020 = rename(data_2020,brazil_cases=Newcases.x,brazil_deaths=Newdeaths.x,
89   chile_cases=Newcases.y,chile_deaths=Newdeaths.y,
90   venezuela_cases=Newcases,venezuela_deaths=Newdeaths)
91
92 temp=list(brazil_2021,chile_2021,vene_2021)
93 data_2021=temp %>% reduce(full_join, by='date')
94 data_2021 = rename(data_2021,brazil_cases=Newcases.x,brazil_deaths=Newdeaths.x,
95   chile_cases=Newcases.y,chile_deaths=Newdeaths.y,
```

```
96         venezuela_cases=Newcases,venezuela_deaths=Newdeaths)
97
98 temp=list(brazil_2022,chile_2022,vene_2022)
99 data_2022=temp %>% reduce(full_join, by='date')
100 data_2022 = rename(data_2022,brazil_cases=Newcases.x,brazil_deaths=Newdeaths.x,
101                    chile_cases=Newcases.y,chile_deaths=Newdeaths.y,
102                    venezuela_cases=Newcases,venezuela_deaths=Newdeaths)
```

1) Biểu đồ thu thập nhiễm bệnh cho từng tháng

- Hiện thực trong R

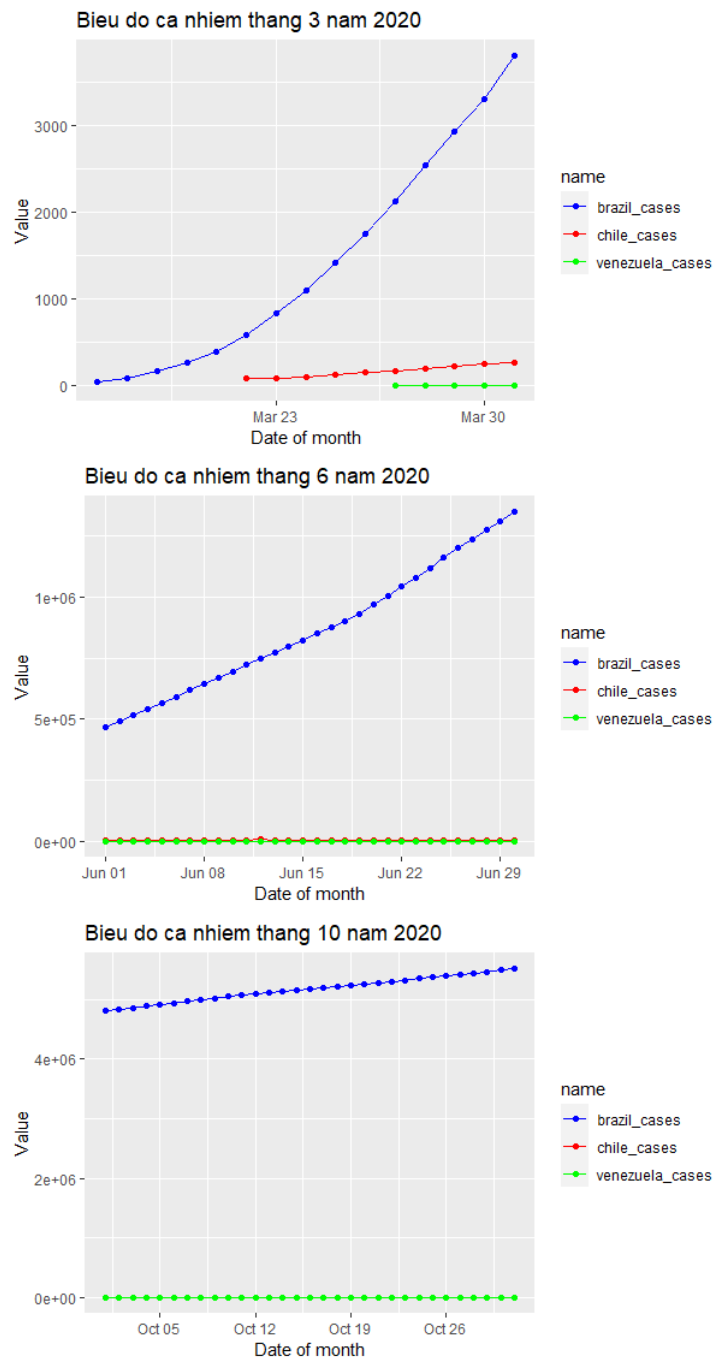
```
1 data_2020_3 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "3")
2 temp <- data_2020_3 %>%
3   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
4 ggplot(data=temp, aes(x=date, y=value, color=name)) +
5   geom_line() +
6   geom_point() +
7   scale_color_manual(values = c("blue", "red", "green")) +
8   labs(title="Biểu đồ ca nhiễm tháng 3 năm 2020", x="Date of month", y="Value")
9
10 data_2020_6 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "6")
11 temp <- data_2020_6 %>%
12   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
13 ggplot(data=temp, aes(x=date, y=value, color=name)) +
14   geom_line() +
15   geom_point() +
16   scale_color_manual(values = c("blue", "red", "green")) +
17   labs(title="Biểu đồ ca nhiễm tháng 6 năm 2020", x="Date of month", y="Value")
18
19 data_2020_10 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "10")
20 temp <- data_2020_10 %>%
21   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
22 ggplot(data=temp, aes(x=date, y=value, color=name)) +
23   geom_line() +
24   geom_point() +
25   scale_color_manual(values = c("blue", "red", "green")) +
26   labs(title="Biểu đồ ca nhiễm tháng 10 năm 2020", x="Date of month", y="Value")
27
28 data_2021_2 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "2")
29 temp <- data_2021_2 %>%
30   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
31 ggplot(data=temp, aes(x=date, y=value, color=name)) +
32   geom_line() +
33   geom_point() +
34   scale_color_manual(values = c("blue", "red", "green")) +
35   labs(title="Biểu đồ ca nhiễm tháng 2 năm 2021", x="Date of month", y="Value")
36
37 data_2021_3 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "3")
38 temp <- data_2021_3 %>%
39   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
40 ggplot(data=temp, aes(x=date, y=value, color=name)) +
41   geom_line() +
42   geom_point() +
43   scale_color_manual(values = c("blue", "red", "green")) +
44   labs(title="Biểu đồ ca nhiễm tháng 3 năm 2021", x="Date of month", y="Value")
45
46 data_2021_6 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "6")
47 temp <- data_2021_6 %>%
48   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
49 ggplot(data=temp, aes(x=date, y=value, color=name)) +
50   geom_line() +
51   geom_point() +
52   scale_color_manual(values = c("blue", "red", "green")) +
53   labs(title="Biểu đồ ca nhiễm tháng 6 năm 2021", x="Date of month", y="Value")
54
55 data_2021_10 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "10")
56 temp <- data_2021_10 %>%
57   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
58 ggplot(data=temp, aes(x=date, y=value, color=name)) +
59   geom_line() +
```

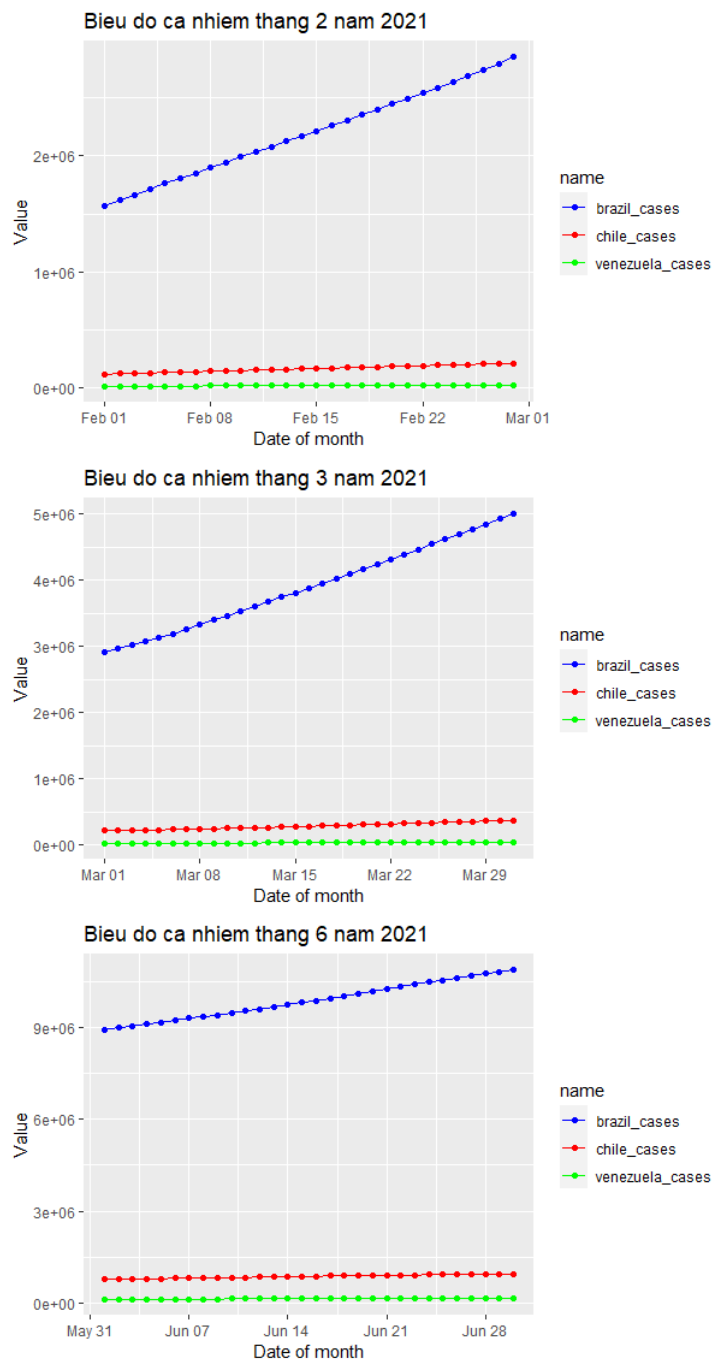
```

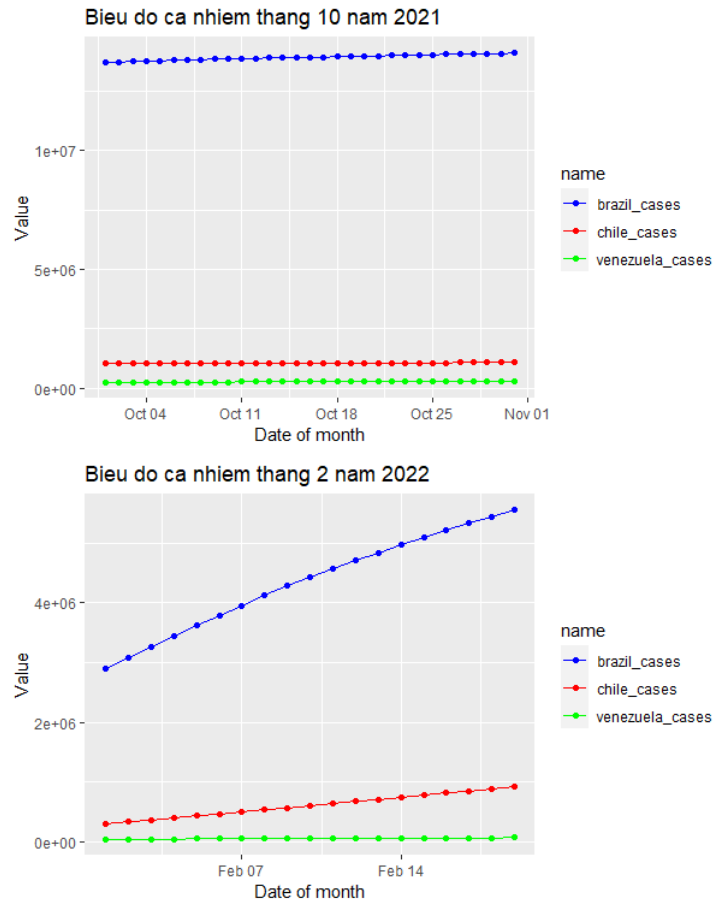
60 geom_point()+
61 scale_color_manual(values =c("blue","red","green"))+
62 labs(title="Bieu do ca nhiem thang 10 nam 2021",x="Date of month",y="Value")
63
64 data_2022_2 = subset(data_2022, as.numeric(format(as.Date(data_2022$date), "%m")) == "2")
65 temp <- data_2022_2>%
66 tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,venezuela_cases))
67 ggplot(data=temp,aes(x=date, y=value,color=name))+
68 geom_line()+
69 geom_point()+
70 scale_color_manual(values =c("blue","red","green"))+
71 labs(title="Bieu do ca nhiem thang 2 nam 2022",x="Date of month",y="Value")

```

• Kết quả:







2) Biểu đồ thu thập tử vong cho từng tháng

- Hiện thực trong R

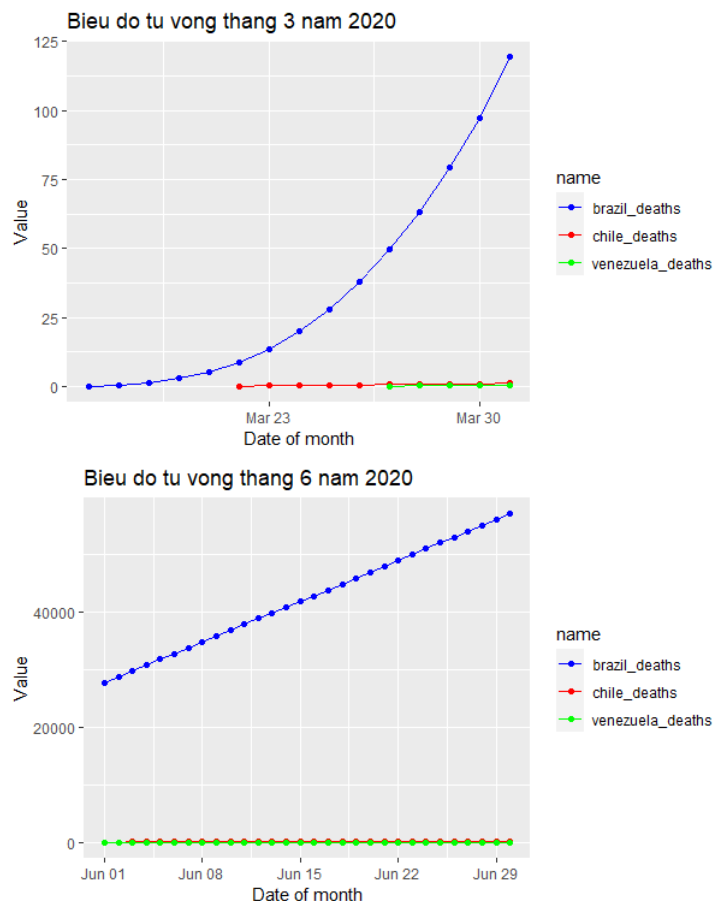
```
1 #Vi_2
2 temp <- data_2020_3%>%
3   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
4   ggplot(data=temp,aes(x=date, y=value,color=name))+
5     geom_line()+
6     geom_point()+
7     scale_color_manual(values =c("blue","red","green"))+
8     labs(title="Bieu do tu vong thang 3 nam 2020",x="Date of month",y="Value")
9
10 temp <- data_2020_6%>%
11   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
12   ggplot(data=temp,aes(x=date, y=value,color=name))+
13     geom_line()+
14     geom_point()+
15     scale_color_manual(values =c("blue","red","green"))+
16     labs(title="Bieu do tu vong thang 6 nam 2020",x="Date of month",y="Value")
17
18 temp <- data_2020_10%>%
19   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
20   ggplot(data=temp,aes(x=date, y=value,color=name))+
21     geom_line()+
22     geom_point()+
23     scale_color_manual(values =c("blue","red","green"))+
24     labs(title="Bieu do tu vong thang 10 nam 2020",x="Date of month",y="Value")
25
26 temp <- data_2021_2%>%
27   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
28   ggplot(data=temp,aes(x=date, y=value,color=name))+
29     geom_line()+
30     geom_point()+
31     scale_color_manual(values =c("blue","red","green"))+
```

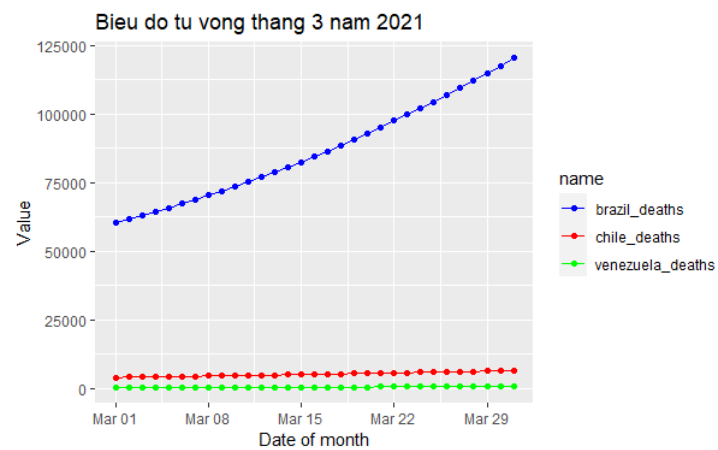
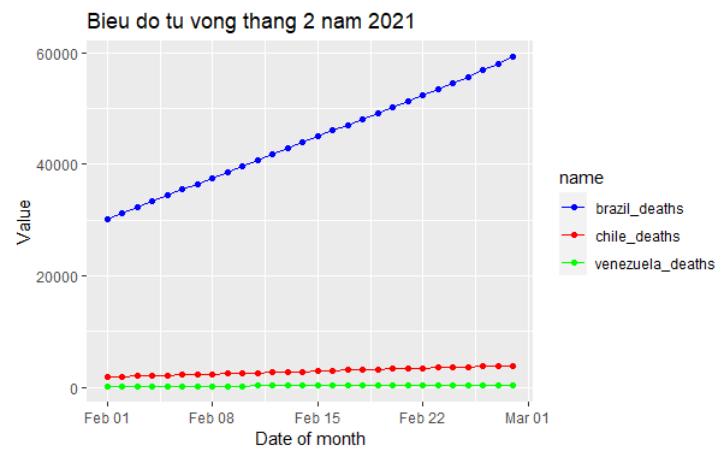
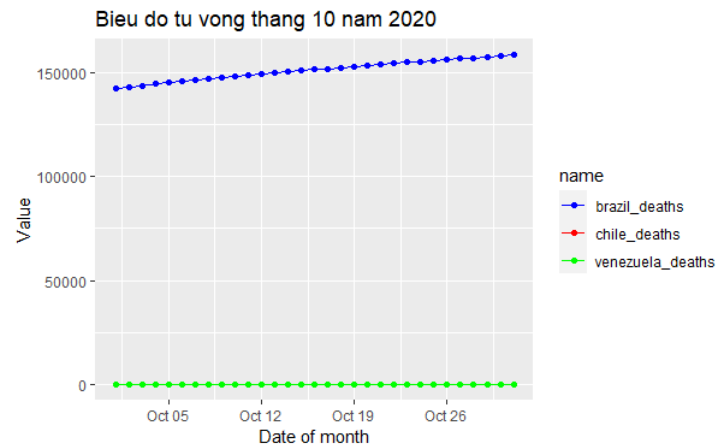
```

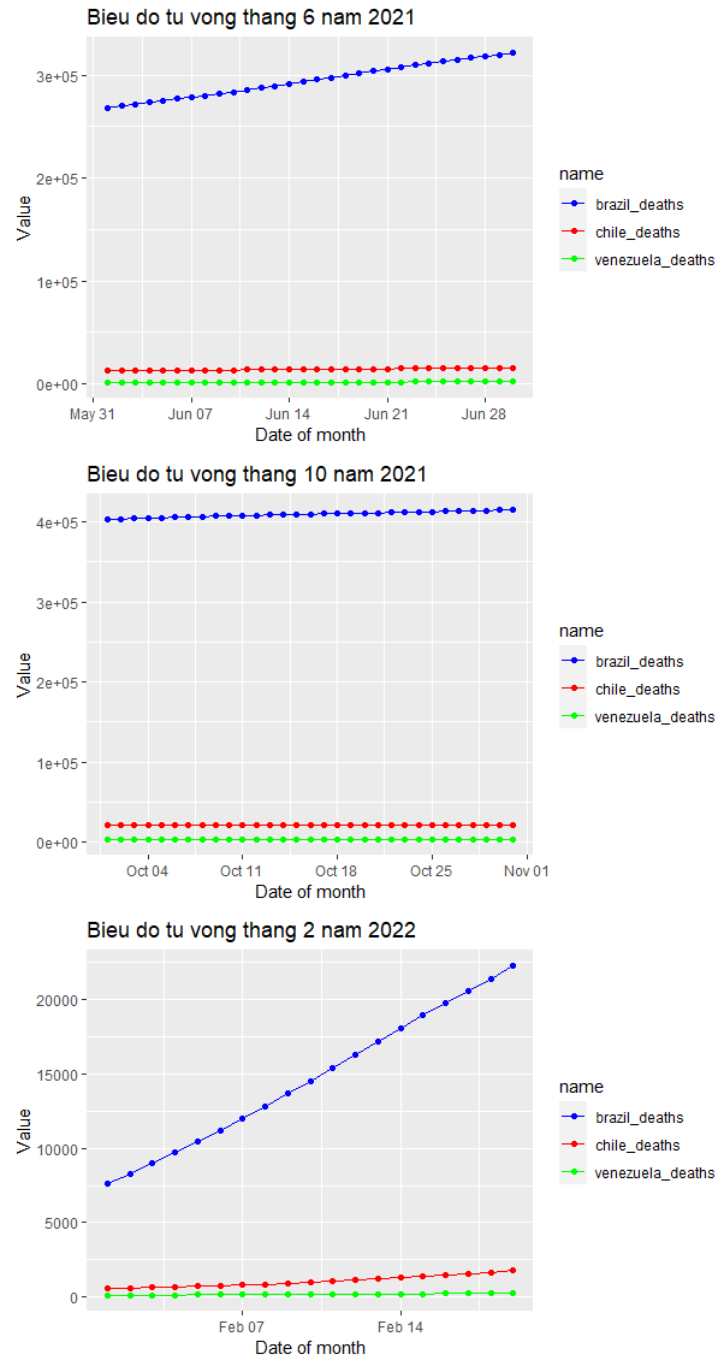
32 labs(title="Bieu do tu vong thang 2 nam 2021",x="Date of month",y="Value")
33
34 temp <- data_2021_3%>%
35   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
36   ggplot(data=temp,aes(x=date, y=value,color=name))+
37     geom_line()+
38     geom_point()+
39     scale_color_manual(values =c("blue","red","green"))+
40     labs(title="Bieu do tu vong thang 3 nam 2021",x="Date of month",y="Value")
41
42 temp <- data_2021_6%>%
43   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
44   ggplot(data=temp,aes(x=date, y=value,color=name))+
45     geom_line()+
46     geom_point()+
47     scale_color_manual(values =c("blue","red","green"))+
48     labs(title="Bieu do tu vong thang 6 nam 2021",x="Date of month",y="Value")
49
50 temp <- data_2021_10%>%
51   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
52   ggplot(data=temp,aes(x=date, y=value,color=name))+
53     geom_line()+
54     geom_point()+
55     scale_color_manual(values =c("blue","red","green"))+
56     labs(title="Bieu do tu vong thang 10 nam 2021",x="Date of month",y="Value")
57
58 temp <- data_2022_2%>%
59   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
60   ggplot(data=temp,aes(x=date, y=value,color=name))+
61     geom_line()+
62     geom_point()+
63     scale_color_manual(values =c("blue","red","green"))+
64     labs(title="Bieu do tu vong thang 2 nam 2022",x="Date of month",y="Value")

```

● Kết quả:







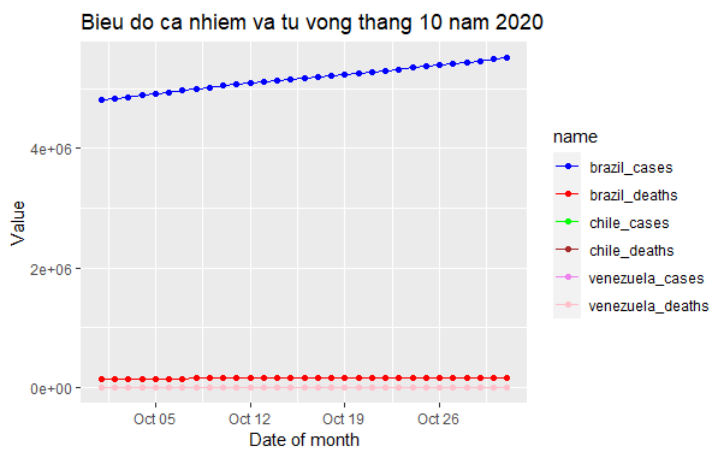
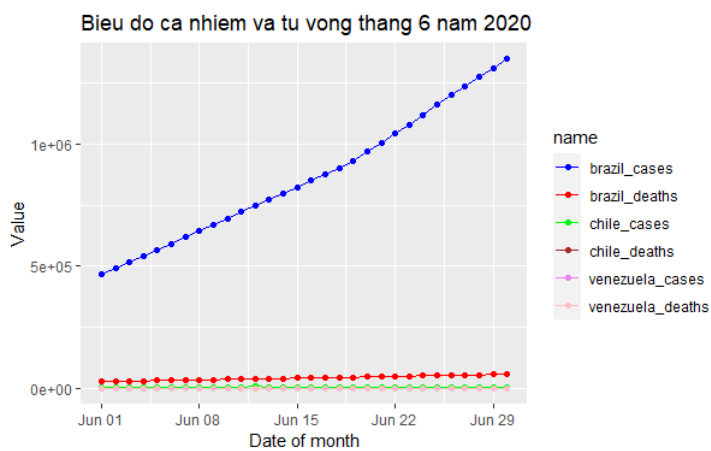
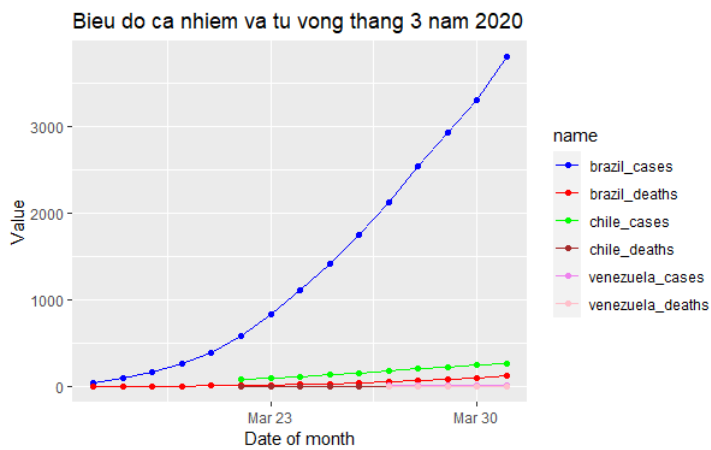
3) Biểu đồ thu thập gồm nhiễm bệnh và tử vong cho từng tháng

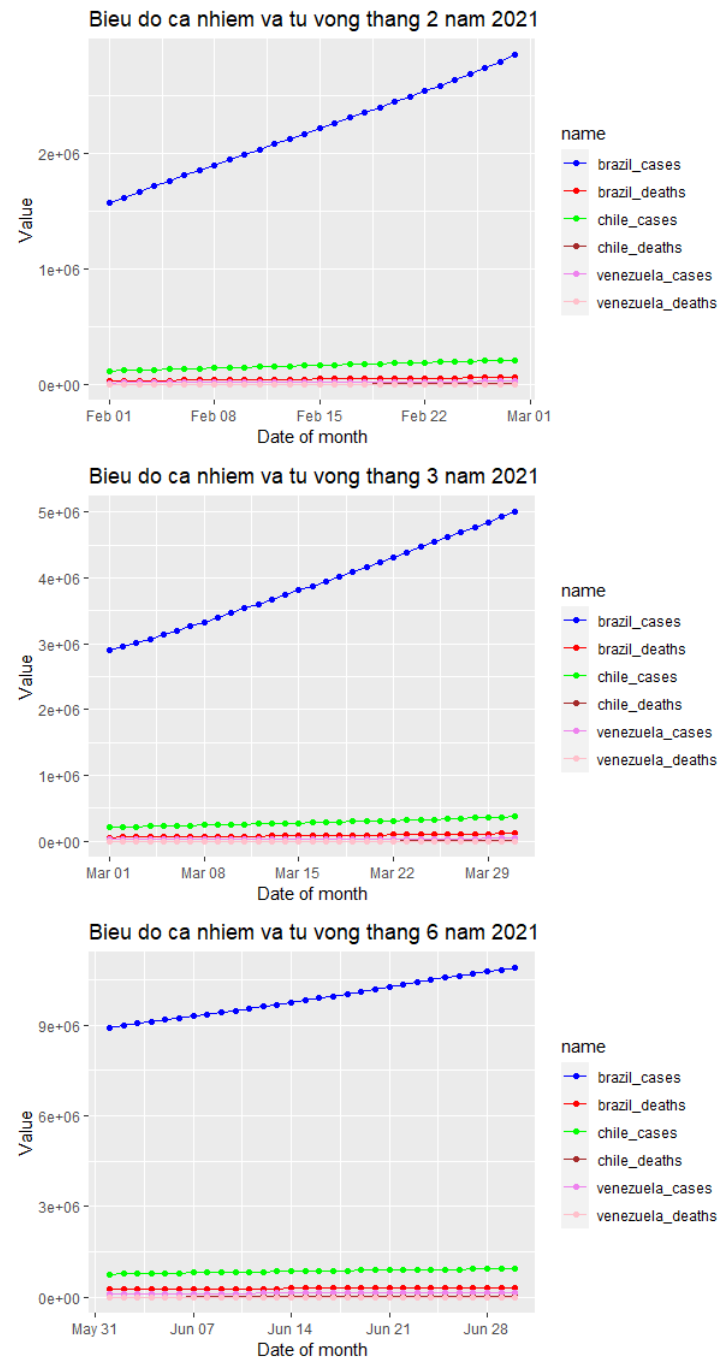
- Hiện thực trong R

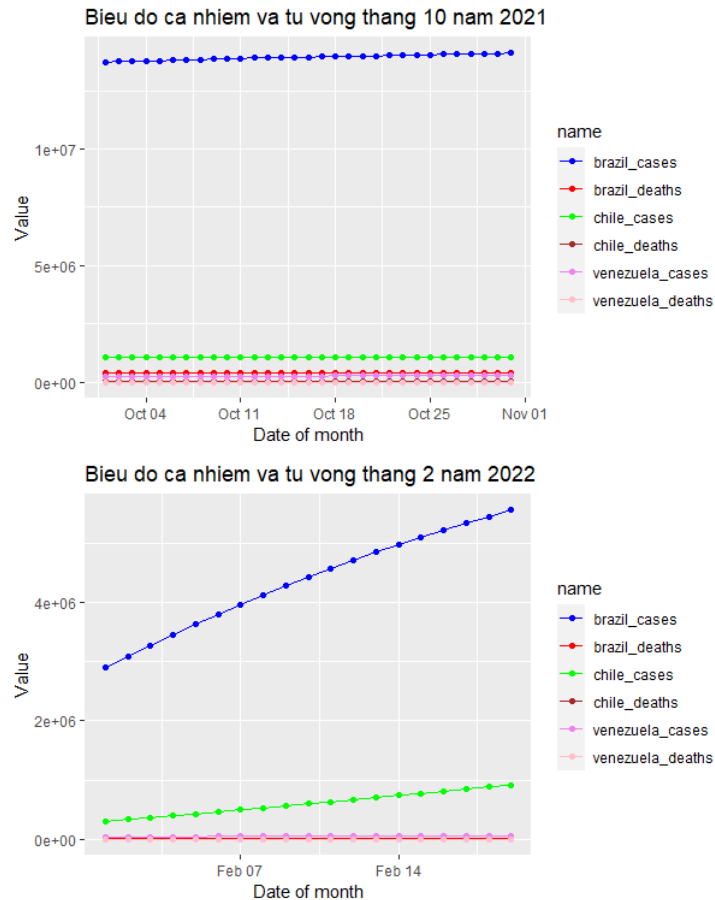
```
1 temp <- data_2020_3%>%
2   tidyr::pivot_longer(cols = c(
3     brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
4 ggplot(data=temp,aes(x=date, y=value,color=name))+
5   geom_line()+
6   geom_point()+
7   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
8   labs(title="Bieu do ca nhien va tu vong thang 3 nam 2020",x="Date of month",y="Value")
9 temp <- data_2020_6%>%
10  tidyr::pivot_longer(cols = c(
11    brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
12 ggplot(data=temp,aes(x=date, y=value,color=name))+
```

```
12 geom_line()+
13 geom_point()+
14 scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
15 labs(title="Bieu do ca nhien va tu vong thang 6 nam 2020",x="Date of month",y="Value")
16
17 temp <- data_2020_10%>%
18 tidyr::pivot_longer(cols = c(
19   brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
20 ggplot(data=temp,aes(x=date, y=value,color=name))+
21   geom_line()+
22   geom_point()+
23   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
24   labs(title="Bieu do ca nhien va tu vong thang 10 nam 2020",x="Date of month",y="Value")
25
26 temp <- data_2021_2%>%
27 tidyr::pivot_longer(cols = c(
28   brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
29 ggplot(data=temp,aes(x=date, y=value,color=name))+
30   geom_line()+
31   geom_point()+
32   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
33   labs(title="Bieu do ca nhien va tu vong thang 2 nam 2021",x="Date of month",y="Value")
34
35 temp <- data_2021_3%>%
36 tidyr::pivot_longer(cols = c(
37   brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
38 ggplot(data=temp,aes(x=date, y=value,color=name))+
39   geom_line()+
40   geom_point()+
41   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
42   labs(title="Bieu do ca nhien va tu vong thang 3 nam 2021",x="Date of month",y="Value")
43
44 temp <- data_2021_6%>%
45 tidyr::pivot_longer(cols = c(
46   brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
47 ggplot(data=temp,aes(x=date, y=value,color=name))+
48   geom_line()+
49   geom_point()+
50   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
51   labs(title="Bieu do ca nhien va tu vong thang 6 nam 2021",x="Date of month",y="Value")
52
53 temp <- data_2021_10%>%
54 tidyr::pivot_longer(cols = c(
55   brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
56 ggplot(data=temp,aes(x=date, y=value,color=name))+
57   geom_line()+
58   geom_point()+
59   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
60   labs(title="Bieu do ca nhien va tu vong thang 10 nam 2021",x="Date of month",y="Value")
61
62 temp <- data_2022_2%>%
63 tidyr::pivot_longer(cols = c(
64   brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
65 ggplot(data=temp,aes(x=date, y=value,color=name))+
66   geom_line()+
67   geom_point()+
68   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
69   labs(title="Bieu do ca nhien va tu vong thang 2 nam 2022",x="Date of month",y="Value")
```

- Kết quả:







4) Biểu đồ thu thập nhiễm bệnh gồm 2 tháng cuối của năm

• Hiện thực trong R

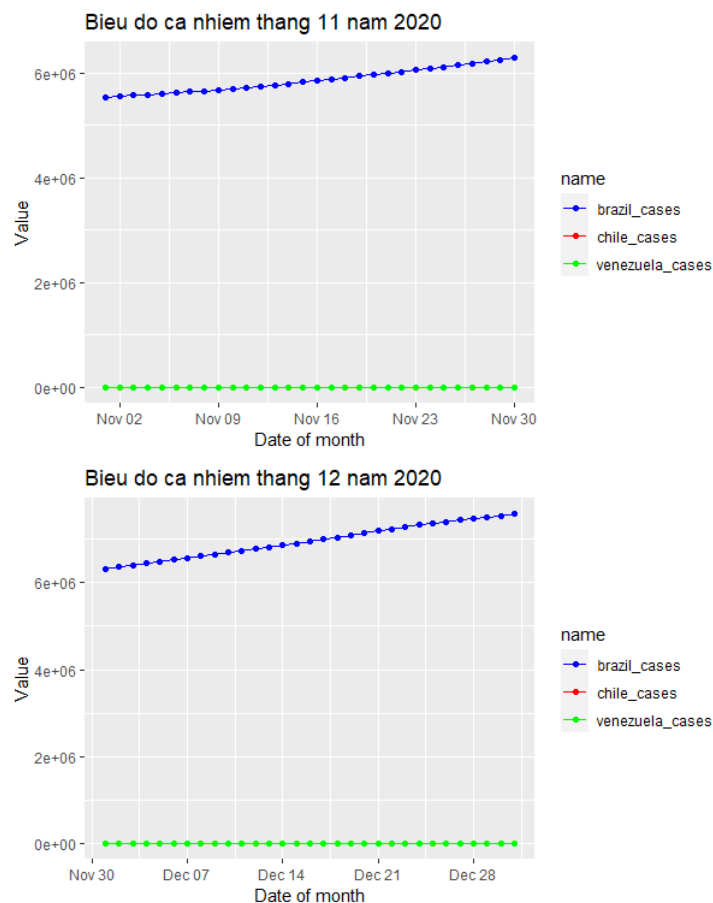
```
1 data_2020_11 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "11")
2 data_2020_12 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "12")
3
4 data_2021_11 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "11")
5 data_2021_12 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "12")
6
7 data_2022_1 = subset(data_2022, as.numeric(format(as.Date(data_2022$date), "%m")) == "1")
8 data_2022_2 = subset(data_2022, as.numeric(format(as.Date(data_2022$date), "%m")) == "2")
9
10 temp <- data_2020_11 %>%
11   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
12 ggplot(data=temp, aes(x=date, y=value, color=name)) +
13   geom_line() +
14   geom_point() +
15   scale_color_manual(values = c("blue", "red", "green")) +
16   labs(title="Bieu do ca nhien thang 11 nam 2020", x="Date of month", y="Value")
17
18 temp <- data_2020_12 %>%
19   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
20 ggplot(data=temp, aes(x=date, y=value, color=name)) +
21   geom_line() +
22   geom_point() +
23   scale_color_manual(values = c("blue", "red", "green")) +
24   labs(title="Bieu do ca nhien thang 12 nam 2020", x="Date of month", y="Value")
25
26 temp <- data_2021_11 %>%
27   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
28 ggplot(data=temp, aes(x=date, y=value, color=name)) +
29   geom_line() +
30   geom_point() +
31   scale_color_manual(values = c("blue", "red", "green")) +
```

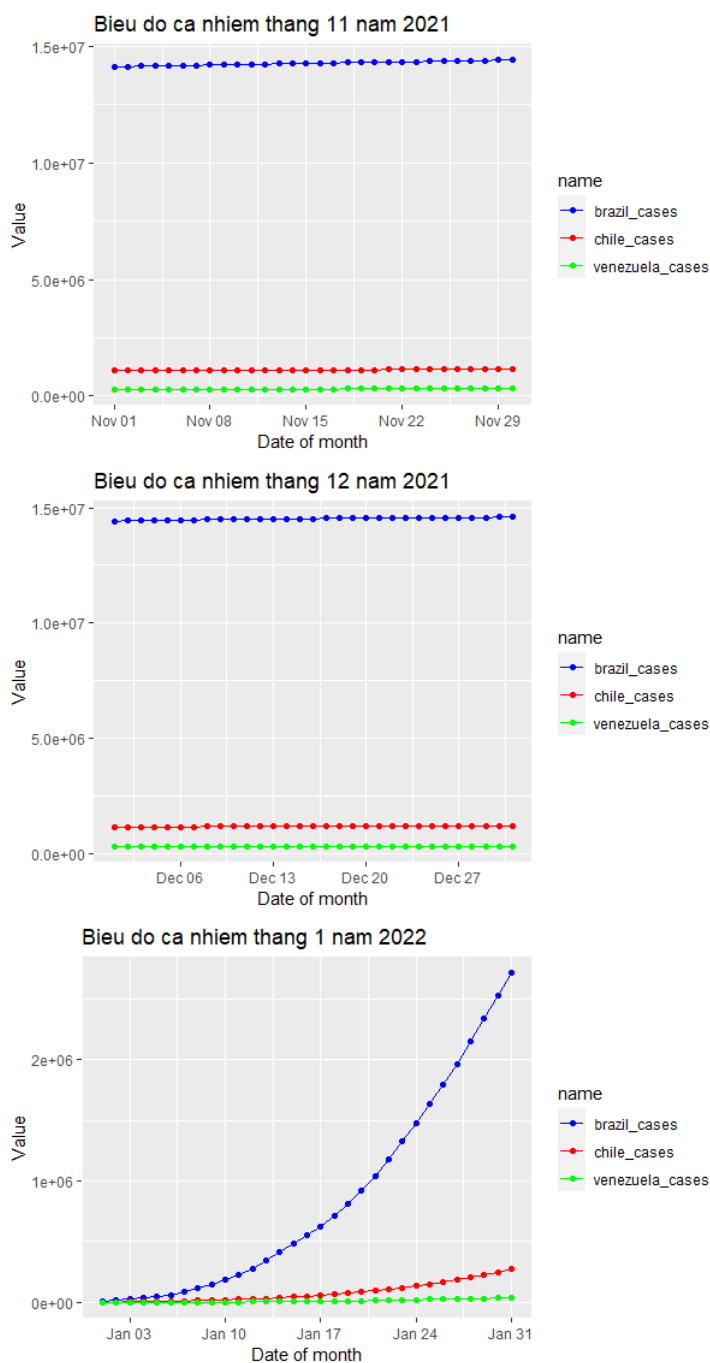
```

32 labs(title="Bieu do ca nhiem thang 11 nam 2021",x="Date of month",y="Value")
33
34 temp <- data_2021_12%>%
35   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,venezuela_cases))
36   ggplot(data=temp,aes(x=date, y=value,color=name))+
37     geom_line()+
38     geom_point()+
39     scale_color_manual(values =c("blue","red","green"))+
40     labs(title="Bieu do ca nhiem thang 12 nam 2021",x="Date of month",y="Value")
41
42 temp <- data_2022_1%>%
43   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,venezuela_cases))
44   ggplot(data=temp,aes(x=date, y=value,color=name))+
45     geom_line()+
46     geom_point()+
47     scale_color_manual(values =c("blue","red","green"))+
48     labs(title="Bieu do ca nhiem thang 1 nam 2022",x="Date of month",y="Value")
49
50 temp <- data_2022_2%>%
51   tidyr::pivot_longer(cols = c(brazil_cases,chile_cases,venezuela_cases))
52   ggplot(data=temp,aes(x=date, y=value,color=name))+
53     geom_line()+
54     geom_point()+
55     scale_color_manual(values =c("blue","red","green"))+
56     labs(title="Bieu do ca nhiem thang 2 nam 2022",x="Date of month",y="Value")

```

- Kết quả:





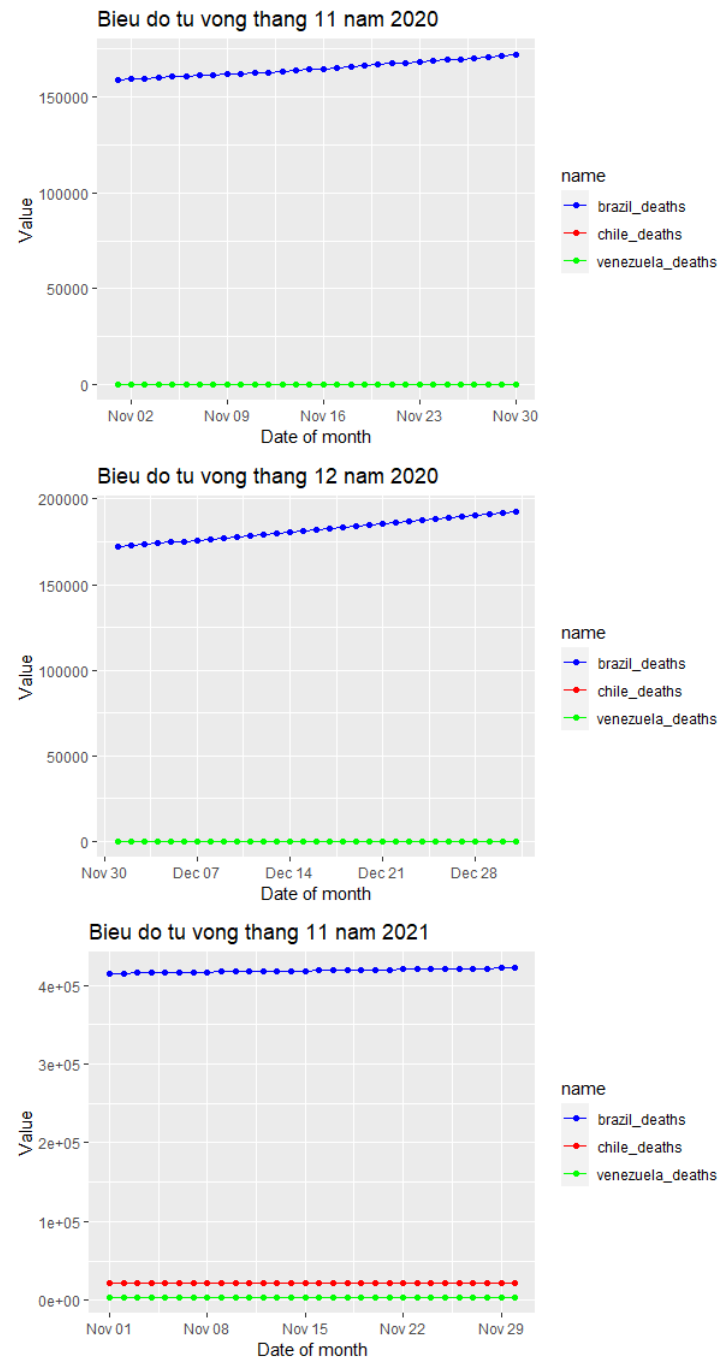


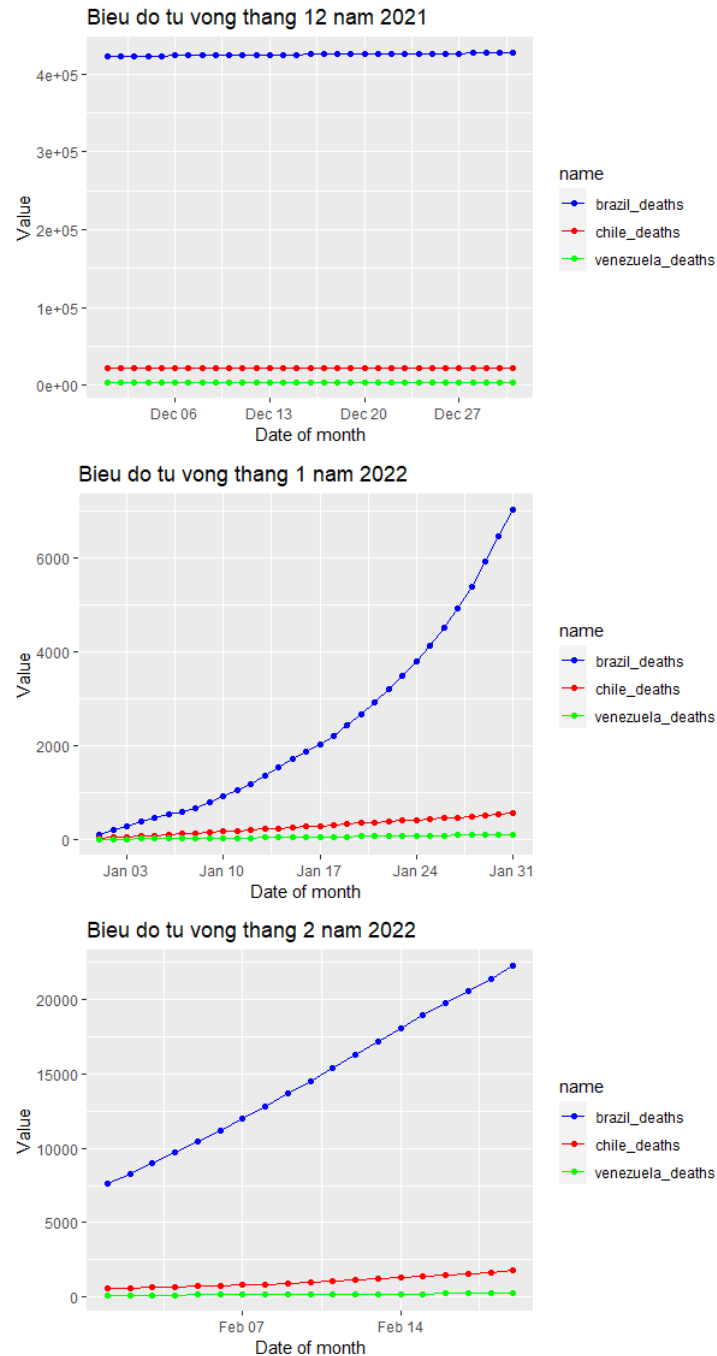
5) Biểu đồ thu thập tử vong gồm 2 tháng cuối của năm

- Hiện thực trong R

```
1 temp <- data_2020_11%>%
2   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
3   ggplot(data=temp,aes(x=date, y=value,color=name))+
4     geom_line()+
5     geom_point()+
6     scale_color_manual(values =c("blue","red","green"))+
7     labs(title="Bieu do tu vong thang 11 nam 2020",x="Date of month",y="Value")
8
9 temp <- data_2020_12%>%
10   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
11   ggplot(data=temp,aes(x=date, y=value,color=name))+
12     geom_line()+
13     geom_point()+
14     scale_color_manual(values =c("blue","red","green"))+
15     labs(title="Bieu do tu vong thang 12 nam 2020",x="Date of month",y="Value")
16
17 temp <- data_2021_11%>%
18   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
19   ggplot(data=temp,aes(x=date, y=value,color=name))+
20     geom_line()+
21     geom_point()+
22     scale_color_manual(values =c("blue","red","green"))+
23     labs(title="Bieu do tu vong thang 11 nam 2021",x="Date of month",y="Value")
24
25 temp <- data_2021_12%>%
26   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
27   ggplot(data=temp,aes(x=date, y=value,color=name))+
28     geom_line()+
29     geom_point()+
30     scale_color_manual(values =c("blue","red","green"))+
31     labs(title="Bieu do tu vong thang 12 nam 2021",x="Date of month",y="Value")
32
33 temp <- data_2022_1%>%
34   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
35   ggplot(data=temp,aes(x=date, y=value,color=name))+
36     geom_line()+
37     geom_point()+
38     scale_color_manual(values =c("blue","red","green"))+
39     labs(title="Bieu do tu vong thang 1 nam 2022",x="Date of month",y="Value")
40
41 temp <- data_2022_2%>%
42   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
43   ggplot(data=temp,aes(x=date, y=value,color=name))+
44     geom_line()+
45     geom_point()+
46     scale_color_manual(values =c("blue","red","green"))+
47     labs(title="Bieu do tu vong thang 2 nam 2022",x="Date of month",y="Value")
```

- Kết quả:





6) Biểu đồ thu thập gồm nhiễm bệnh và tử vong gồm 2 tháng cuối của năm

- Hiện thực trong R

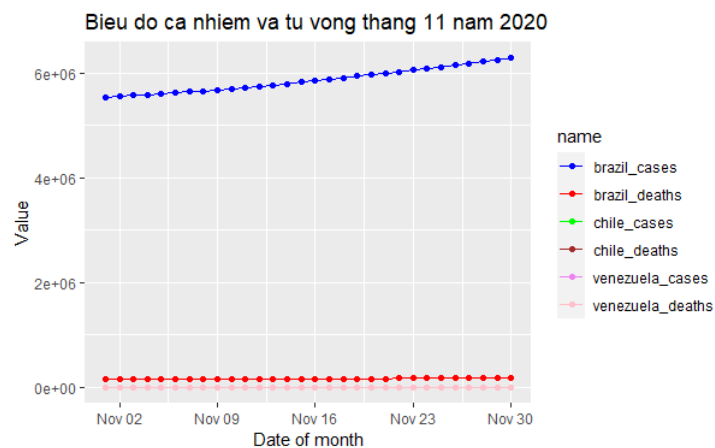
```
1 temp <- data_2020_11%>%
2   tidyr::pivot_longer(cols = c(
3     brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
4 ggplot(data=temp,aes(x=date, y=value,color=name))+
5   geom_line()+
6   geom_point()+
7   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
8   labs(title="Bieu do ca nhien va tu vong thang 11 nam 2020",x="Date of month",y="Value")
9 temp <- data_2020_12%>%
10  tidyr::pivot_longer(cols = c(
11    brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
12 ggplot(data=temp,aes(x=date, y=value,color=name))+
```

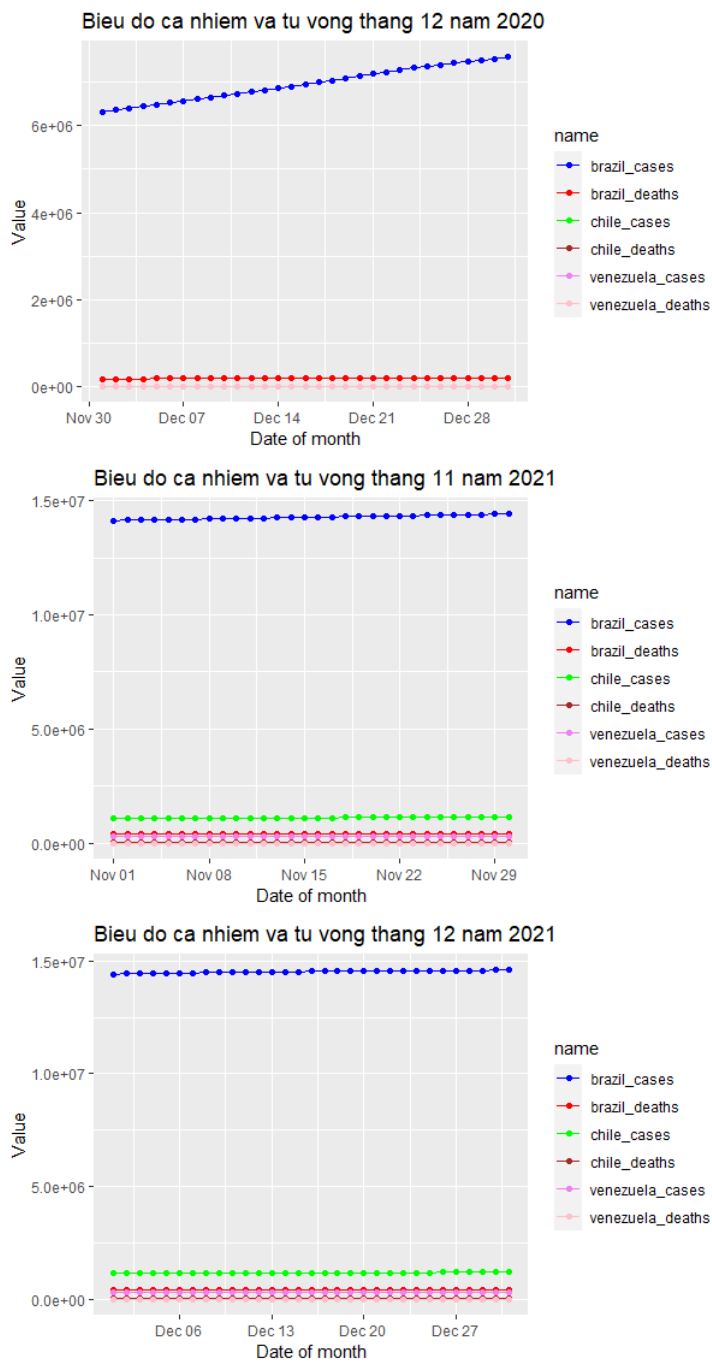
```

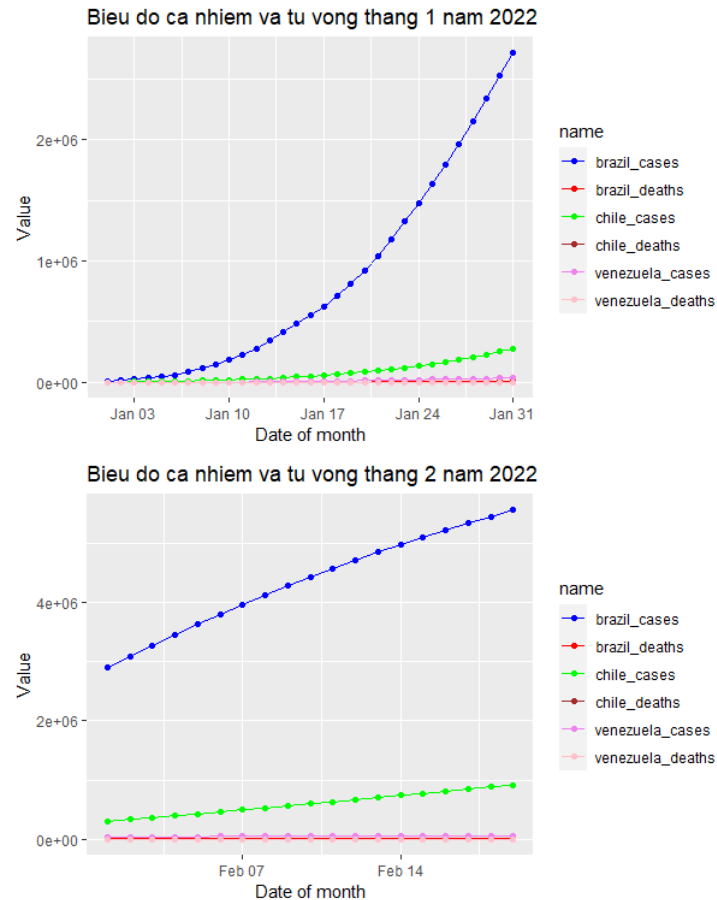
12 geom_line()+
13 geom_point()+
14 scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
15 labs(title="Bieu do ca nhien va tu vong thang 12 nam 2020",x="Date of month",y="Value")
16
17 temp <- data_2021_11%>%
18   tidyr::pivot_longer(cols = c(
19     brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
19 ggplot(data=temp,aes(x=date, y=value,color=name))+
20   geom_line()+
21   geom_point()+
22   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
23   labs(title="Bieu do ca nhien va tu vong thang 11 nam 2021",x="Date of month",y="Value")
24
25 temp <- data_2021_12%>%
26   tidyr::pivot_longer(cols = c(
27     brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
27 ggplot(data=temp,aes(x=date, y=value,color=name))+
28   geom_line()+
29   geom_point()+
30   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
31   labs(title="Bieu do ca nhien va tu vong thang 12 nam 2021",x="Date of month",y="Value")
32
33 temp <- data_2022_1%>%
34   tidyr::pivot_longer(cols = c(
35     brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
35 ggplot(data=temp,aes(x=date, y=value,color=name))+
36   geom_line()+
37   geom_point()+
38   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
39   labs(title="Bieu do ca nhien va tu vong thang 1 nam 2022",x="Date of month",y="Value")
40
41 temp <- data_2022_2%>%
42   tidyr::pivot_longer(cols = c(
43     brazil_cases,brazil_deaths,chile_cases,chile_deaths,venezuela_cases,venezuela_deaths))
43 ggplot(data=temp,aes(x=date, y=value,color=name))+
44   geom_line()+
45   geom_point()+
46   scale_color_manual(values =c("blue","red","green","brown","violet","pink"))+
47   labs(title="Bieu do ca nhien va tu vong thang 2 nam 2022",x="Date of month",y="Value")

```

• Kết quả:







7) Biểu đồ thu thập nhiễm bệnh tích lũy cho từng tháng

- Hiện thực trong R

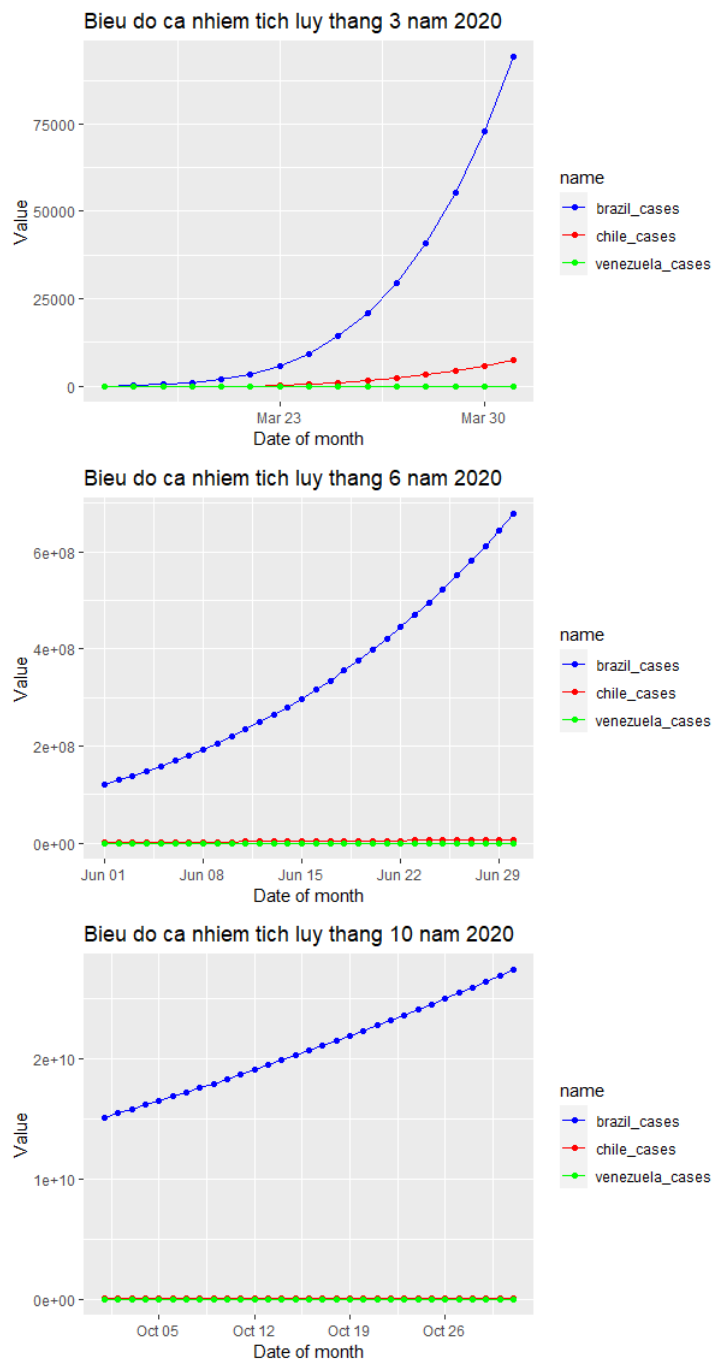
```

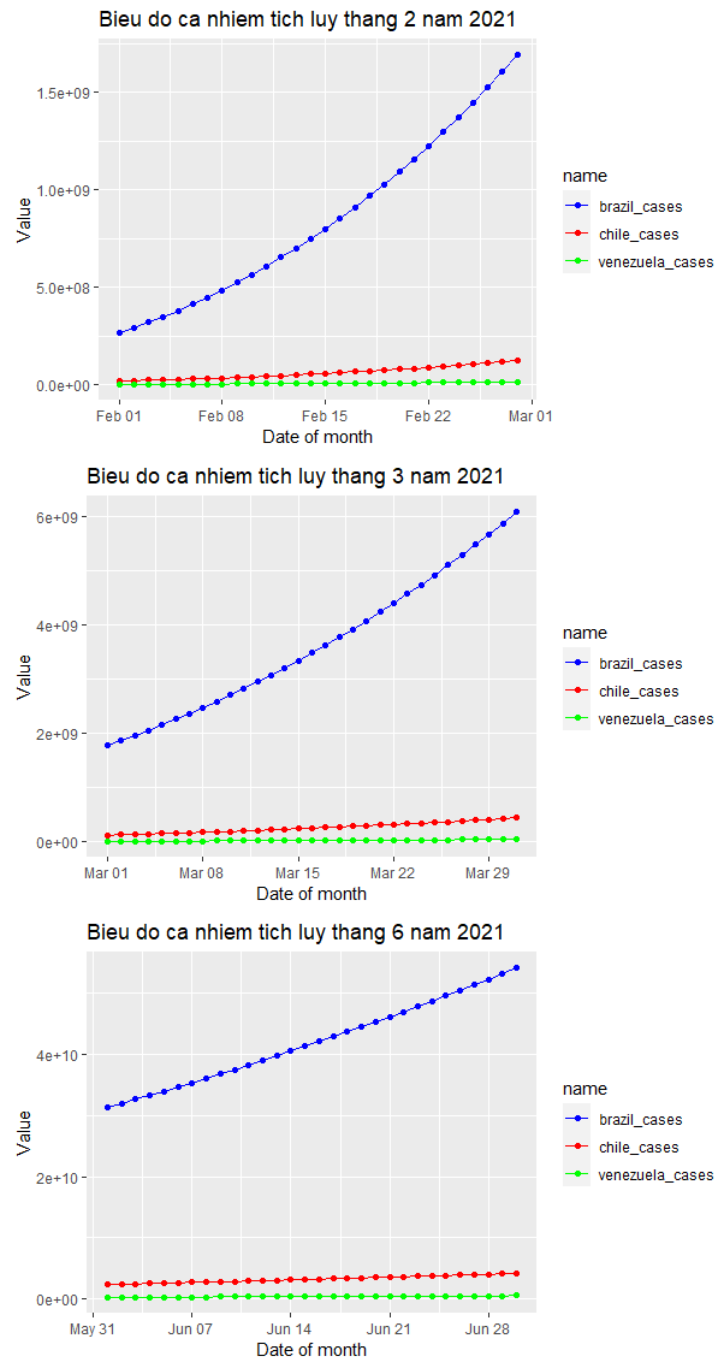
1 data_2020[is.na(data_2020)] = 0
2 data_2020$brazil_cases = cumsum(data_2020$brazil_cases)
3 data_2020$brazil_deaths = cumsum(data_2020$brazil_deaths)
4 data_2021$brazil_cases = cumsum(data_2021$brazil_cases)
5 data_2021$brazil_deaths = cumsum(data_2021$brazil_deaths)
6 data_2022$brazil_cases = cumsum(data_2022$brazil_cases)
7 data_2022$brazil_deaths = cumsum(data_2022$brazil_deaths)
8
9 data_2020$chile_cases = cumsum(data_2020$chile_cases)
10 data_2020$chile_deaths = cumsum(data_2020$chile_deaths)
11 data_2021$chile_cases = cumsum(data_2021$chile_cases)
12 data_2021$chile_deaths = cumsum(data_2021$chile_deaths)
13 data_2022$chile_cases = cumsum(data_2022$chile_cases)
14 data_2022$chile_deaths = cumsum(data_2022$chile_deaths)
15
16 data_2020$venezuela_cases = cumsum(data_2020$venezuela_cases)
17 data_2020$venezuela_deaths = cumsum(data_2020$venezuela_deaths)
18 data_2021$venezuela_cases = cumsum(data_2021$venezuela_cases)
19 data_2021$venezuela_deaths = cumsum(data_2021$venezuela_deaths)
20 data_2022$venezuela_cases = cumsum(data_2022$venezuela_cases)
21 data_2022$venezuela_deaths = cumsum(data_2022$venezuela_deaths)
22
23 cumsum_2020_3 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "3")
24 cumsum_2020_6 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "6")
25 cumsum_2020_10 = subset(data_2020, as.numeric(format(as.Date(data_2020$date), "%m")) == "10")
26
27 cumsum_2021_2 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "2")
28 cumsum_2021_3 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "3")
29 cumsum_2021_6 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "6")
30 cumsum_2021_10 = subset(data_2021, as.numeric(format(as.Date(data_2021$date), "%m")) == "10")
31

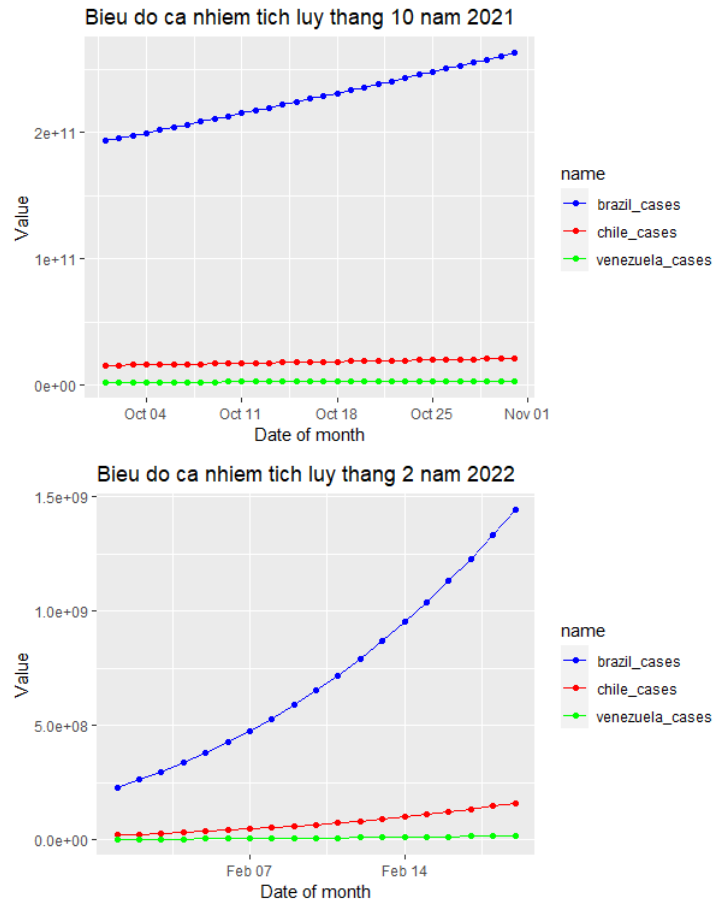
```

```
32 cumsum_2022_2 = subset(data_2022, as.numeric(format(as.Date(data_2022$date), "%m")) == "2")
33
34 temp <- cumsum_2020_3%>%
35   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
36   ggplot(data=temp, aes(x=date, y=value, color=name))+
37     geom_line()+
38     geom_point()+
39     scale_color_manual(values = c("blue", "red", "green"))+
40     labs(title="Bieu do ca nhien tích lũy thang 3 nam 2020", x="Date of month", y="Value")
41
42 temp <- cumsum_2020_6%>%
43   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
44   ggplot(data=temp, aes(x=date, y=value, color=name))+
45     geom_line()+
46     geom_point()+
47     scale_color_manual(values = c("blue", "red", "green"))+
48     labs(title="Bieu do ca nhien tích lũy thang 6 nam 2020", x="Date of month", y="Value")
49
50 temp <- cumsum_2020_10%>%
51   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
52   ggplot(data=temp, aes(x=date, y=value, color=name))+
53     geom_line()+
54     geom_point()+
55     scale_color_manual(values = c("blue", "red", "green"))+
56     labs(title="Bieu do ca nhien tích lũy thang 10 nam 2020", x="Date of month", y="Value")
57
58 temp <- cumsum_2021_2%>%
59   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
60   ggplot(data=temp, aes(x=date, y=value, color=name))+
61     geom_line()+
62     geom_point()+
63     scale_color_manual(values = c("blue", "red", "green"))+
64     labs(title="Bieu do ca nhien tích lũy thang 2 nam 2021", x="Date of month", y="Value")
65
66 temp <- cumsum_2021_3%>%
67   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
68   ggplot(data=temp, aes(x=date, y=value, color=name))+
69     geom_line()+
70     geom_point()+
71     scale_color_manual(values = c("blue", "red", "green"))+
72     labs(title="Bieu do ca nhien tích lũy thang 3 nam 2021", x="Date of month", y="Value")
73
74 temp <- cumsum_2021_6%>%
75   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
76   ggplot(data=temp, aes(x=date, y=value, color=name))+
77     geom_line()+
78     geom_point()+
79     scale_color_manual(values = c("blue", "red", "green"))+
80     labs(title="Bieu do ca nhien tích lũy thang 6 nam 2021", x="Date of month", y="Value")
81
82 temp <- cumsum_2021_10%>%
83   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
84   ggplot(data=temp, aes(x=date, y=value, color=name))+
85     geom_line()+
86     geom_point()+
87     scale_color_manual(values = c("blue", "red", "green"))+
88     labs(title="Bieu do ca nhien tích lũy thang 10 nam 2021", x="Date of month", y="Value")
89
90 temp <- cumsum_2022_2%>%
91   tidyr::pivot_longer(cols = c(brazil_cases, chile_cases, venezuela_cases))
92   ggplot(data=temp, aes(x=date, y=value, color=name))+
93     geom_line()+
94     geom_point()+
95     scale_color_manual(values = c("blue", "red", "green"))+
96     labs(title="Bieu do ca nhien tích lũy thang 2 nam 2022", x="Date of month", y="Value")
```

- Kết quả:







8) Biểu đồ thu thập tử vong tích lũy cho từng tháng

- Hiện thực trong R

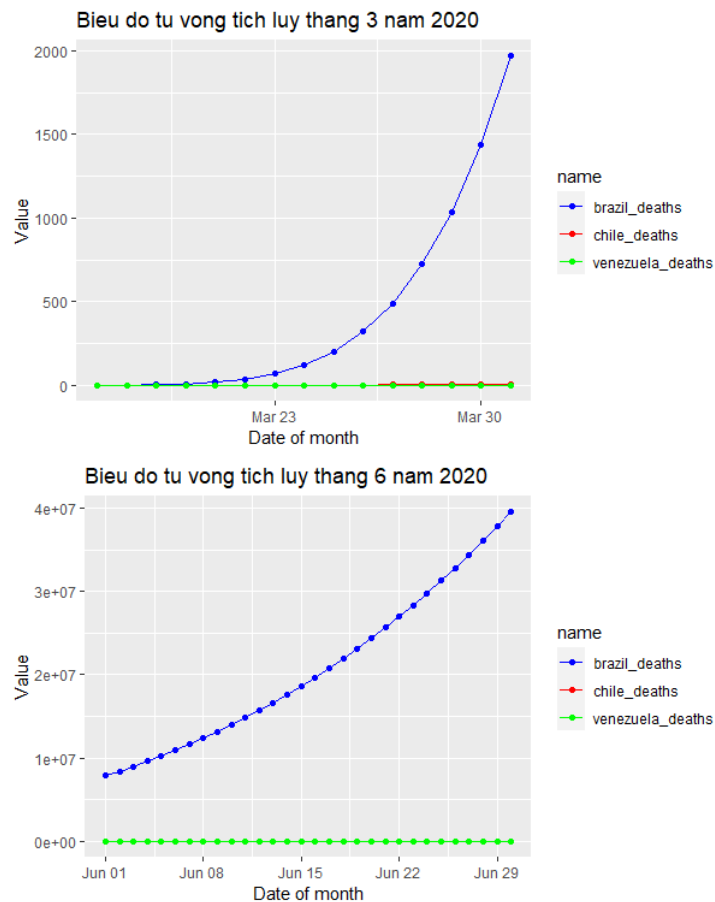
```
1 temp <- cumsum_2020_3%>%
2   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
3   ggplot(data=temp,aes(x=date, y=value,color=name))+
4     geom_line()+
5     geom_point()+
6     scale_color_manual(values =c("blue","red","green"))+
7     labs(title="Bieu do tu vong tích lũy tháng 3 năm 2020",x="Date of month",y="Value")
8
9 temp <- cumsum_2020_6%>%
10   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
11   ggplot(data=temp,aes(x=date, y=value,color=name))+
12     geom_line()+
13     geom_point()+
14     scale_color_manual(values =c("blue","red","green"))+
15     labs(title="Bieu do tu vong tích lũy tháng 6 năm 2020",x="Date of month",y="Value")
16
17 temp <- cumsum_2020_10%>%
18   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
19   ggplot(data=temp,aes(x=date, y=value,color=name))+
20     geom_line()+
21     geom_point()+
22     scale_color_manual(values =c("blue","red","green"))+
23     labs(title="Bieu do tu vong tích lũy tháng 10 năm 2020",x="Date of month",y="Value")
24
25 temp <- cumsum_2021_2%>%
26   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
27   ggplot(data=temp,aes(x=date, y=value,color=name))+
28     geom_line()+
29     geom_point()+
30     scale_color_manual(values =c("blue","red","green"))+
31     labs(title="Bieu do tu vong tích lũy tháng 2 năm 2021",x="Date of month",y="Value")
```

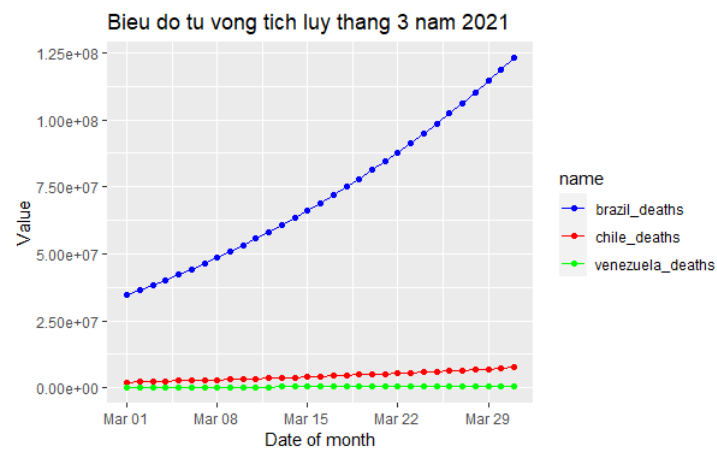
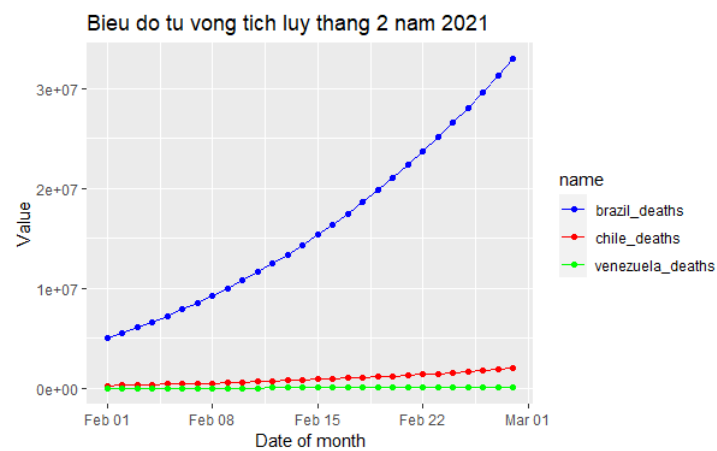
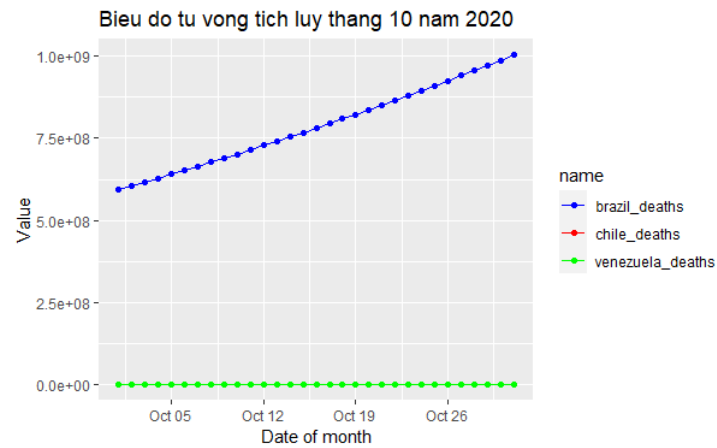
```

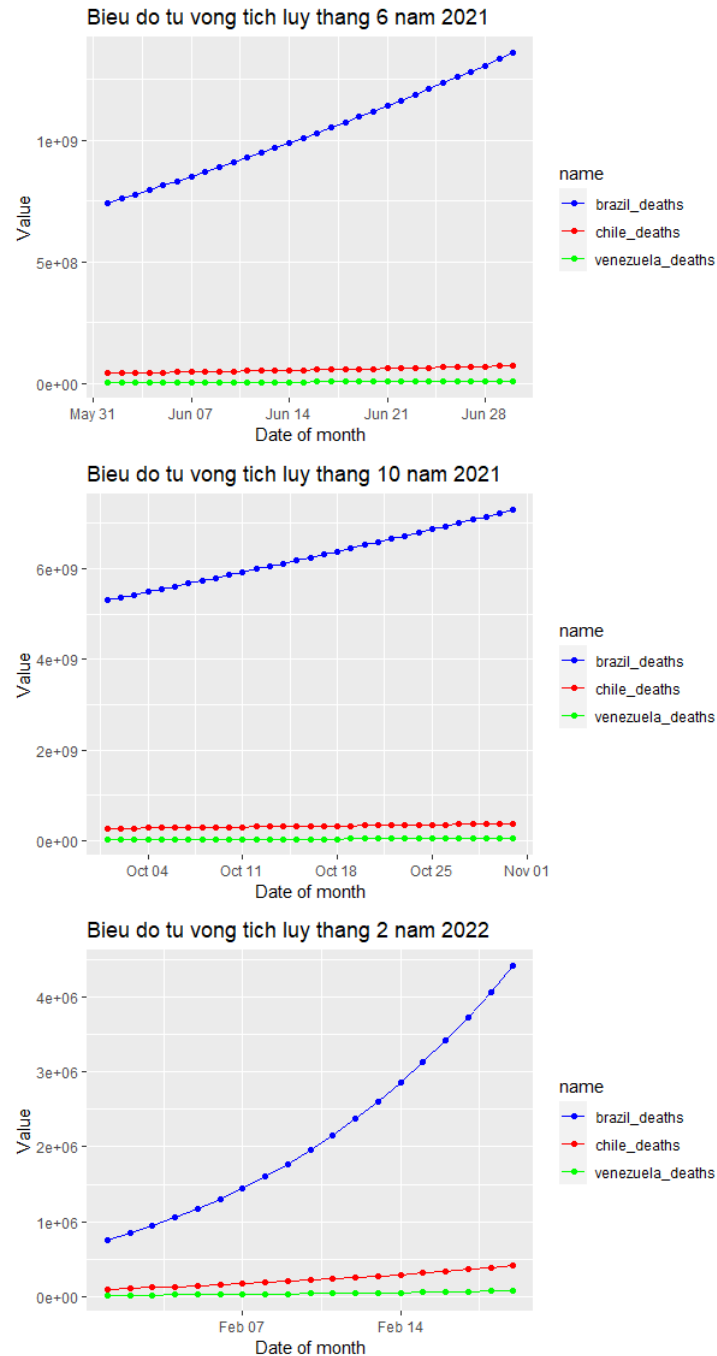
32
33 temp <- cumsum_2021_3%>%
34   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
35 ggplot(data=temp,aes(x=date, y=value,color=name))+
36   geom_line()+
37   geom_point()+
38   scale_color_manual(values =c("blue","red","green"))+
39   labs(title="Bieu do tu vong tích luy thang 3 nam 2021",x="Date of month",y="Value")
40
41 temp <- cumsum_2021_6%>%
42   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
43 ggplot(data=temp,aes(x=date, y=value,color=name))+
44   geom_line()+
45   geom_point()+
46   scale_color_manual(values =c("blue","red","green"))+
47   labs(title="Bieu do tu vong tích luy thang 6 nam 2021",x="Date of month",y="Value")
48
49 temp <- cumsum_2021_10%>%
50   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
51 ggplot(data=temp,aes(x=date, y=value,color=name))+
52   geom_line()+
53   geom_point()+
54   scale_color_manual(values =c("blue","red","green"))+
55   labs(title="Bieu do tu vong tích luy thang 10 nam 2021",x="Date of month",y="Value")
56
57 temp <- cumsum_2022_2%>%
58   tidyr::pivot_longer(cols = c(brazil_deaths,chile_deaths,venezuela_deaths))
59 ggplot(data=temp,aes(x=date, y=value,color=name))+
60   geom_line()+
61   geom_point()+
62   scale_color_manual(values =c("blue","red","green"))+
63   labs(title="Bieu do tu vong tích luy thang 2 nam 2022",x="Date of month",y="Value")

```

• Kết quả:







vii) Nhóm câu hỏi liên quan đến tất cả quốc gia theo thời gian là tháng

- Trên từng năm hãy vẽ biểu đồ thể hiện trục Ox là thời gian, trục Oy là nhiễm bệnh/tử vong. Hãy dùng 4 ký số của mã đề để vẽ 4 tháng tương ứng theo ký số đó. Nếu ký số là 0 thì lấy tháng là 10. Số đường thể hiện trên đồ thị là số quốc gia.

- Chuẩn bị dữ liệu cho toàn bộ phần vii (phần code dùng chung cho cả phần vii)

```
1 library(tidyverse)
2 library(lubridate)
3 library(dplyr)
4 library(xlsx)
5 setwd("E:/BTL_CTRR")
6 data_covid = read.csv("owid-covid-data.csv", header = TRUE)
7 library("ggplot2")
8 data_covid = mutate(data_covid, new_cases = abs(new_cases))
9 data_covid = mutate(data_covid, new_deaths = abs(new_deaths))
```

```

10 data_covid <- subset(data_covid,data_covid$continent != "")
11 data_covid[is.na(data_covid)] <- 0
12 data <- data_covid %>%
13   separate(date, into = c("month", "day", "year"), sep = "/")
14 data <- transform(data,month=as.numeric(month))
15
16 year_2020 <- data[data$year=="2020",]
17 year_2021 <- data[data$year=="2021",]
18 year_2022 <- data[data$year=="2022",]
19
20 year_2020 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2020, FUN = sum)
21 year_2021 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2021, FUN = sum)
22 year_2022 <- aggregate(cbind(new_cases,new_deaths) ~ month, data = year_2022, FUN = sum)
23
24 data_2020<-subset(year_2020,year_2020$month==2 | year_2020$month==10 |
25   year_2020$month==3 | year_2020$month==6)
26 data_2021<-subset(year_2021,year_2021$month==2 | year_2021$month==10 |
27   year_2021$month==3 | year_2021$month==6)
28 data_2022<-subset(year_2022,year_2022$month==2 | year_2022$month==10 |
29   year_2022$month==3 | year_2022$month==6)
30
31 data_2020_1<-subset(year_2020,year_2020$month==11 | year_2020$month==12)
32 data_2021_1<-subset(year_2021,year_2021$month==11 | year_2021$month==12)
33 data_2022_1<-subset(year_2022,year_2022$month==2 | year_2022$month==1)

```

1) Biểu đồ thể hiện thu thập dữ liệu nhiễm bệnh theo thời gian là tháng của tất cả quốc gia

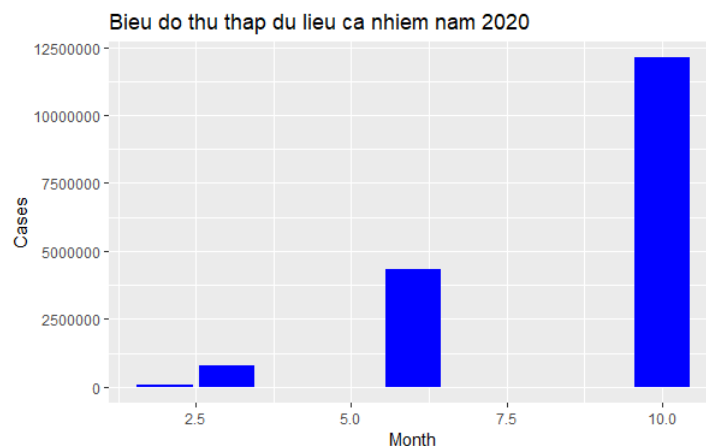
- Hiện thực trong R

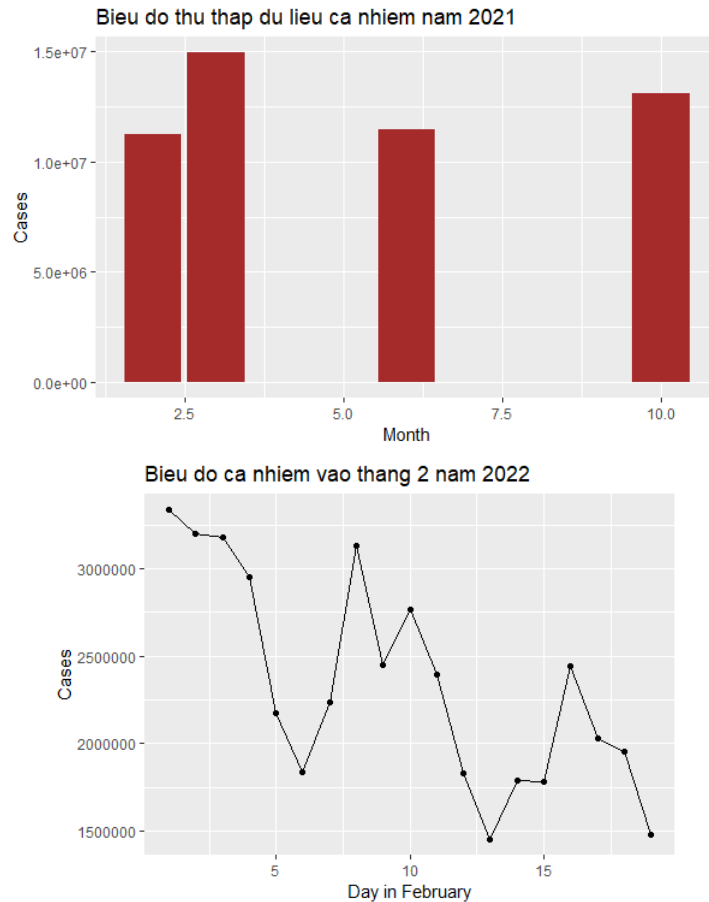
```

1 ggplot(data_2020, aes(y=new_cases, x=month)) +
2   geom_bar(position='dodge', stat='identity', fill ="blue") +
3   labs(x='Month', y='Cases', title='Biểu đồ thu thập dữ liệu ca nhiễm năm 2020')
4
5 ggplot(data_2021, aes(y=new_cases, x=month)) +
6   geom_bar(position='dodge', stat='identity', fill ="brown") +
7   labs(x='Month', y='Cases', title='Biểu đồ thu thập dữ liệu ca nhiễm năm 2021')
8
9 temp = data[data$year=="2022",]
10 temp = temp[temp$month==2,]
11 data_2020_new<- aggregate(cbind(new_cases,new_deaths) ~ day, data = temp, FUN = sum)
12 data_2020_new <- arrange(data_2020_new,day=as.numeric(day))
13 ggplot(data_2020_new,aes(x=as.numeric(day), y=new_cases))+
14   geom_line()+
15   geom_point()+
16   labs(title="Biểu đồ ca nhiễm vào tháng 2 năm 2022",x="Day in February",y="Cases")

```

- Kết quả:



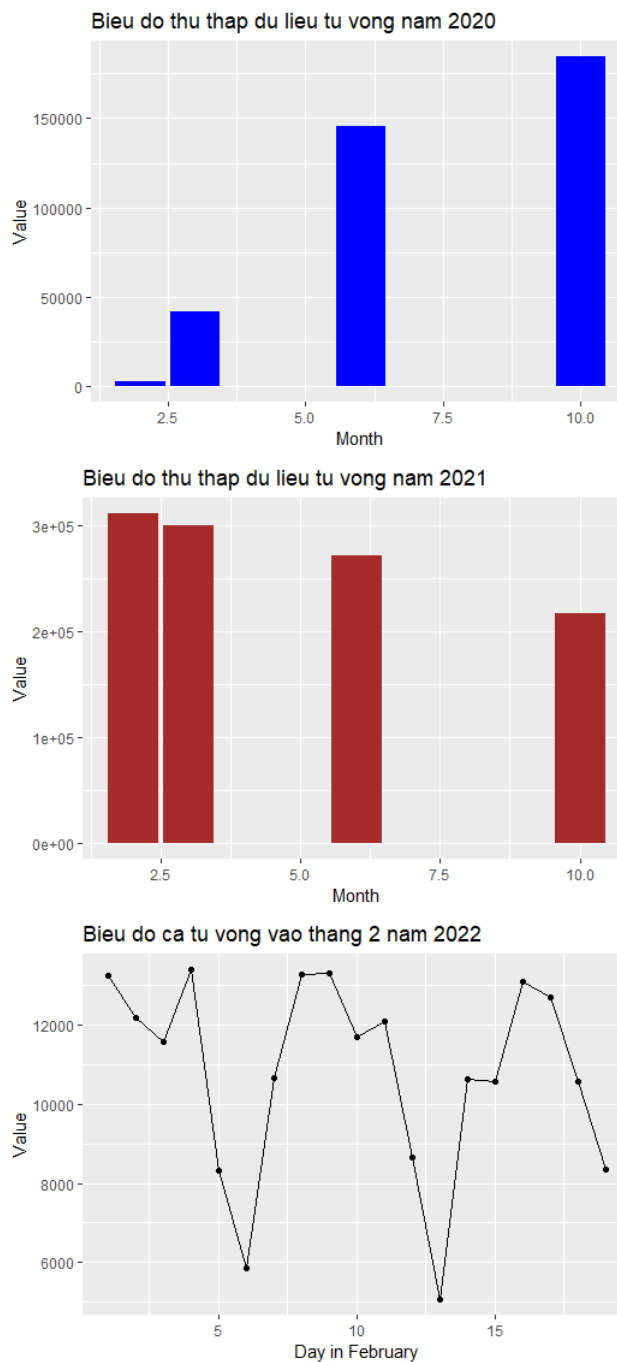


2) Biểu đồ thể hiện thu thập dữ liệu tử vong theo thời gian là tháng của tất cả quốc gia

- Hiện thực trong R

```
1 ggplot(data_2020, aes(y=new_deaths, x=month)) +
2   geom_bar(position='dodge', stat='identity', fill="blue") +
3   labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong nam 2020')
4
5 ggplot(data_2021, aes(y=new_deaths, x=month)) +
6   geom_bar(position='dodge', stat='identity', fill="brown") +
7   labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong nam 2021')
8
9 ggplot(data=data_2020_new, aes(x=as.numeric(day), y=new_deaths))+
10  geom_line()+
11  geom_point()+
12  labs(title="Bieu do ca tu vong vào tháng 2 năm 2022", x="Day in February", y="Value")
```

- Kết quả:

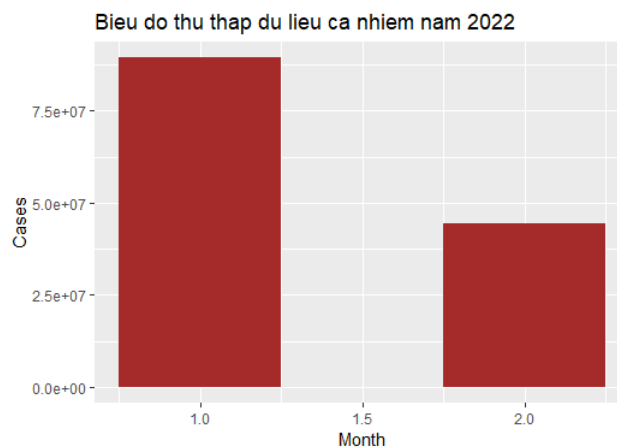
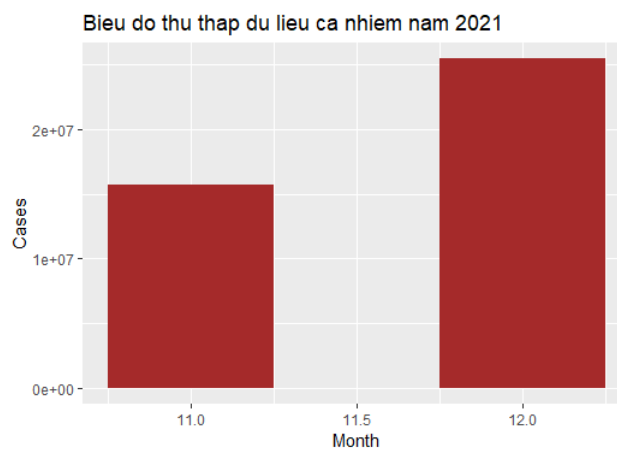
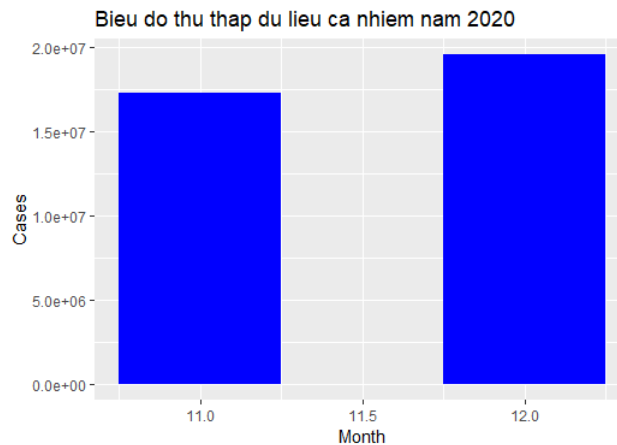


3) Biểu đồ thể hiện thu thập dữ liệu nhiễm bệnh theo thời gian là 2 tháng cuối của năm của tất cả quốc gia

- Hiện thực trong R

```
1 ggplot(data_2020_1, aes(y=new_cases, x=month,width = 0.5)) +
2   geom_bar(position='dodge', stat='identity', fill ="blue") +
3   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem nam 2020')
4
5 ggplot(data_2021_1, aes(y=new_cases, x=month,width = 0.5)) +
6   geom_bar(position='dodge', stat='identity', fill ="brown") +
7   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem nam 2021')
8
9 ggplot(data_2022_1, aes(y=new_cases, x=month,width = 0.5)) +
10  geom_bar(position='dodge', stat='identity', fill ="brown") +
11  labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem nam 2022')
```


- Kết quả:

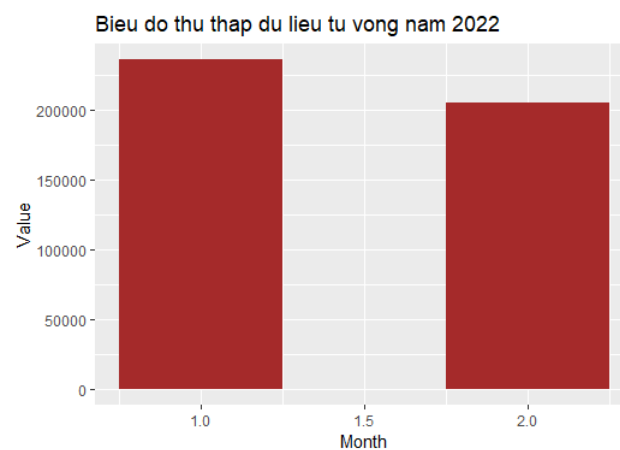
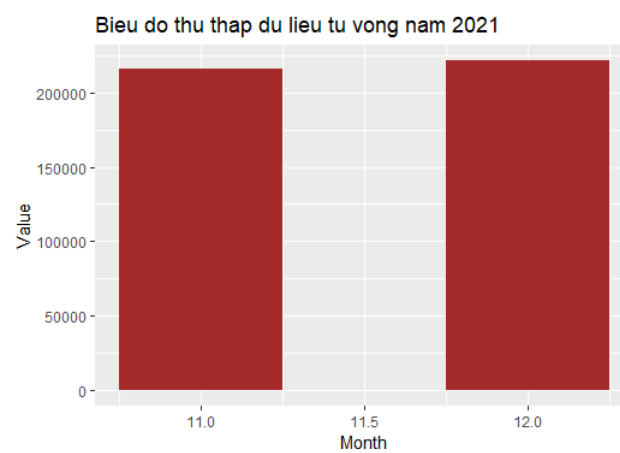
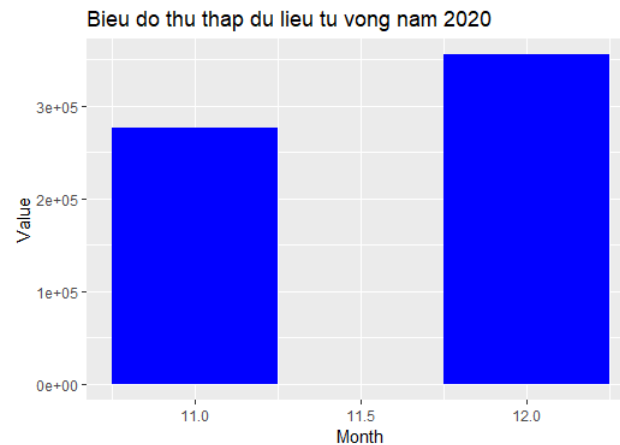


- 4) Biểu đồ thể hiện thu thập dữ liệu tử vong theo thời gian là 2 tháng cuối của năm của tất cả quốc gia

- Hiện thực trong R

```
1 ggplot(data_2020_1, aes(y=new_deaths, x=month, width = 0.5)) +
2   geom_bar(position='dodge', stat='identity', fill ="blue") +
3   labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong nam 2020')
4
5 ggplot(data_2021_1, aes(y=new_deaths, x=month,width = 0.5)) +
6   geom_bar(position='dodge', stat='identity', fill ="brown") +
7   labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong nam 2021')
8
9 ggplot(data_2022_1, aes(y=new_deaths, x=month,width = 0.5)) +
10  geom_bar(position='dodge', stat='identity', fill ="brown") +
11  labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong nam 2022')
```

- Kết quả:



- 5) Biểu đồ thể hiện thu thập dữ liệu nhiễm bệnh tương đối tích lũy theo thời gian là 2 tháng cuối của năm của tất cả quốc gia

- Hiện thực trong R

```

1  cumsum_2020=as.data.frame(lapply(year_2020,cumsum))
2  cumsum_2020$month = c(1:12)
3  cumsum_2021=as.data.frame(lapply(year_2021,cumsum))
4  cumsum_2021$month = c(1:12)
5  cumsum_2022=as.data.frame(lapply(year_2022,cumsum))
6  cumsum_2022$month = c(1,2)
7
8  cumsum_2020_1<-subset(cumsum_2020,cumsum_2020$month==11 | cumsum_2020$month==12)
9  cumsum_2021_1<-subset(cumsum_2021,cumsum_2021$month==11 | cumsum_2021$month==12)
10 cumsum_2022_1<-subset(cumsum_2022,cumsum_2022$month==2 | cumsum_2022$month==1)
11

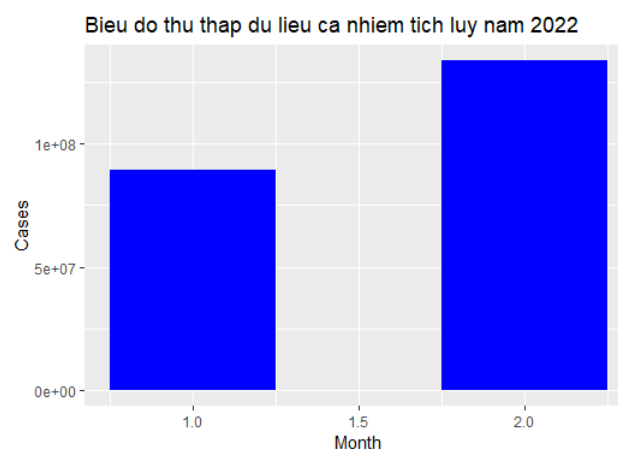
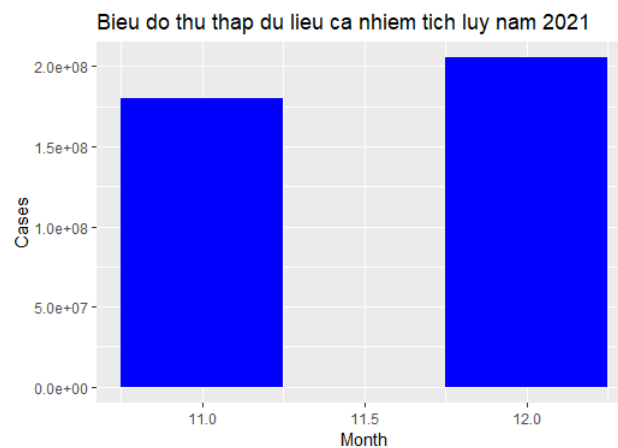
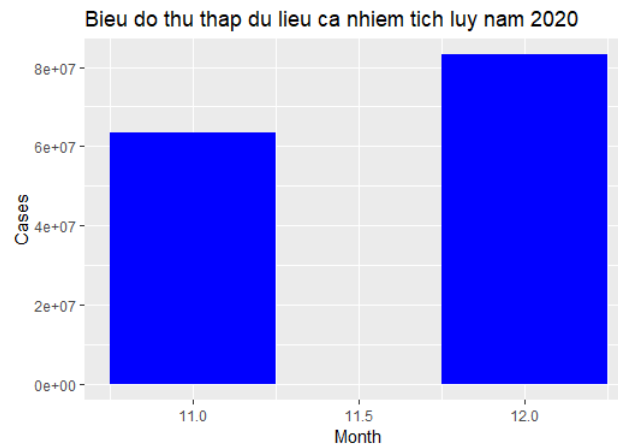
```

```

12 ggplot(cumsum_2020_1, aes(y=new_cases, x=month, width = 0.5)) +
13   geom_bar(position='dodge', stat='identity', fill = "blue") +
14   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2020')
15
16 ggplot(cumsum_2021_1, aes(y=new_cases, x=month, width = 0.5)) +
17   geom_bar(position='dodge', stat='identity', fill = "blue") +
18   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2021')
19
20 ggplot(cumsum_2022_1, aes(y=new_cases, x=month, width = 0.5)) +
21   geom_bar(position='dodge', stat='identity', fill = "blue") +
22   labs(x='Month', y='Cases', title='Bieu do thu thap du lieu ca nhiem tích luy nam 2022')

```

- Kết quả:

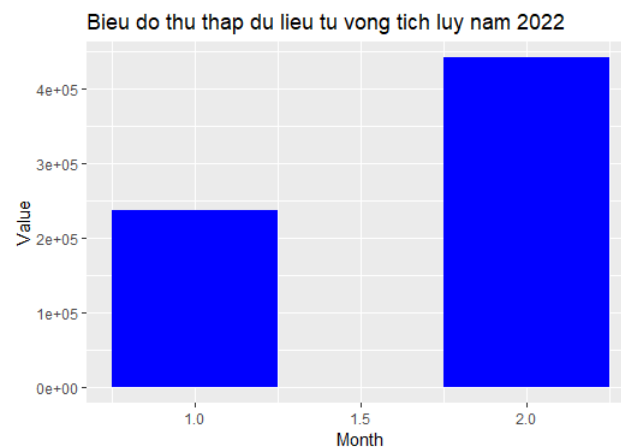
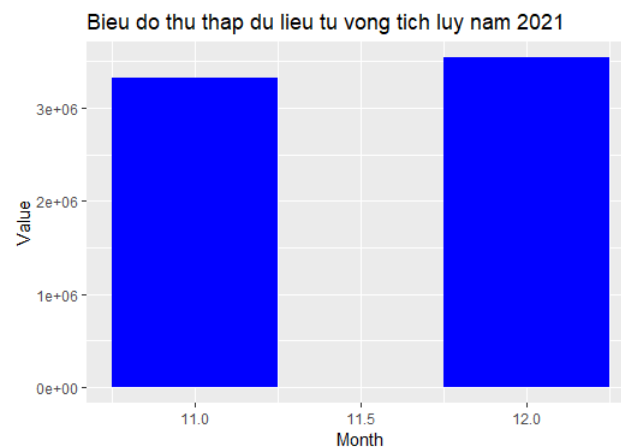
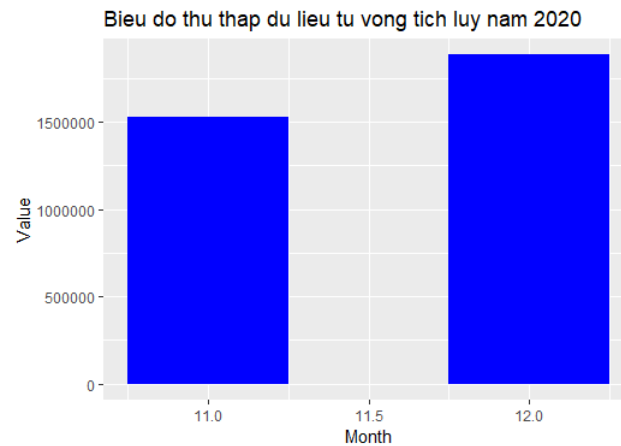


6) Biểu đồ thu thập tử vong tương đối tích lũy theo thời gian là 2 tháng cuối của tất cả quốc gia

- Hiện thực trong R

```
1 ggplot(cumsum_2020_1, aes(y=new_deaths, x=month, width = 0.5)) +
2   geom_bar(position='dodge', stat='identity', fill = "blue") +
3   labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong tích luy nam 2020')
4
5 ggplot(cumsum_2021_1, aes(y=new_deaths, x=month, width = 0.5)) +
6   geom_bar(position='dodge', stat='identity', fill = "blue") +
7   labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong tích luy nam 2021')
8
9 ggplot(cumsum_2022_1, aes(y=new_deaths, x=month, width = 0.5)) +
10  geom_bar(position='dodge', stat='identity', fill = "blue") +
11  labs(x='Month', y='Value', title='Bieu do thu thap du lieu tu vong tích luy nam 2022')
```

- Kết quả:



viii) Nhóm câu hỏi liên quan đến tất cả quốc gia theo trung bình 7 ngày gần nhất

Trên từng năm hãy vẽ biểu đồ thể hiện trục Ox là thời gian, trục Oy là nhiễm bệnh/tử vong. Hãy

dùng 4 ký số của mã để vẽ 4 tháng tương ứng theo ký số đó. Nếu ký số là 0 thì lấy tháng là 10.

- Chuẩn bị dữ liệu cho toàn bộ phần viii (phần code dùng chung cho cả phần viii)

```
1 library(ggplot2)
2 library(readr)
3 library(runner)
4 library(dplyr)
5
6 covid <- read_csv("owid-covid-data.csv")
7 covid$new_cases <- abs(covid$new_cases)
8 covid$new_deaths <- abs(covid$new_deaths)
9 covid$date <- as.Date(covid$date, "%m/%d/%Y")
10 covid$month <- as.numeric(format(covid$date, '%m'))
11 covid$year <- as.numeric(format(covid$date, '%Y'))
12 covid <- subset(covid, !is.na(continent))
13 covid$new_cases[is.na(covid$new_cases)] <- 0
14 covid$new_deaths[is.na(covid$new_deaths)] <- 0
```

- 1) Biểu đồ thể hiện thu thập dữ liệu nhiễm bệnh theo thời gian là tháng của tất cả quốc gia theo trung bình 7 ngày gần nhất

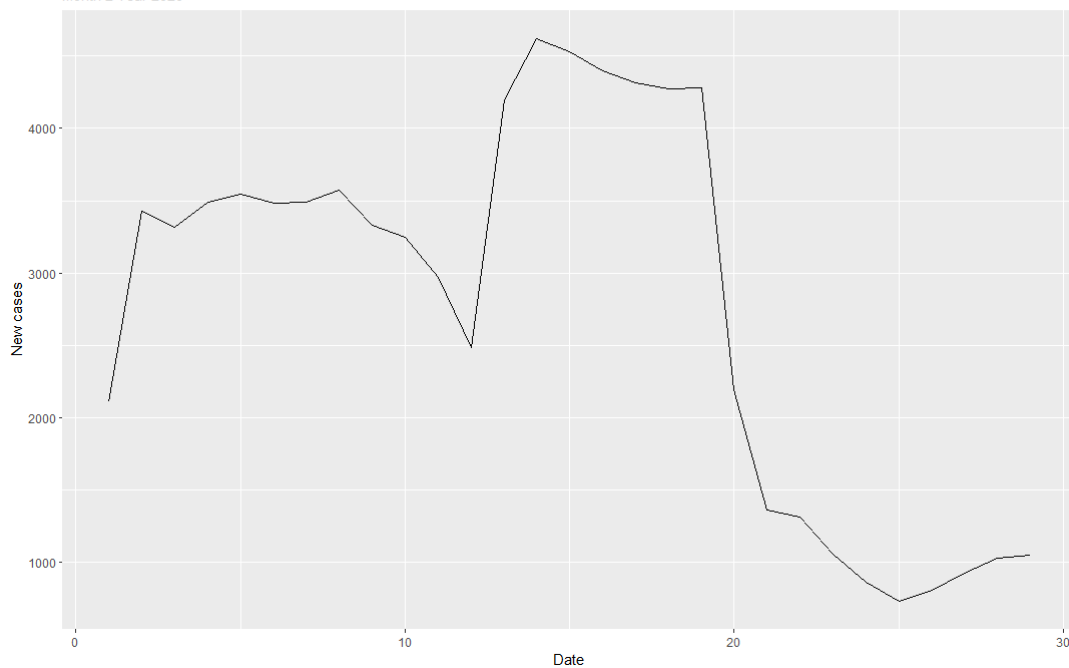
- Hiện thực trong R

```
1 figure <- function(m,y,name)
2 {
3   task <- subset(covid, month == m & year == y)
4   task <- aggregate(task$new_cases, list(task$date), FUN=sum)
5   colnames(task) <- c("date", "new_cases")
6   task$rec <- 1
7   task$rec <- sum_run(x = task$rec, k = 7, i = as.Date(task$date, format = "%m/%d/%Y"))
8   task$sum <- sum_run(x = task$new_cases, k = 7, i = as.Date(task$date, format = "%m/%d/%Y"))
9   task$sum <- task$sum/task$rec
10  task$day <- as.numeric(format(task$date, '%d'))
11
12  graph <- ggplot(task, aes(x=day, y=sum)) + geom_line() + theme(plot.title = element_text(
13    size = 15, face = "bold")) + labs(title = "Average for past 7 days", subtitle=name, x= "
14    Date", y= "New cases")
15  graph
16 }
17
18 figure(2,2020,"Month 2 Year 2020")
19 figure(3,2020,"Month 3 Year 2020")
20 figure(6,2020,"Month 6 Year 2020")
21 figure(10,2020,"Month 10 Year 2020")
22 figure(2,2021,"Month 2 Year 2021")
23 figure(3,2021,"Month 3 Year 2021")
24 figure(6,2021,"Month 6 Year 2021")
25 figure(10,2021,"Month 10 Year 2021")
26 figure(2,2022,"Month 2 Year 2021")
```

- Kết quả:

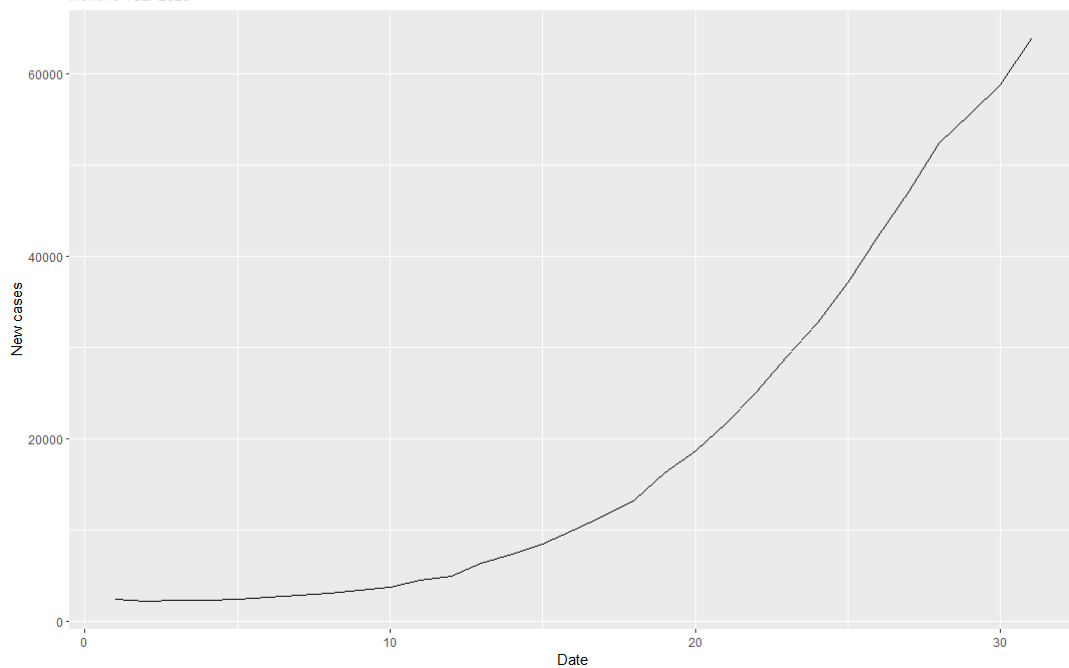
Average for past 7 days

Month 2 Year 2020



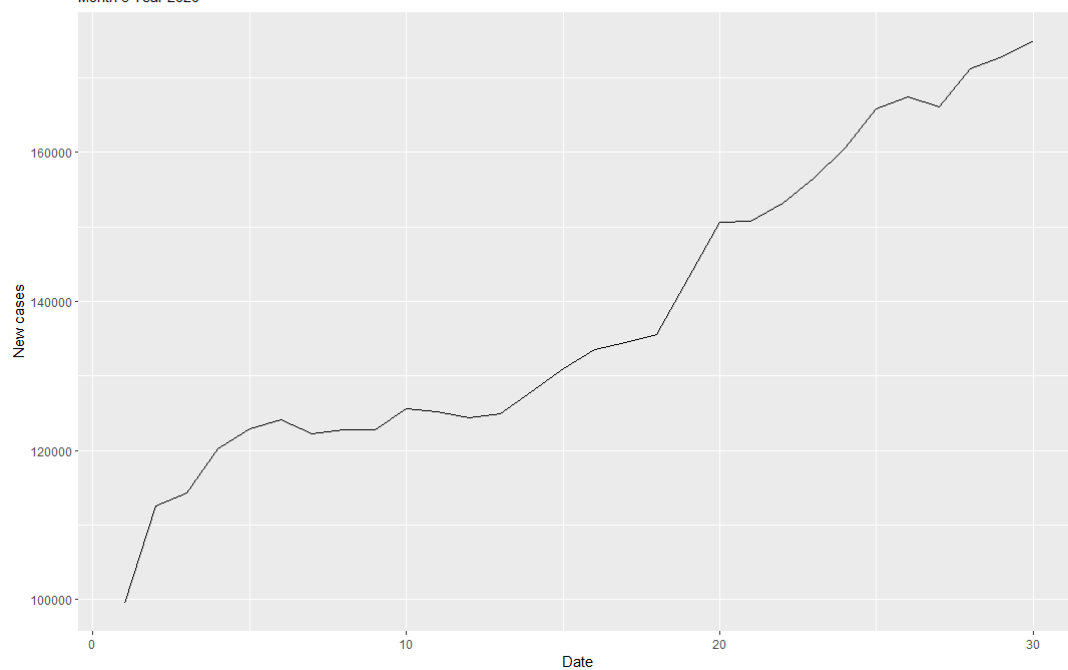
Average for past 7 days

Month 3 Year 2020



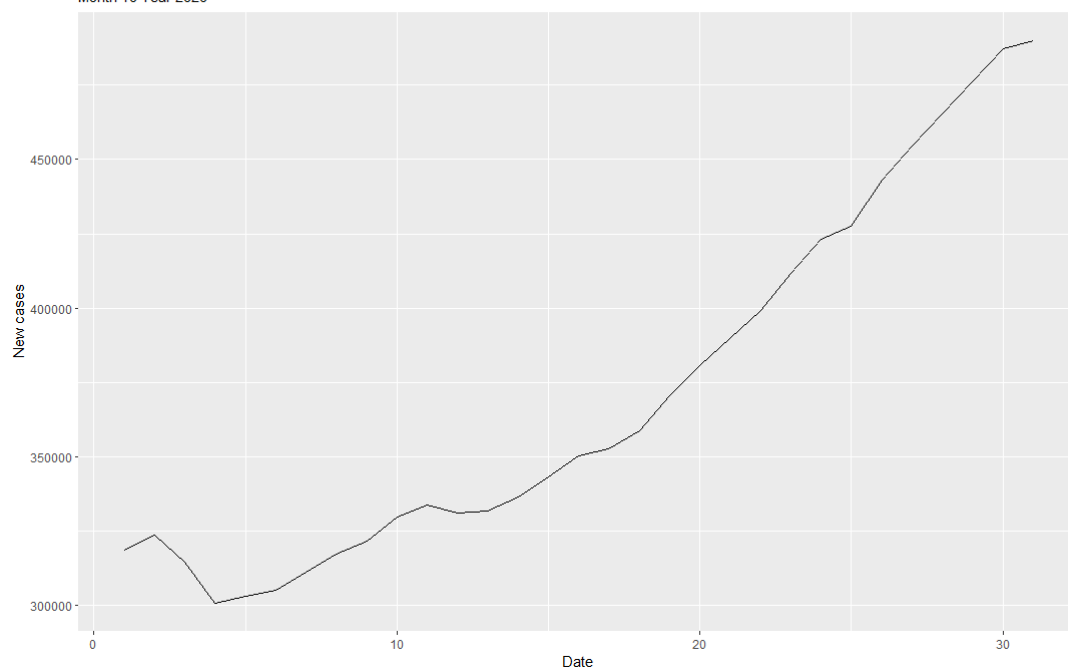
Average for past 7 days

Month 6 Year 2020



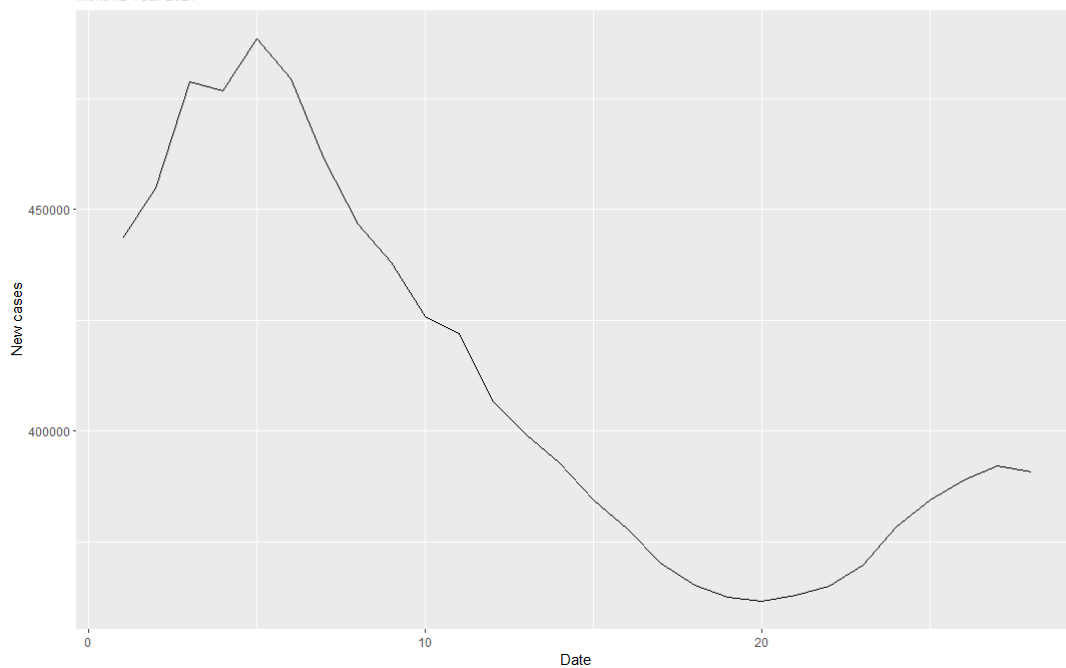
Average for past 7 days

Month 10 Year 2020



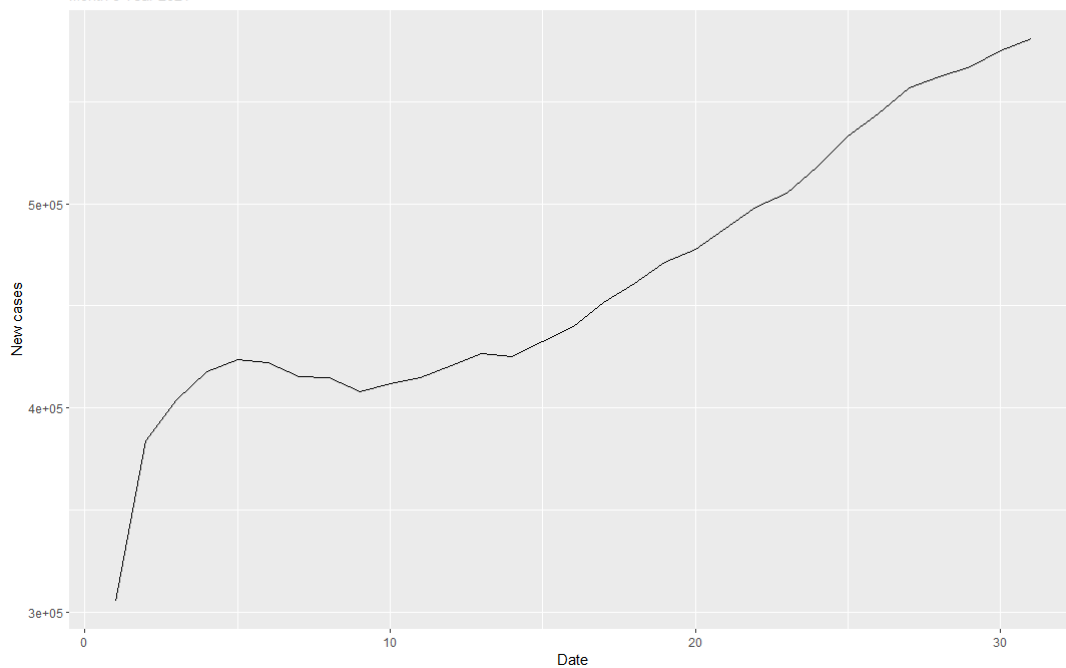
Average for past 7 days

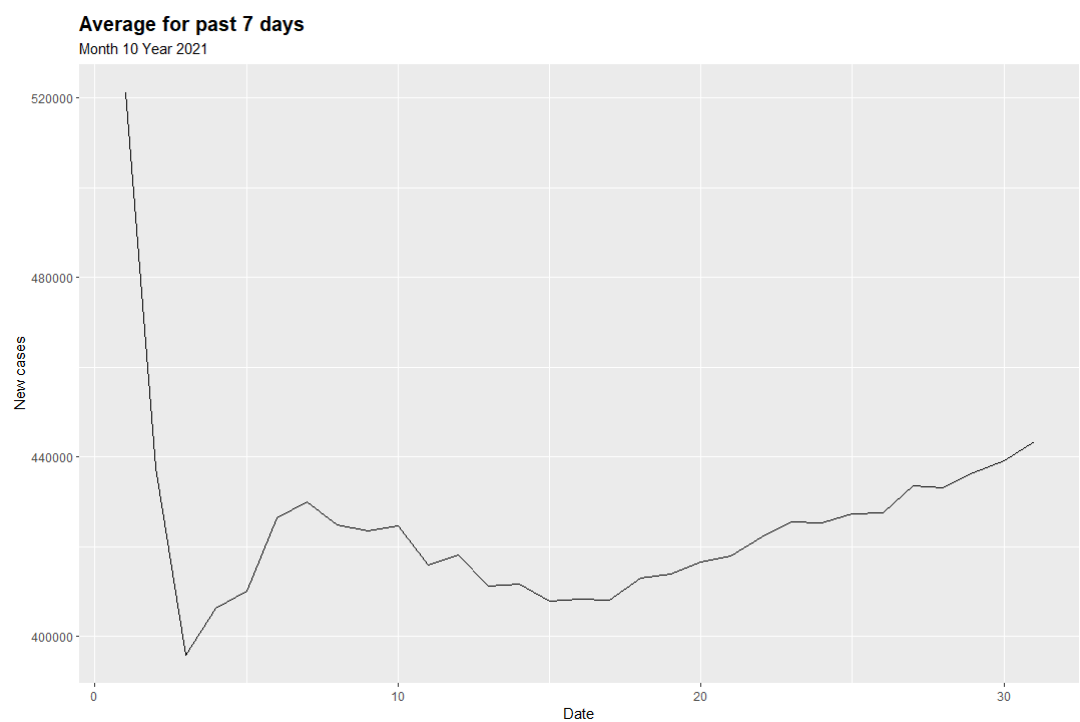
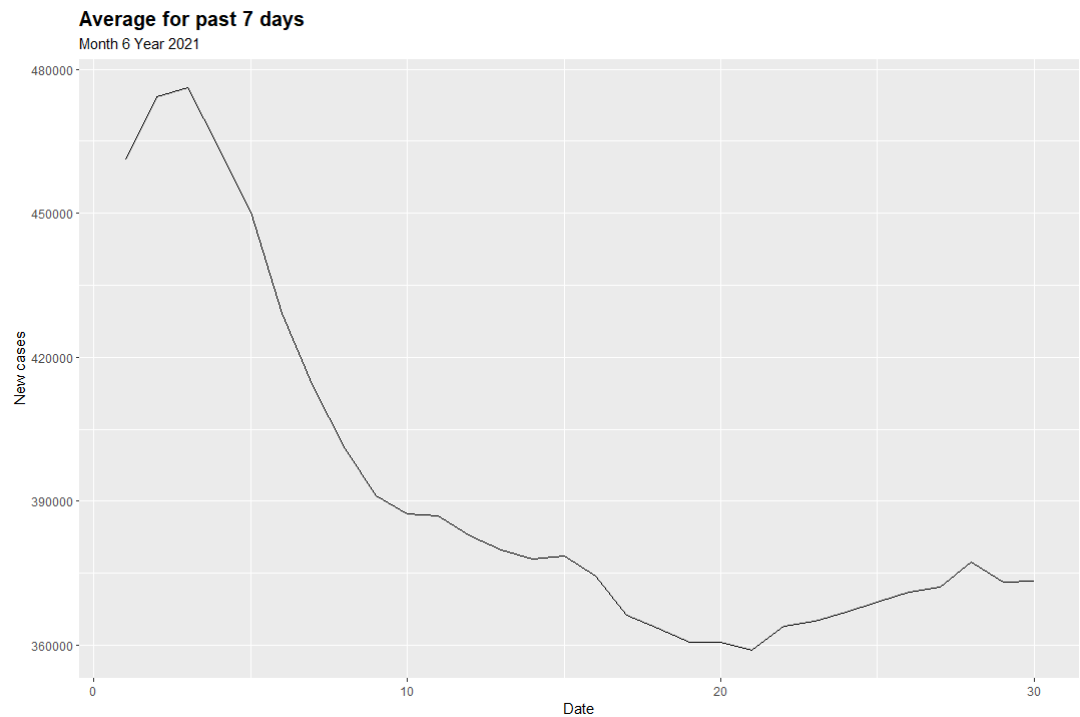
Month 2 Year 2021

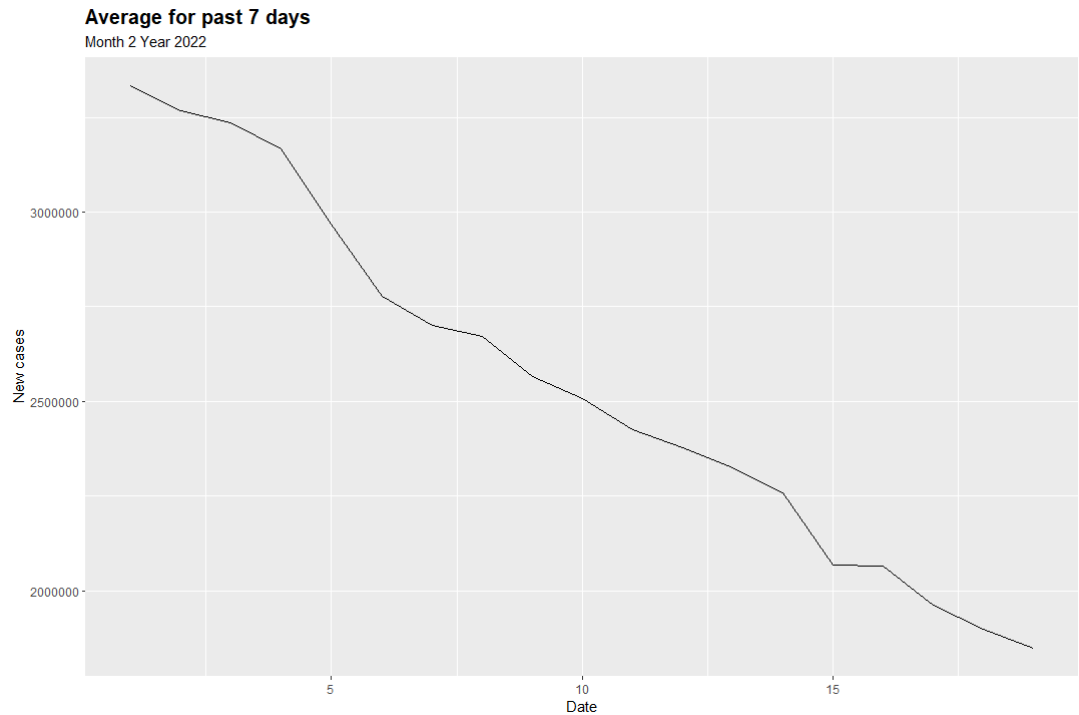


Average for past 7 days

Month 3 Year 2021







2) Biểu đồ thể hiện thu thập dữ liệu tử vong theo thời gian là tháng của tất cả quốc gia theo trung bình 7 ngày gần nhất

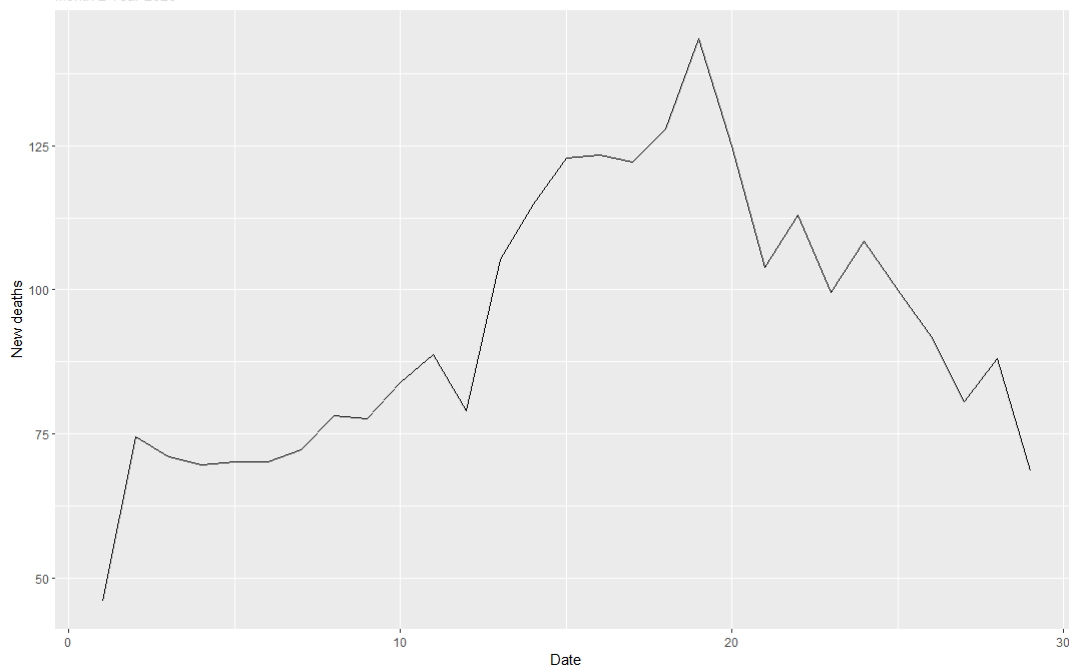
- Hiện thực trong R

```
1 figure <- function(m,y,name)
2 {
3   task <- subset(covid, month == m & year == y)
4   task <- aggregate(task$new_deaths, list(task$date), FUN=sum)
5   colnames(task) <- c("date", "new_deaths")
6   task$rec <- 1
7   task$rec <- sum_run(x = task$rec, k = 7, i = as.Date(task$date, format = "%m/%d/%Y"))
8   task$sum <- sum_run(x = task$new_deaths, k = 7, i = as.Date(task$date, format = "%m/%d/%Y"))
9   task$sum <- task$sum/task$rec
10  task$day <- as.numeric(format(task$date, '%d'))
11
12  graph <- ggplot(task, aes(x=day, y=sum)) + geom_line() + theme(plot.title = element_text(
13    size = 15, face = "bold")) + labs(title="Average for past 7 days", subtitle=name, x= "Date",
14    y= "New deaths")
15  graph
16 }
17 figure(2,2020,"Month 2 Year 2020")
18 figure(3,2020,"Month 3 Year 2020")
19 figure(6,2020,"Month 6 Year 2020")
20 figure(10,2020,"Month 10 Year 2020")
21 figure(2,2021,"Month 2 Year 2021")
22 figure(3,2021,"Month 3 Year 2021")
23 figure(6,2021,"Month 6 Year 2021")
24 figure(10,2021,"Month 10 Year 2021")
25 figure(2,2022,"Month 2 Year 2022")
```

- Kết quả:

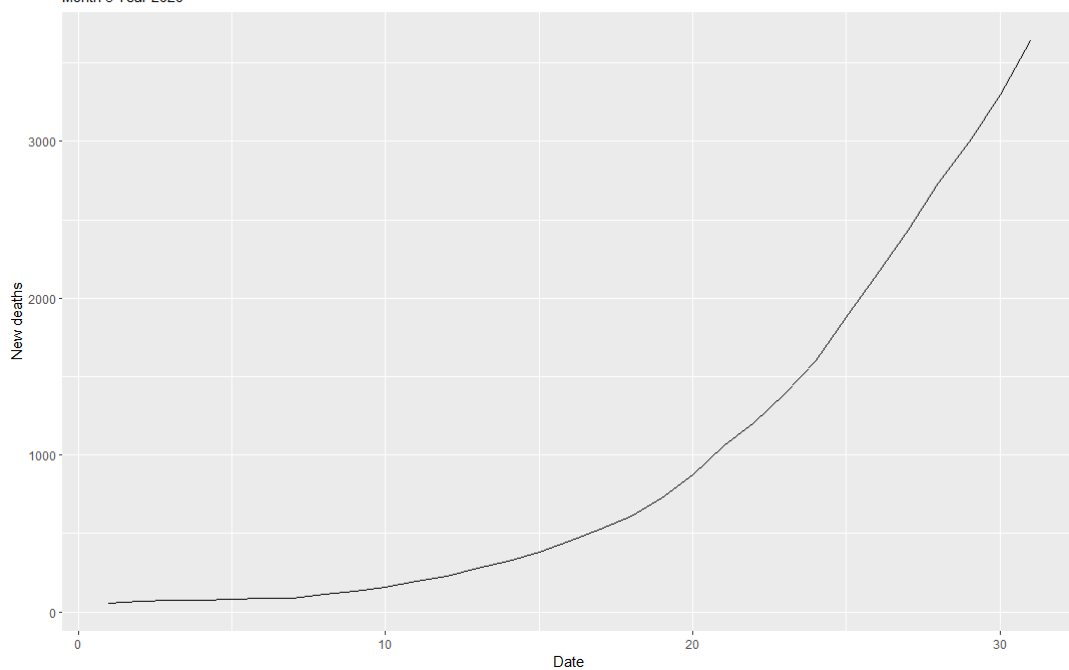
Average for past 7 days

Month 2 Year 2020



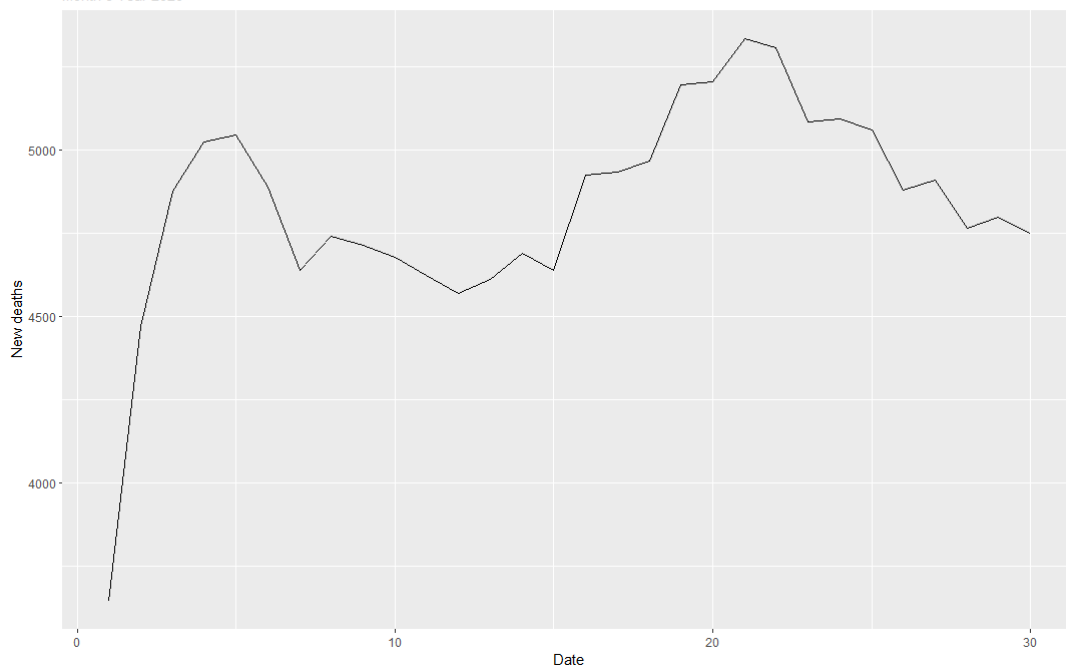
Average for past 7 days

Month 3 Year 2020



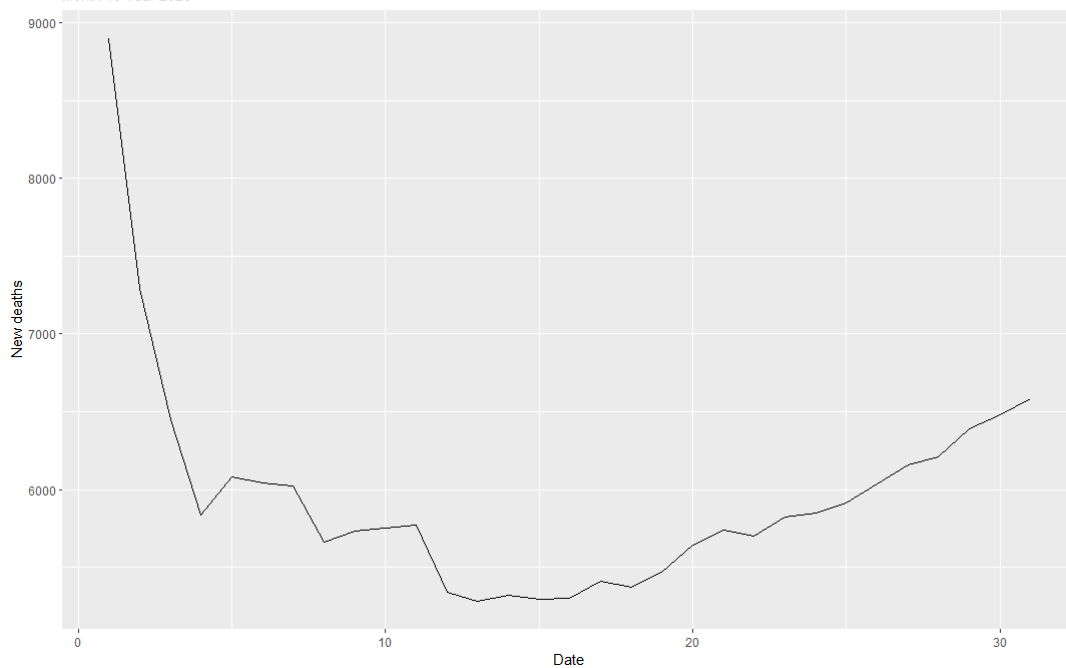
Average for past 7 days

Month 6 Year 2020



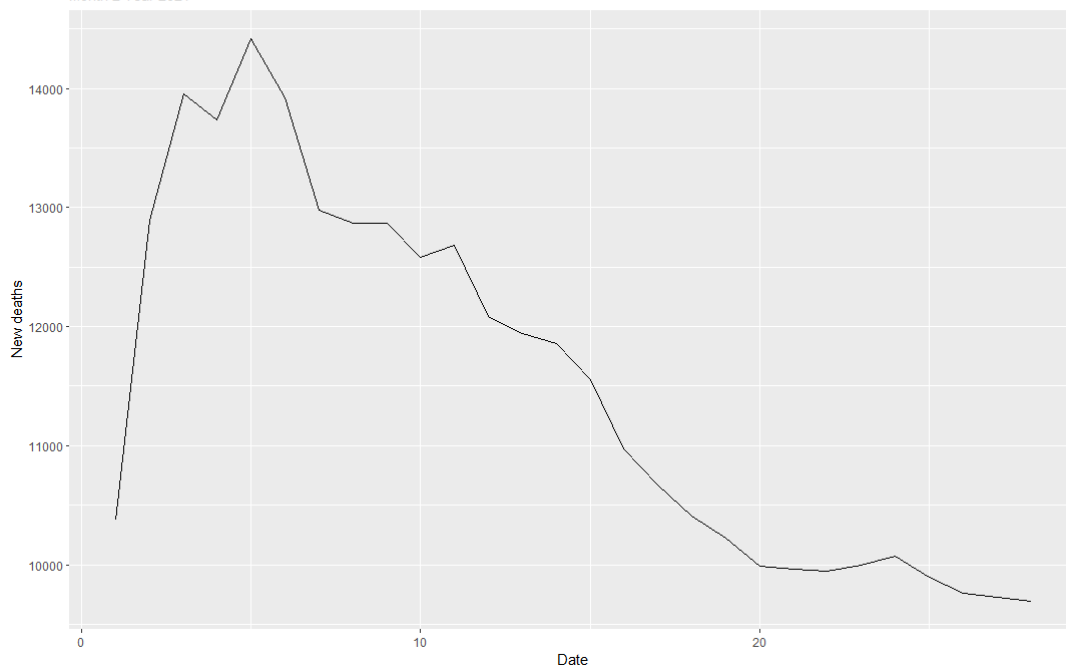
Average for past 7 days

Month 10 Year 2020



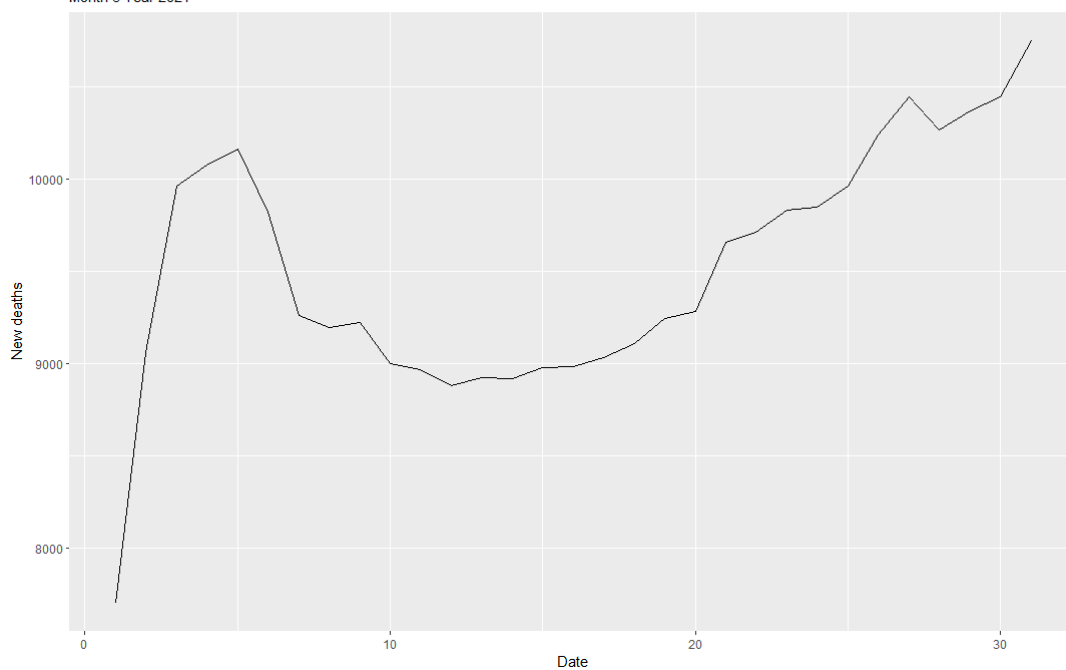
Average for past 7 days

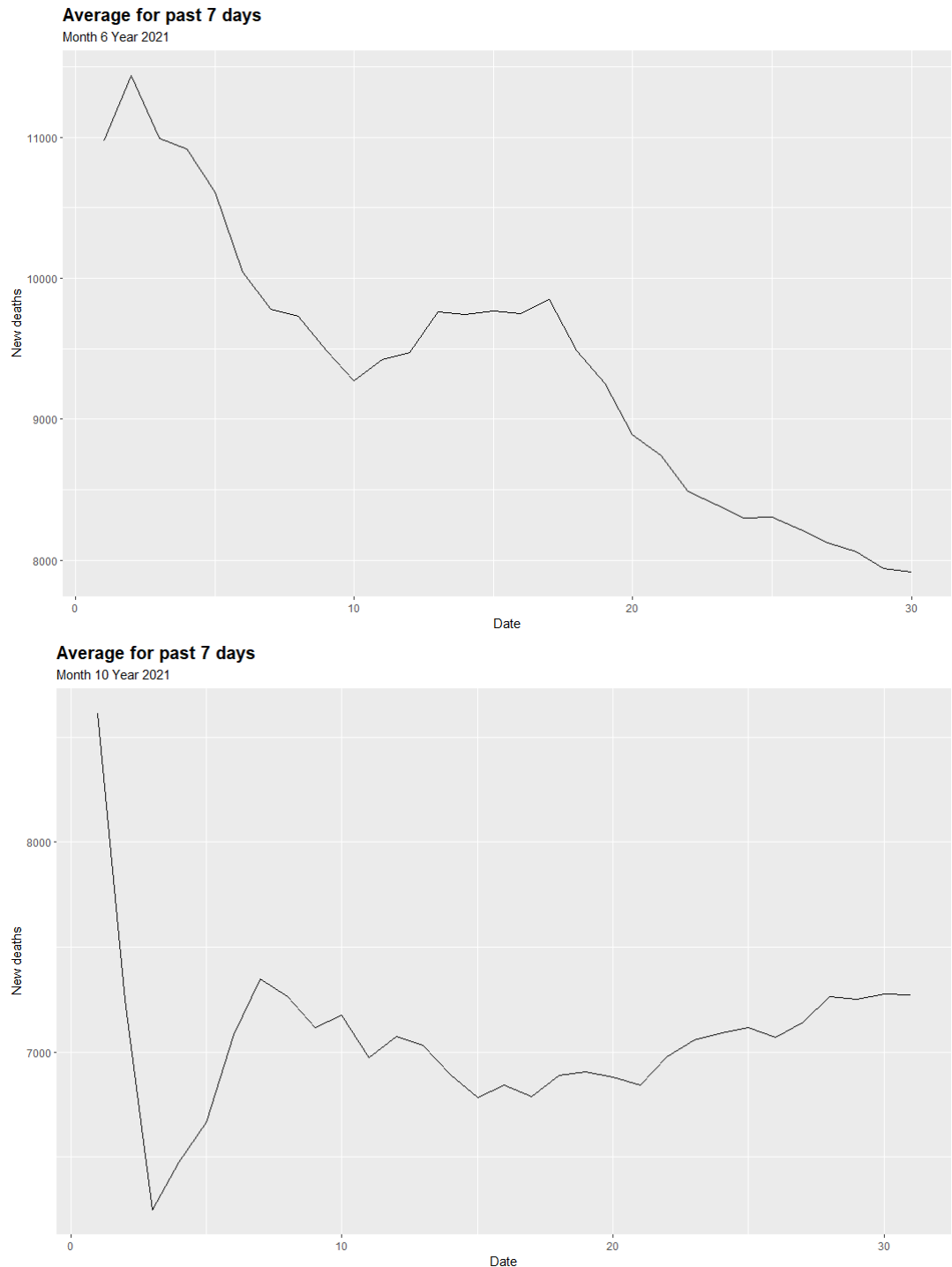
Month 2 Year 2021



Average for past 7 days

Month 3 Year 2021





3) Biểu đồ thể hiện thu thập dữ liệu nhiễm bệnh theo thời gian là 2 tháng của năm của tất cả quốc gia theo trung bình 7 ngày gần nhất

- Hiện thực trong R

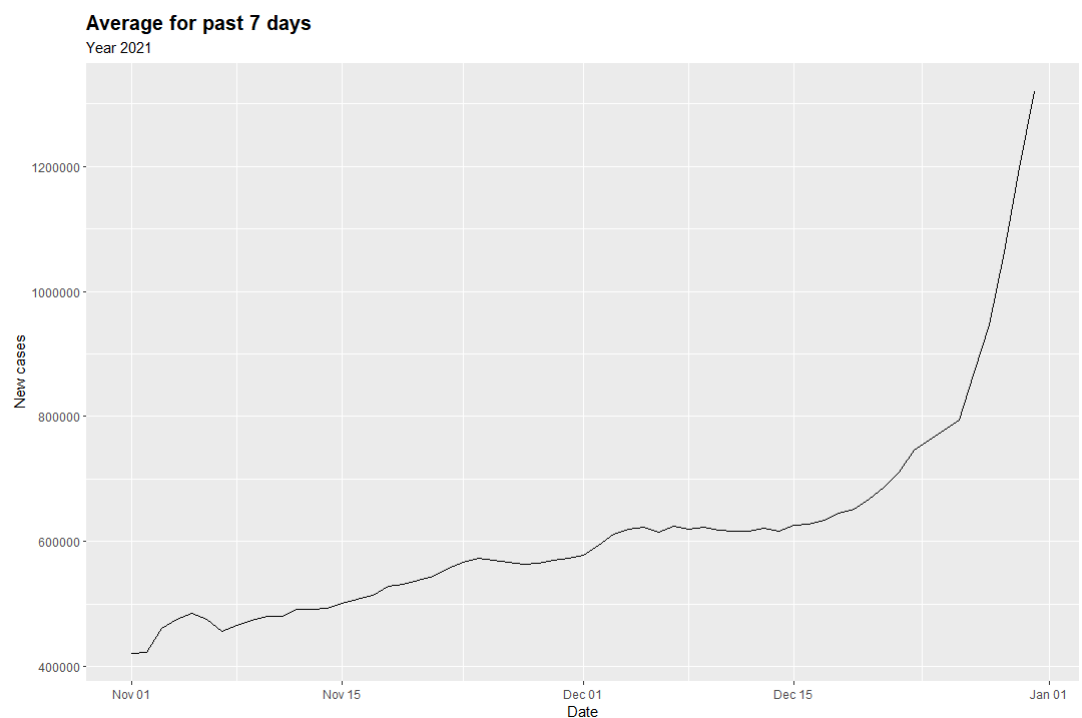
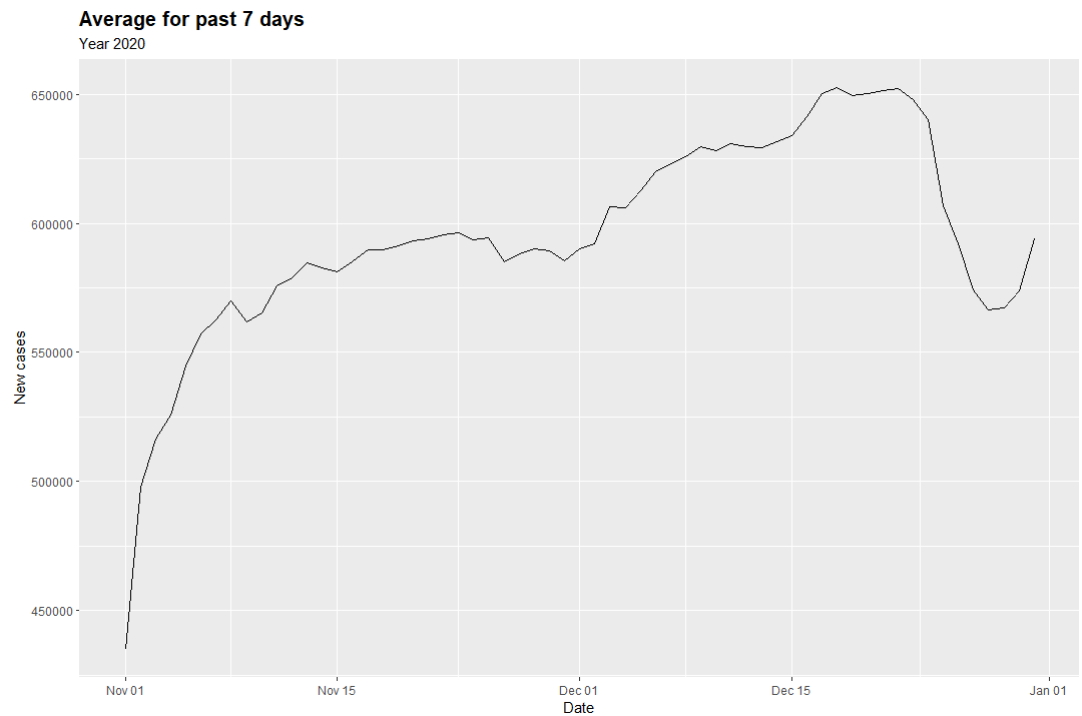
```
1 figure <- function(y,name)
2 {
3   task <- subset(covid, (month == 11|month == 12) & year == y)
4   task <- aggregate(task$new_cases, list(task$date), FUN=sum)
5   colnames(task) <- c("date","new_cases")
6
7   task$rec <- 1
8   task$rec <- sum_run(x = task$rec,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
9   task$sum <- sum_run(x = task$new_cases,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
```

```

10 task$sum <- task$sum/task$rec
11
12 graph <- ggplot(task, aes(x=date, y=sum)) + geom_line() + theme(plot.title = element_text(
13   size = 15, face = "bold")) + labs(title="Average for past 7 days", subtitle=name, x= "Date",
14   y= "New cases")
15 graph
16 }
17
18 figure(2020, "Year 2020")
19 figure(2021, "Year 2021")

```

- Kết quả:

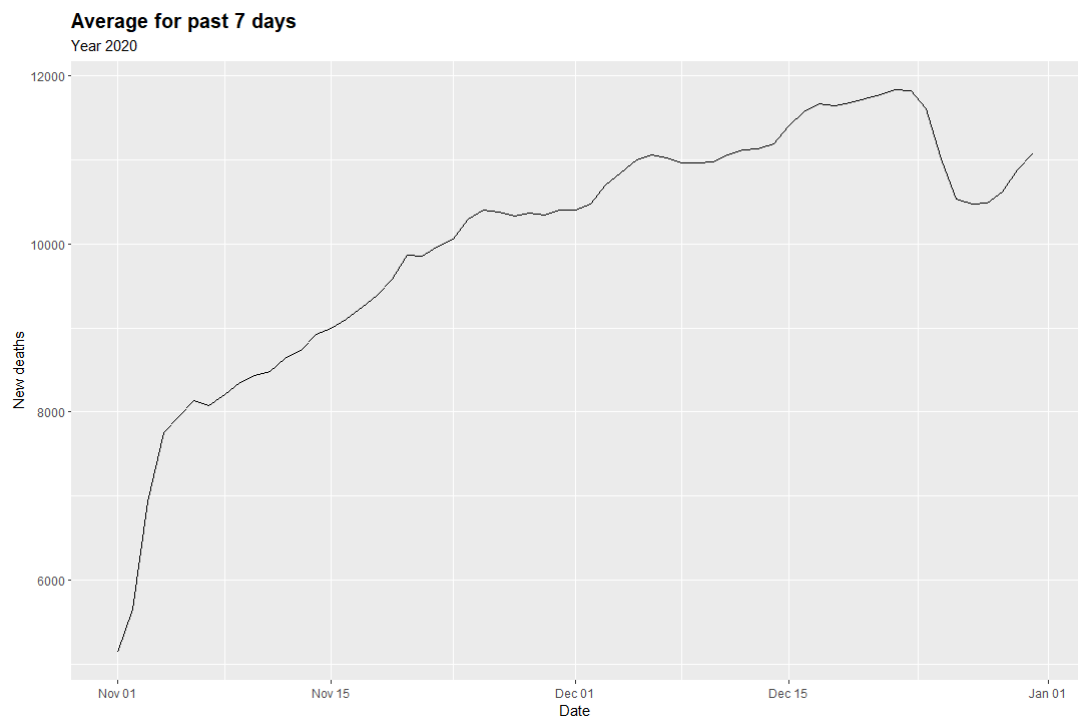


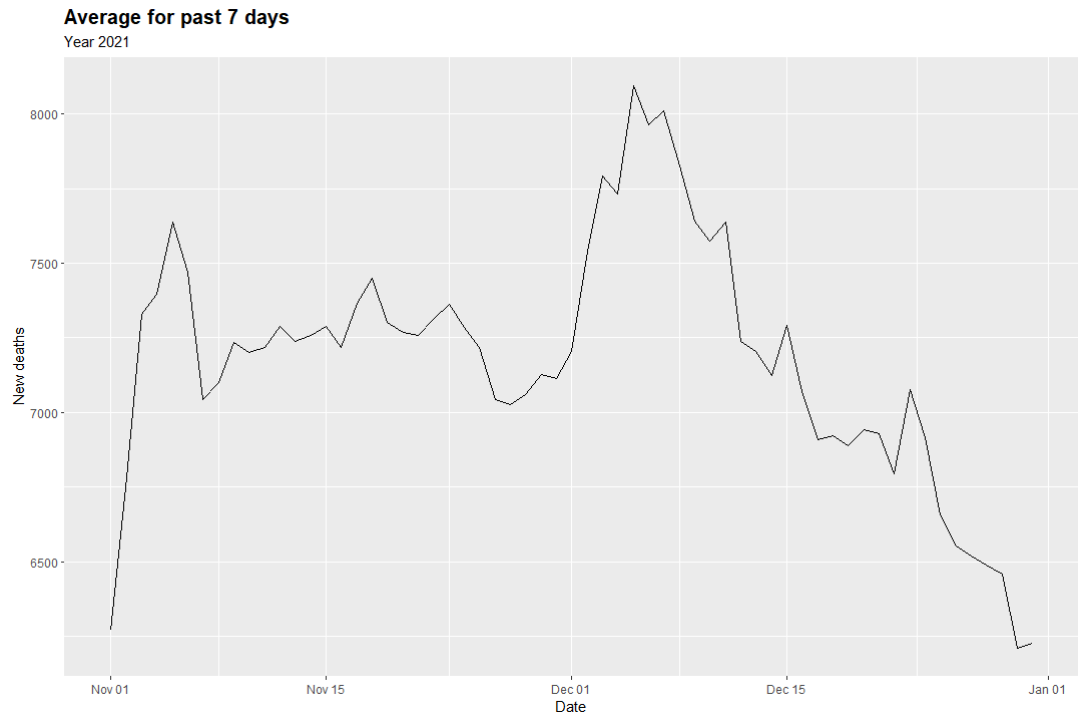
- 4) Biểu đồ thể hiện thu thập dữ liệu tử vong theo thời gian là 2 tháng của năm của tất cả quốc gia theo trung bình 7 ngày gần nhất

- Hiện thực trong R

```
1 figure <- function(y,name)
2 {
3   task <- subset(covid, (month == 11|month == 12) & year == y)
4   task <- aggregate(task$new_deaths, list(task$date), FUN=sum)
5   colnames(task) <- c("date","new_deaths")
6   task$rec <- 1
7   task$rec <- sum_run(x = task$rec,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
8   task$sum <- sum_run(x = task$new_deaths,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
9   task$sum <- task$sum/task$rec
10
11  graph <- ggplot(task, aes(x=date, y=sum)) + geom_line() + theme(plot.title = element_text(
12    size = 15, face = "bold")) + labs(title="Average for past 7 days", subtitle=name, x= "Date",
13    y= "New deaths")
14  graph
15 }
16 figure(2020,"Year 2020")
17 figure(2021,"Year 2021")
```

- Kết quả:



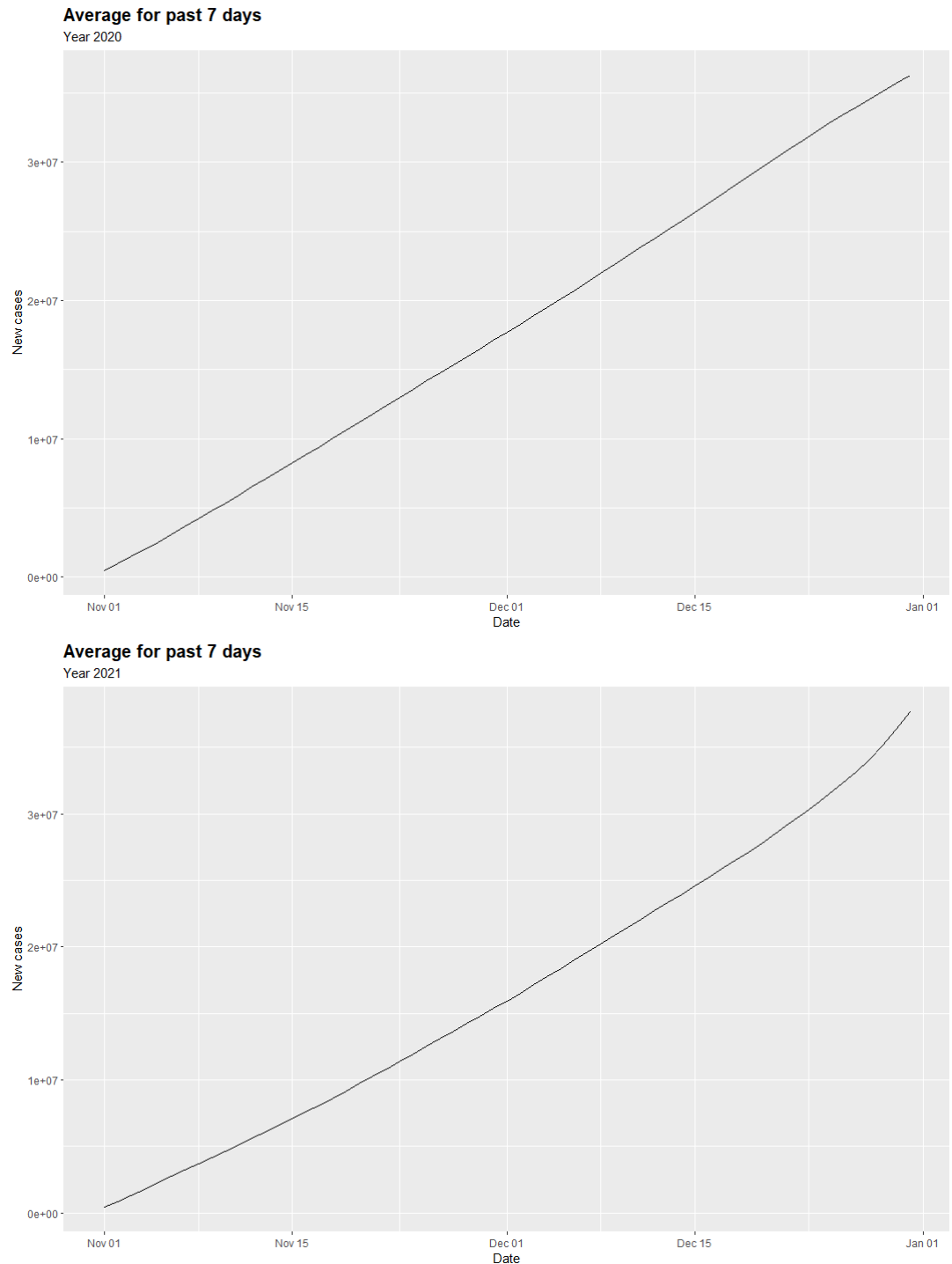


- 5) Biểu đồ thể hiện thu thập dữ liệu nhiễm bệnh tích lũy theo thời gian là 2 tháng của năm của tất cả quốc gia theo trung bình 7 ngày gần nhất

- Hiện thực trong R

```
1 figure <- function(y,name)
2 {
3   task <- subset(covid, (month == 11|month == 12) & year == y)
4   task <- aggregate(task$new_cases, list(task$date), FUN=sum)
5   colnames(task) <- c("date","new_cases")
6
7   task$rec <- 1
8   task$rec <- sum_run(x = task$rec,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
9   task$sum <- sum_run(x = task$new_cases,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
10  task$ave <- task$sum/task$rec
11
12  task$sum <- cumsum(task$ave)
13
14  graph <- ggplot(task, aes(x=date, y=sum)) + geom_line() + theme(plot.title = element_text(
15    size = 15, face = "bold")) + labs(title="Average for past 7 days", subtitle=name, x= "Date",
16    y= "New cases")
17  graph
18 }
19 figure(2020,"Year 2020")
20 figure(2021,"Year 2021")
```

- Kết quả:



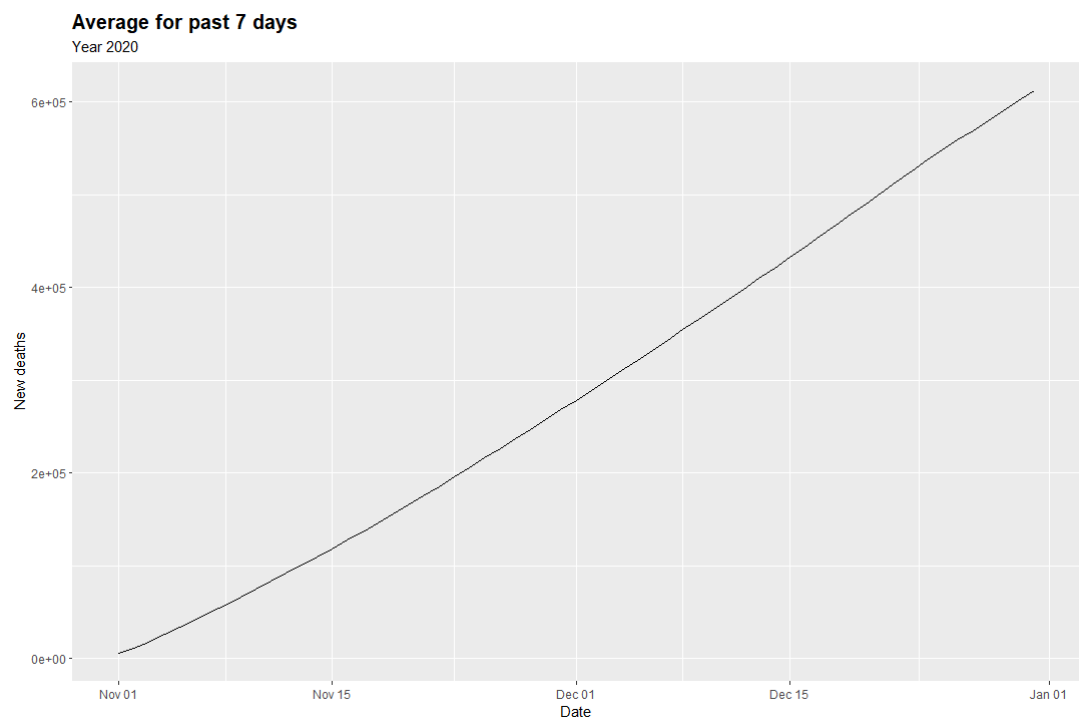
6) Biểu đồ thể hiện thu thập dữ liệu tử vong tích lũy theo thời gian là 2 tháng của năm của tất cả quốc gia theo trung bình 7 ngày gần nhất

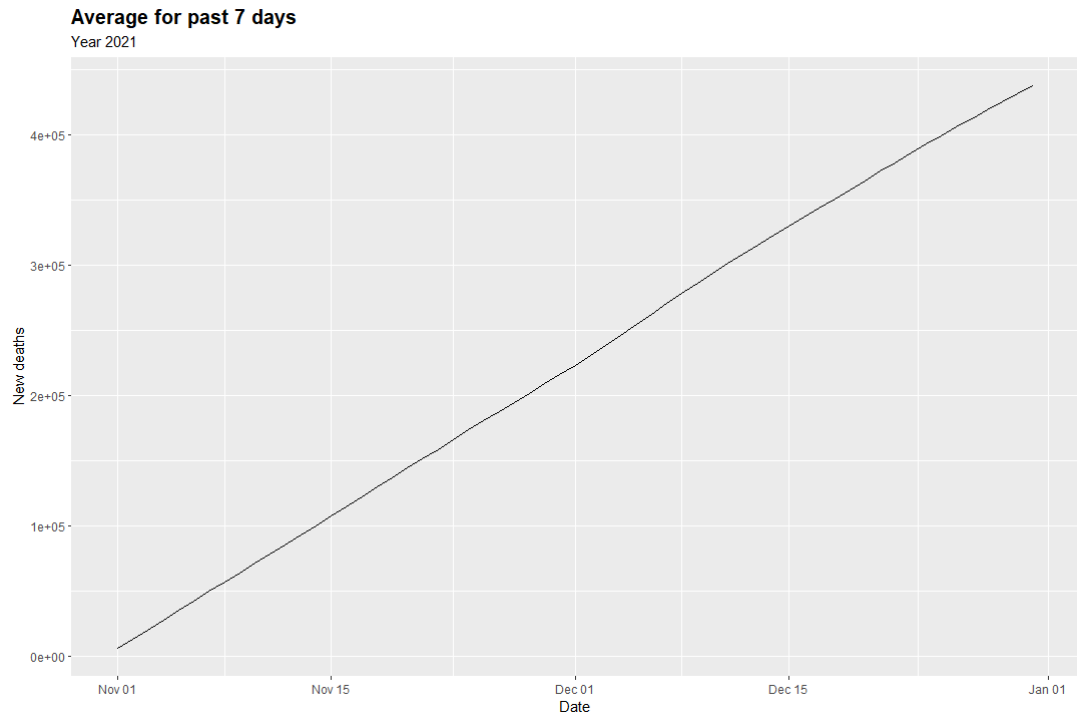
- Hiện thực trong R

```
1 figure <- function(y,name)
2 {
3   task <- subset(covid, (month == 11|month == 12) & year == y)
4   task <- aggregate(task$new_deaths, list(task$date), FUN=sum)
5   colnames(task) <- c("date","new_deaths")
6
7   task$rec <- 1
8   task$rec <- sum_run(x = task$rec,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
9   task$sum <- sum_run(x = task$new_deaths,k = 7,i = as.Date(task$date, format = "%m/%d/%Y"))
```

```
10 task$ave <- task$sum/task$rec
11
12 task$sum <- cumsum(task$ave)
13
14 graph <- ggplot(task, aes(x=date, y=sum)) + geom_line() + theme(plot.title = element_text(
15   size = 15, face = "bold")) + labs(title="Average for past 7 days", subtitle=name, x= "Date",
16   y= "New deaths")
17 graph
18 }
19
20 figure(2020, "Year 2020")
21 figure(2021, "Year 2021")
```

- Kết quả:





ix) Nhóm câu hỏi liên quan đến sự tương quan giữa nhiễm bệnh và tử vong

- Chuẩn bị dữ liệu cho toàn bộ phần ix (phần code dùng chung cho cả phần ix)

```
1 library(ggplot2)
2 library(readr)
3 library(runner)
4 library(dplyr)
5 library(zoo)
6
7 covid <- read_csv("C:/Users/Asus/Documents/covidData.csv")
8 covid$new_cases <- abs(covid$new_cases)
9 covid$new_deaths <- abs(covid$new_deaths)
10 covid$date <- as.Date(covid$date, "%m/%d/%Y")
11 covid$month <- as.numeric(format(covid$date, '%m'))
12 covid$year <- as.numeric(format(covid$date, '%Y'))
13 covid <- subset(covid, !is.na(continent))
14 covid$new_cases[is.na(covid$new_cases)] <- 0
15 covid$new_deaths[is.na(covid$new_deaths)] <- 0
16 year <- unique(format(covid$date, format= "%Y"))
17 lcountry <- c("Brazil", "Chile", "Venezuela")
18 lmonth <- c("02", "03", "06", "10")
```

- Vẽ biểu đồ thể hiện phần trăm giữa nhiễm bệnh tích lũy trên tổng nhiễm bệnh và phần trăm tử vong tích lũy trên tổng số tử vong cho từng quốc gia theo thời gian. Vẽ 2 đường trên cùng biểu đồ

Trên từng quốc gia hãy vẽ biểu đồ thể hiện trục Ox là nhiễm bệnh, trục Oy là tử vong. Hãy lấy 4 tháng theo 4 ký số mã đề thể hiện. Nếu ký số là 0 thì lấy tháng là 10.

- Hiện thực trong R:

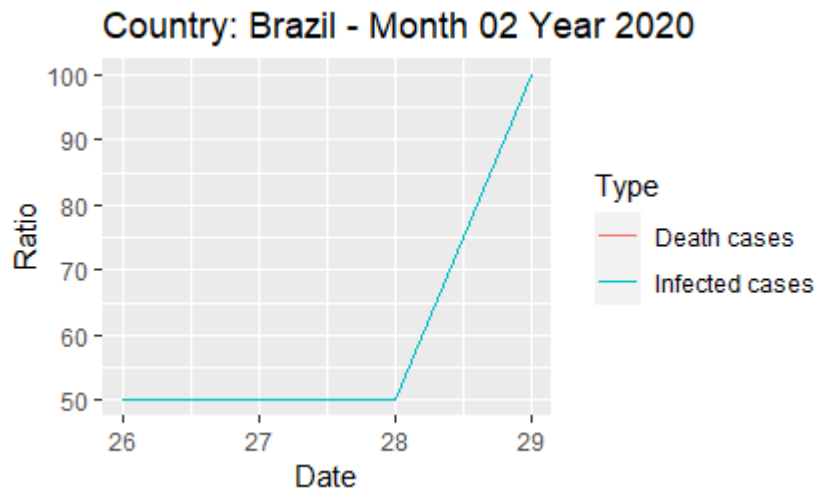
```
1 for (i in year){
2   for (k in lmonth){
3     for (c in lcountry){
4       l1 <- NULL
5       l2 <- NULL
6       tab <- subset(covid, format(covid$date, format= "%Y") == i & format(covid$date, format=
7         "%m") == k & covid$location == c)
8       if (nrow(tab) != 0){
9         l1 <- tab
```

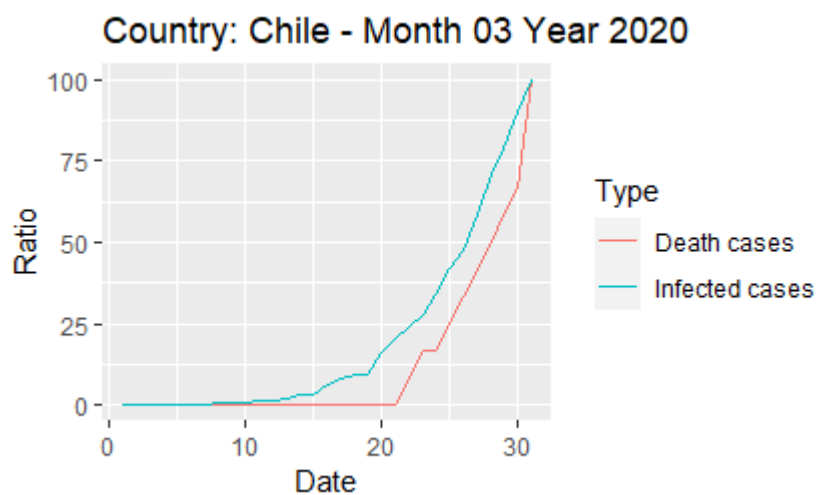
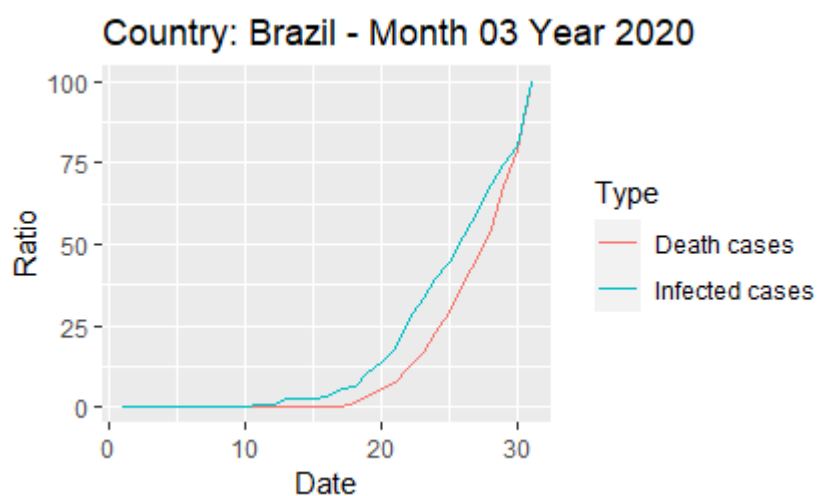
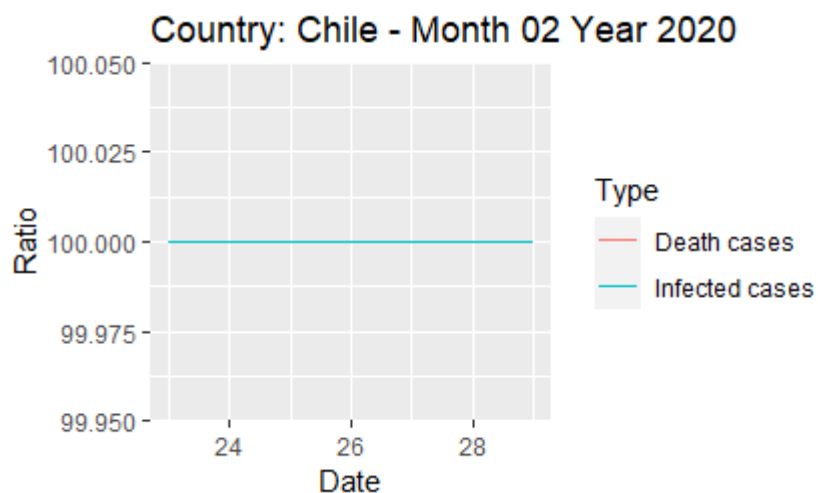
```

9      l1$day <- format(l1$date,"%d")
10     l1$date <- "Infected cases"
11     l1$new <- l1$new_cases
12
13     l2 <- tab
14     l2$day <- format(l2$date,"%d")
15     l2$date <- "Death cases"
16
17     l2$new <- l2$new_deaths
18     keeps <- c("day","date","new")
19     l1 <- subset(l1, select = keeps)
20     l2 <- subset(l2, select = keeps)
21     df <- rbind(l1,l2)
22     print(ggplot(df,aes(x=as.integer(day), y=new, fill = date)) +
23           geom_point(aes(col=date, size = new)) +
24           geom_smooth(formula = y ~ x, method = "lm") +
25           labs(title = paste("Country:",c,"- Month",k,"Year",i), x = "Date", y = "Number
of cases",fill="Type",size = "Number of cases") +
26           guides(col = FALSE))
27     cat(paste("Country:",c,"- Month",k,"Year",i))
28     cat("\n")
29     cat("Correlation coefficients of infected cases: ")
30     cat(cor(as.integer(df$day[df$date=="Infected cases"]),df$new[df$date=="Infected cases"
]))
31     cat("\n")
32     cat("Correlation coefficients of death cases: ")
33     cat(cor(as.integer(df$day[df$date=="Death cases"]),df$new[df$date=="Death cases"]))
34     cat("\n")
35   }
36 }
37 }
38 }

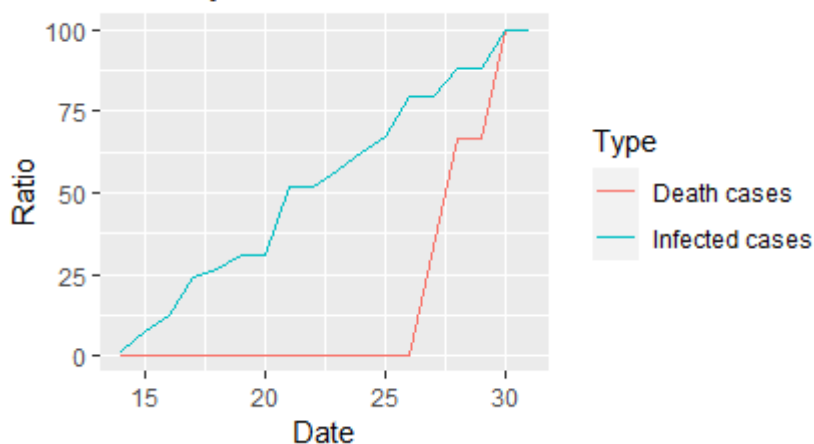
```

- Kết quả

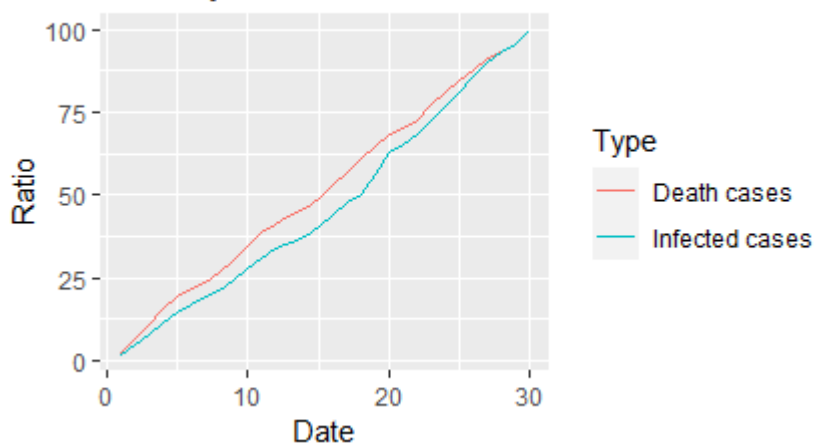




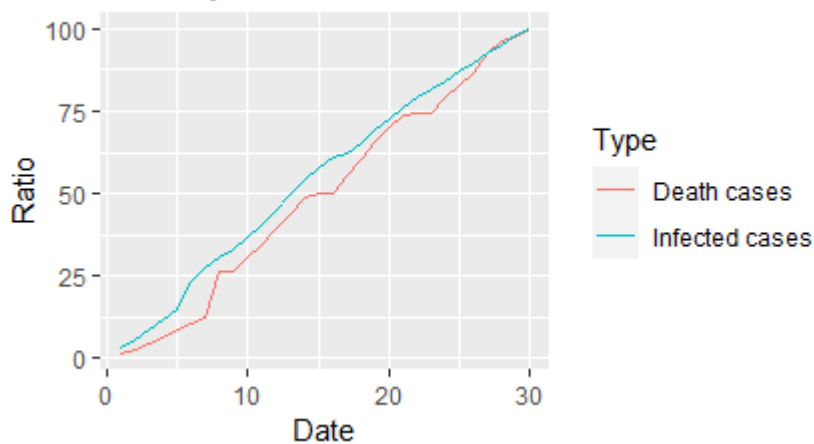
Country: Venezuela - Month 03 Year 2020



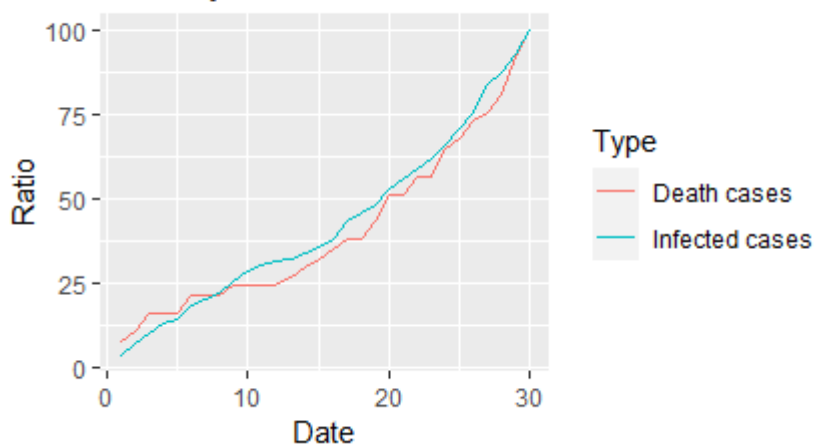
Country: Brazil - Month 06 Year 2020



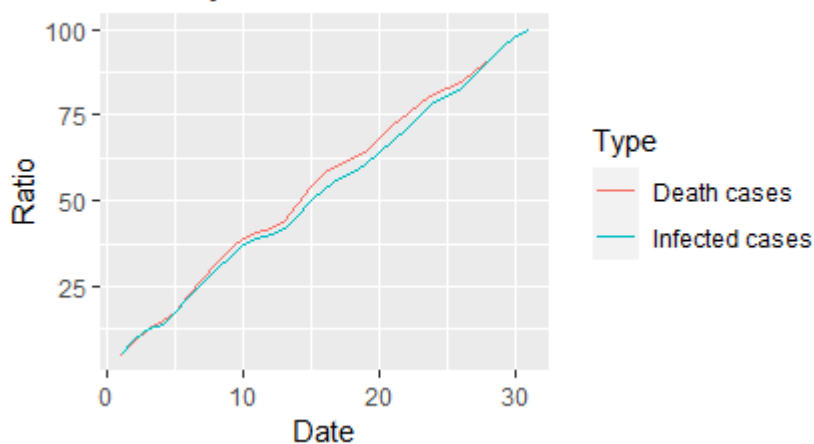
Country: Chile - Month 06 Year 2020



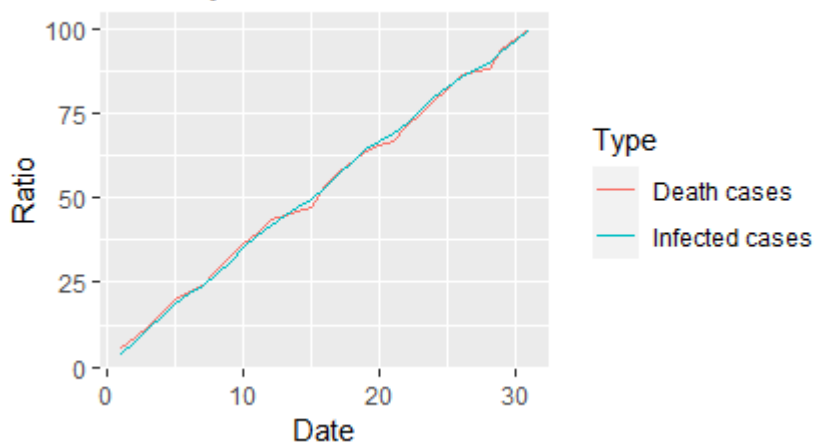
Country: Venezuela - Month 06 Year 2020



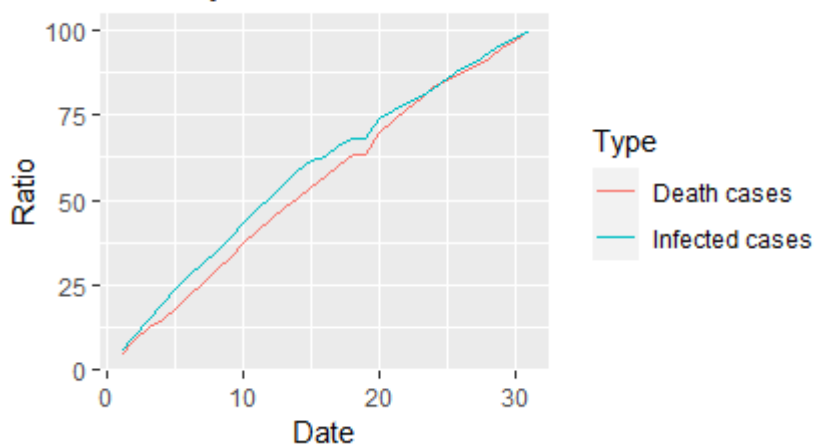
Country: Brazil - Month 10 Year 2020



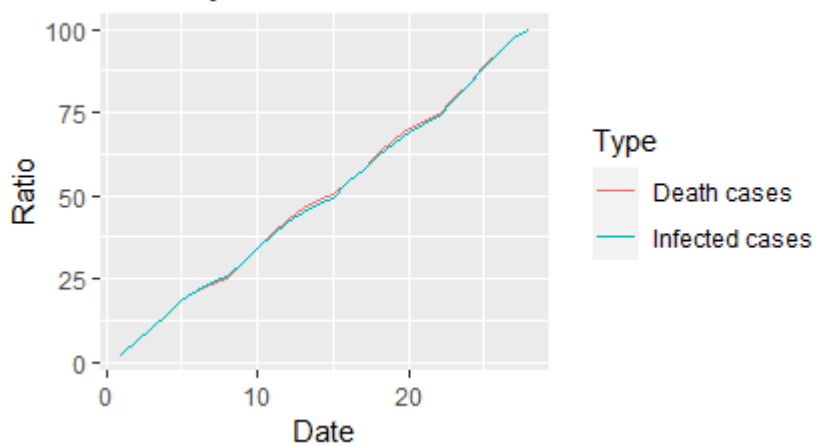
Country: Chile - Month 10 Year 2020



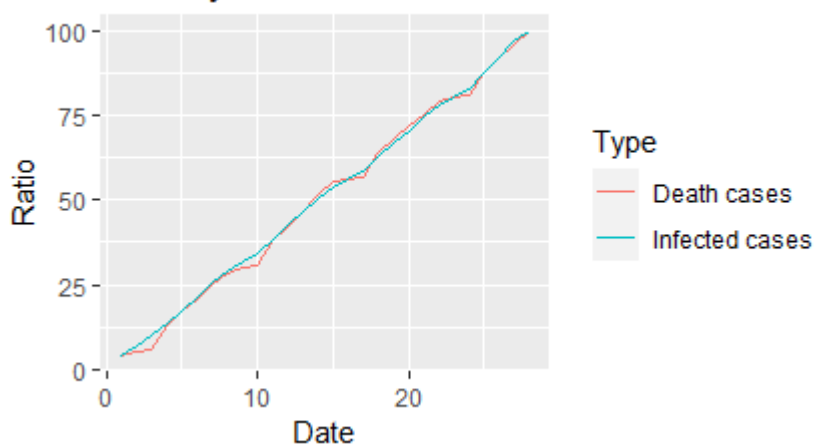
Country: Venezuela - Month 10 Year 2020



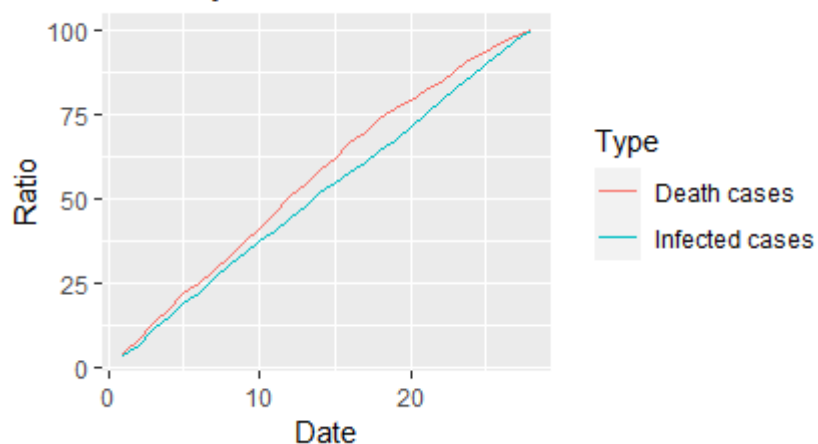
Country: Brazil - Month 02 Year 2021



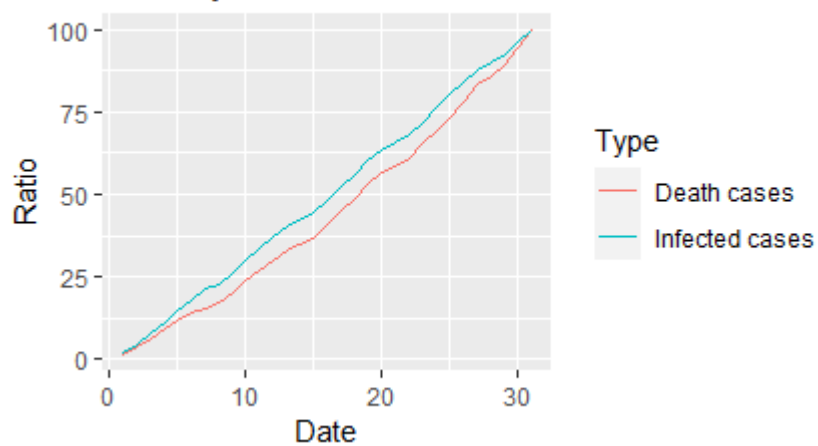
Country: Chile - Month 02 Year 2021



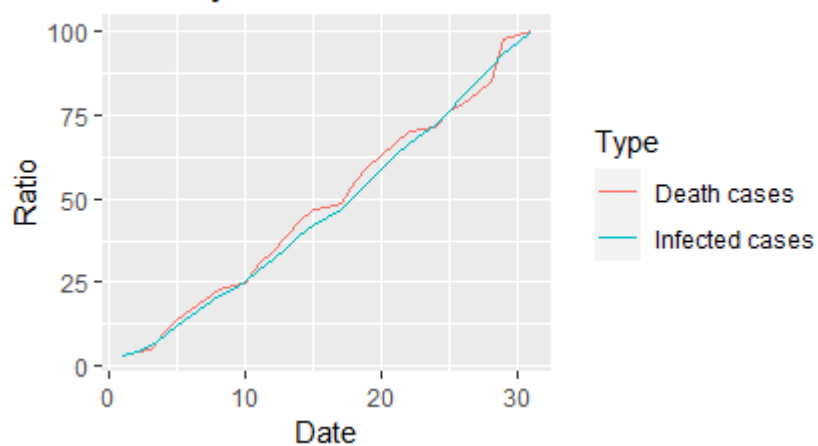
Country: Venezuela - Month 02 Year 2021



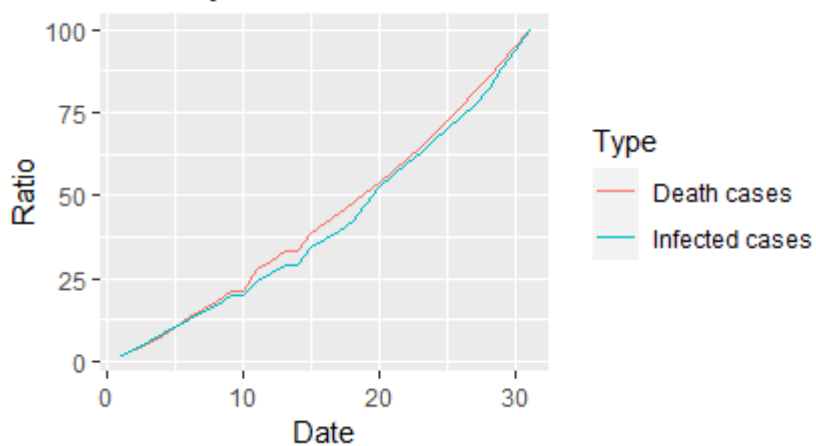
Country: Brazil - Month 03 Year 2021



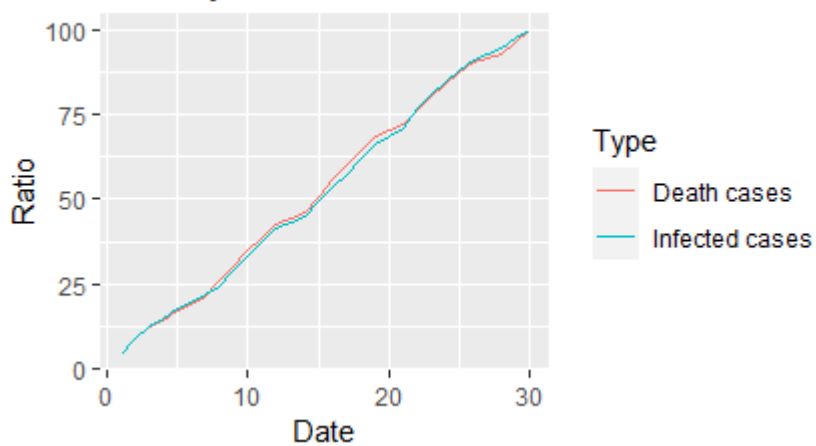
Country: Chile - Month 03 Year 2021



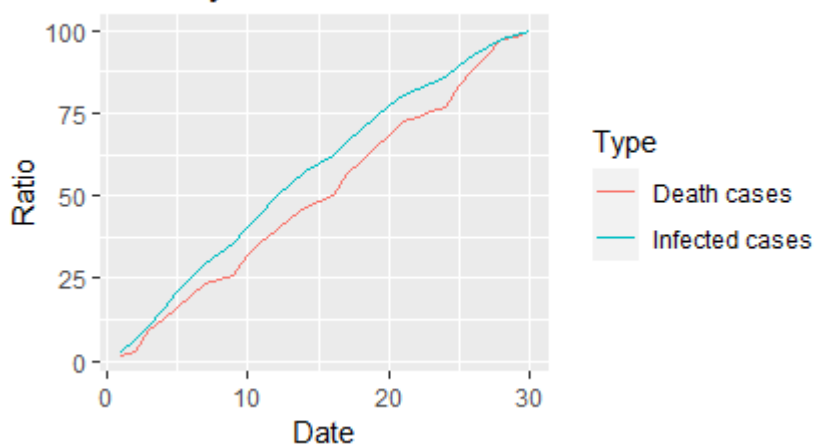
Country: Venezuela - Month 03 Year 2021



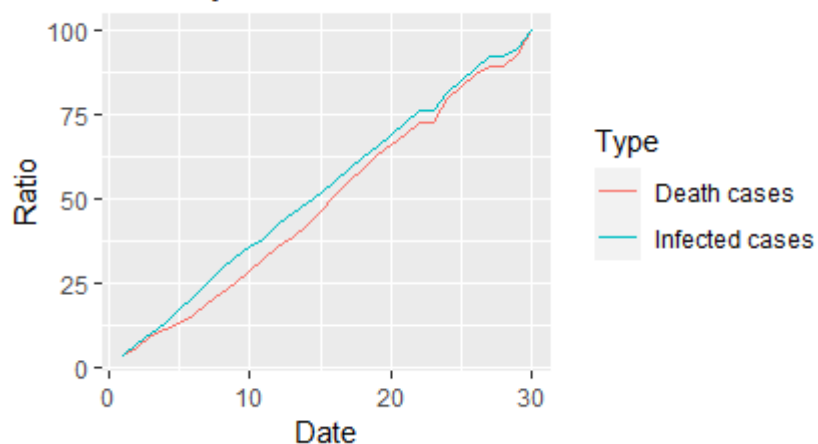
Country: Brazil - Month 06 Year 2021



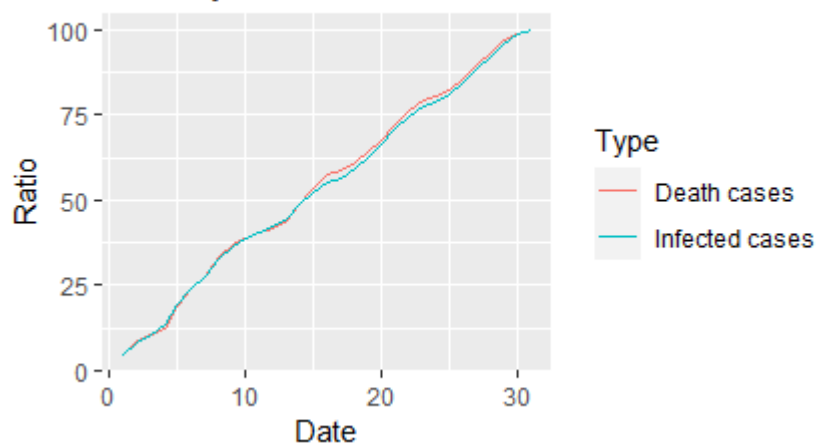
Country: Chile - Month 06 Year 2021



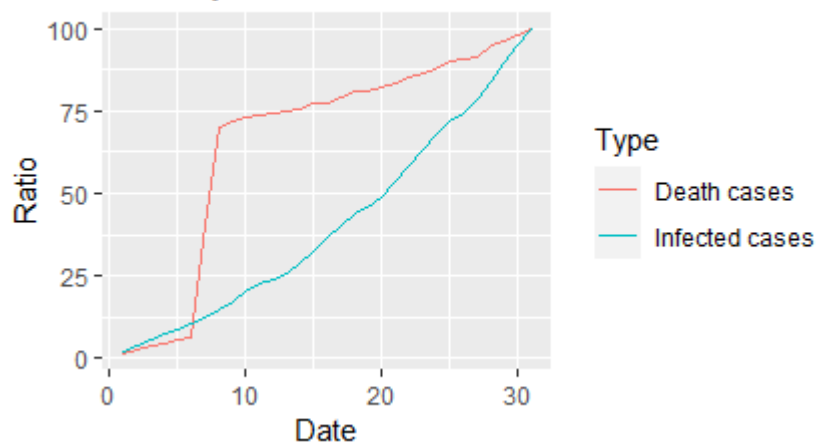
Country: Venezuela - Month 06 Year 2021



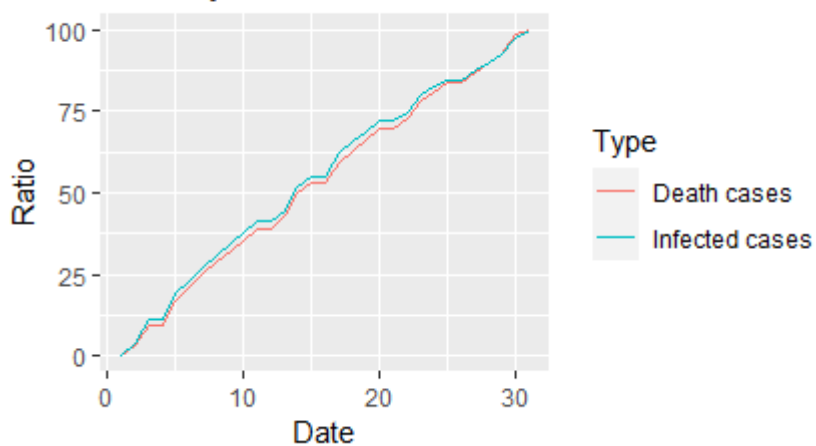
Country: Brazil - Month 10 Year 2021



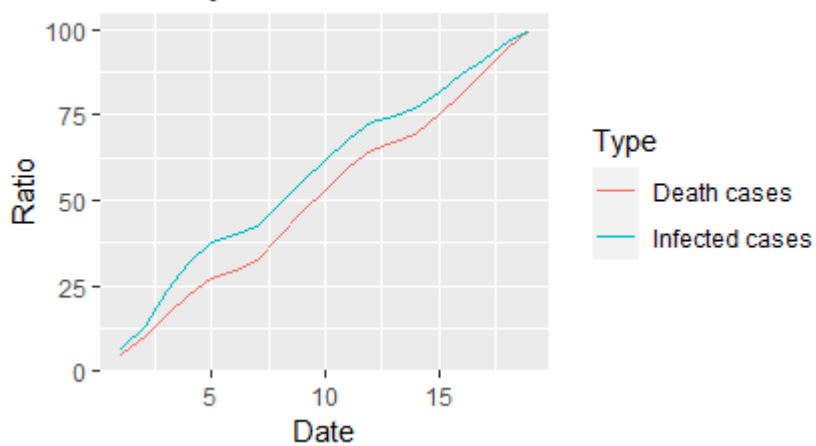
Country: Chile - Month 10 Year 2021



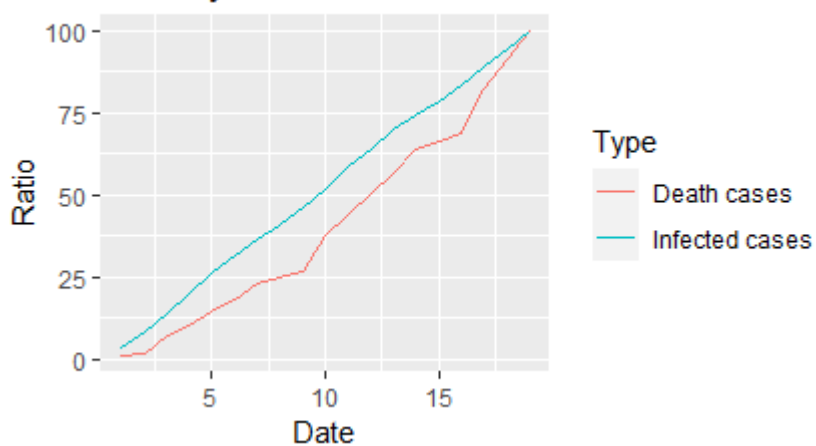
Country: Venezuela - Month 10 Year 2021

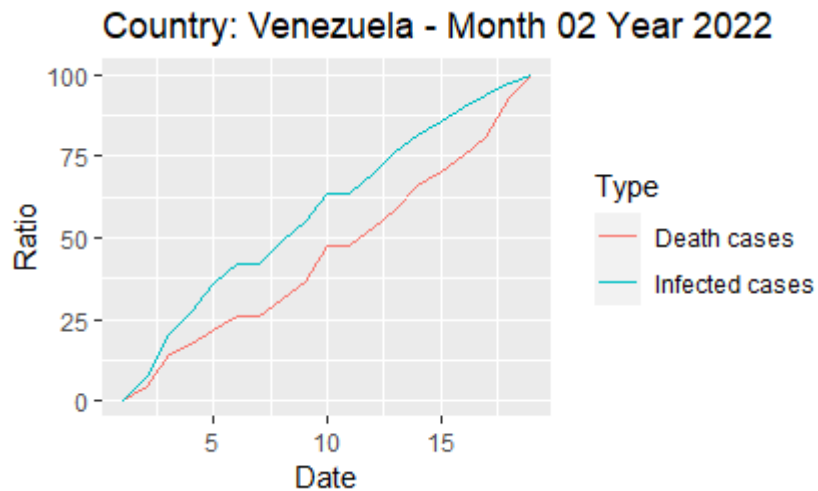


Country: Brazil - Month 02 Year 2022



Country: Chile - Month 02 Year 2022



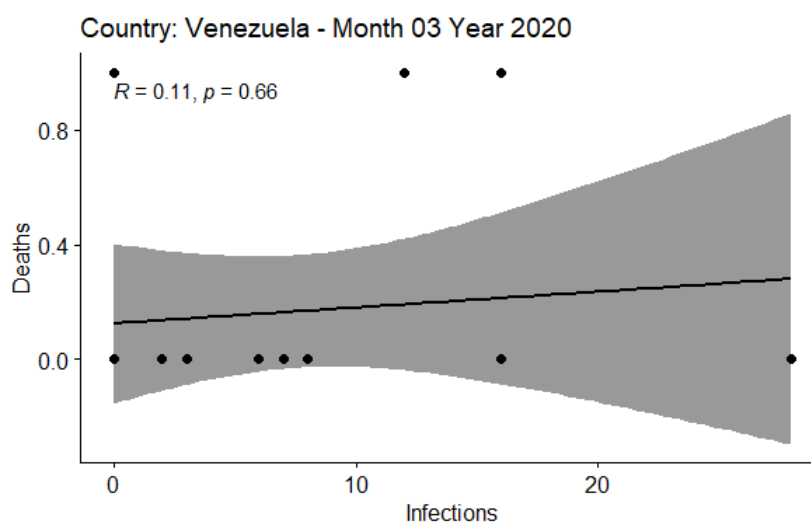
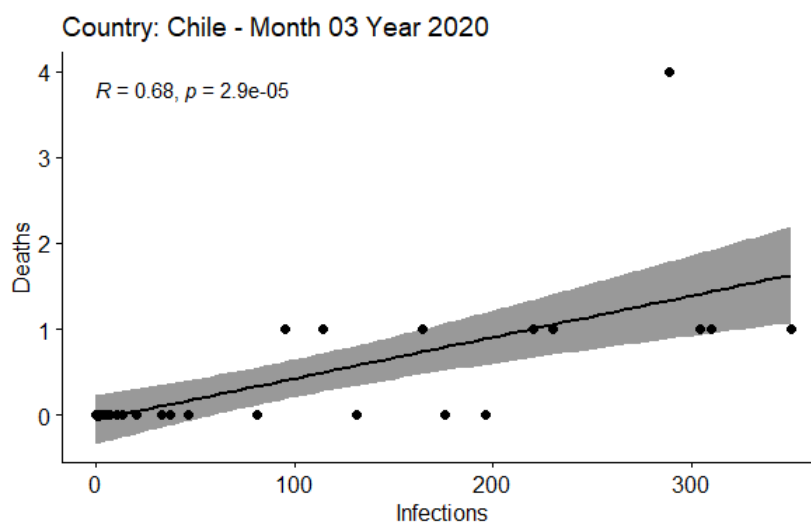
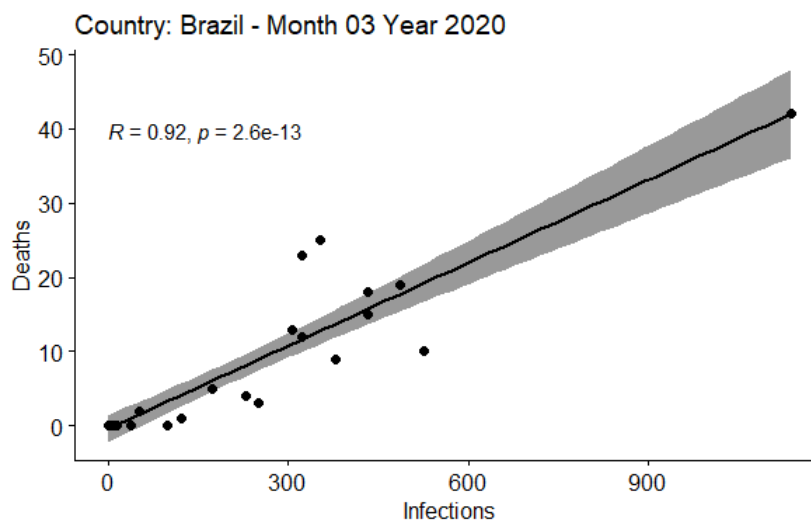


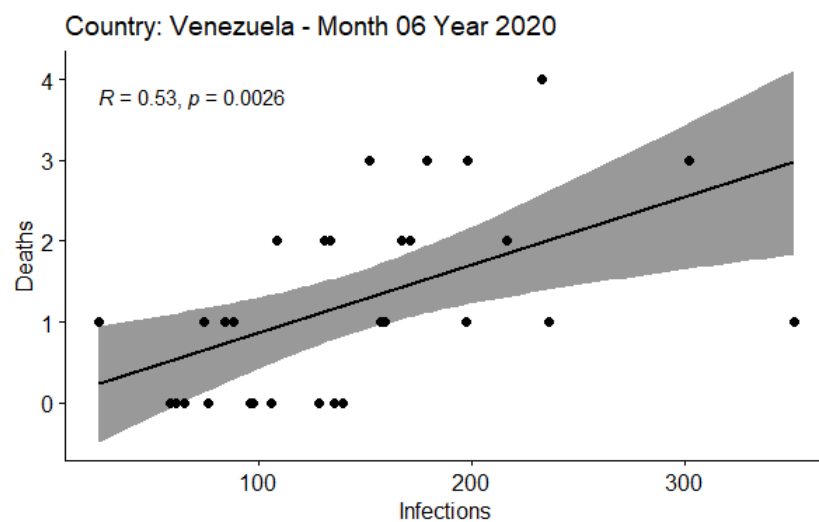
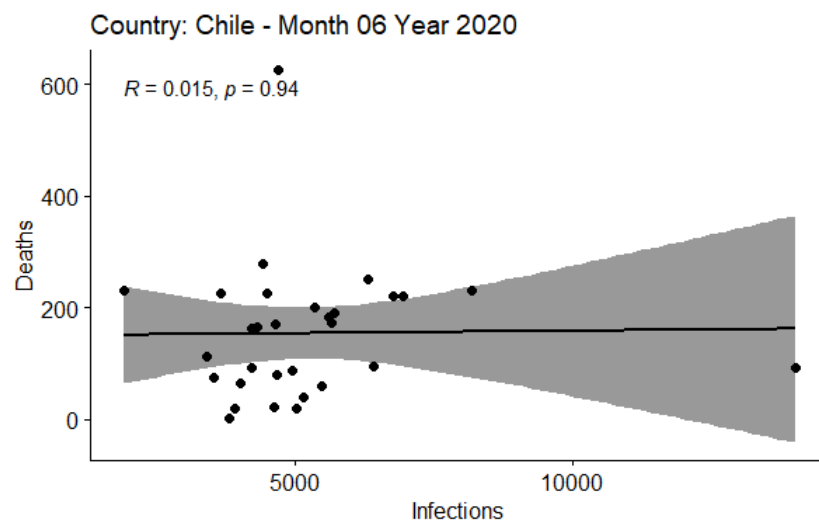
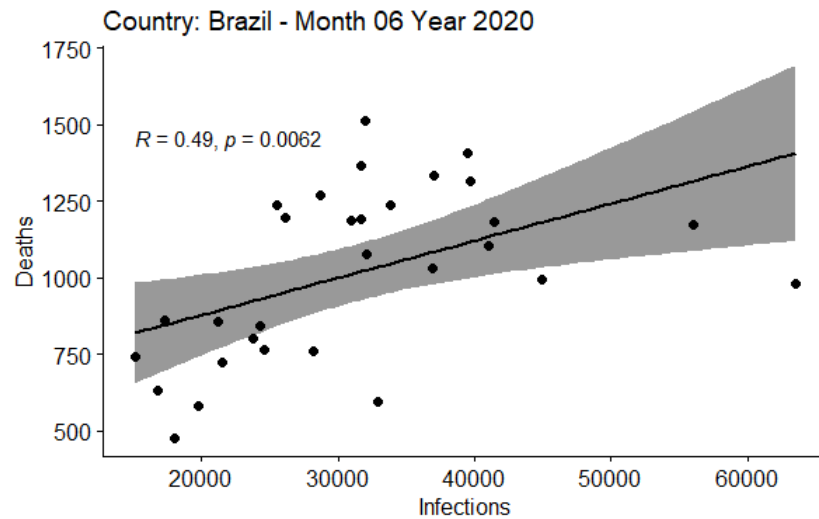
2) Xét tương quan trong mỗi tháng,

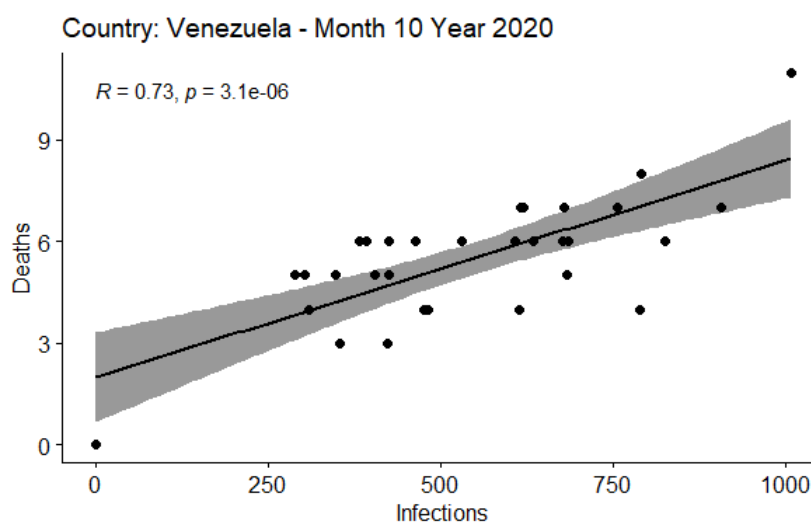
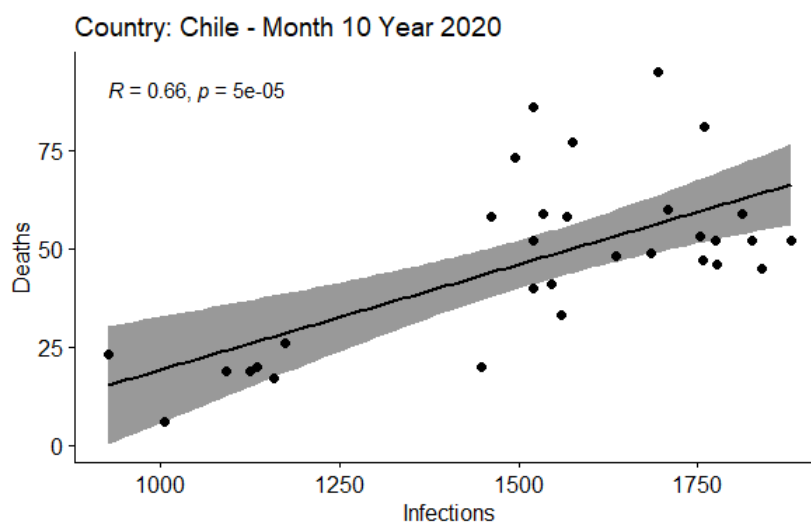
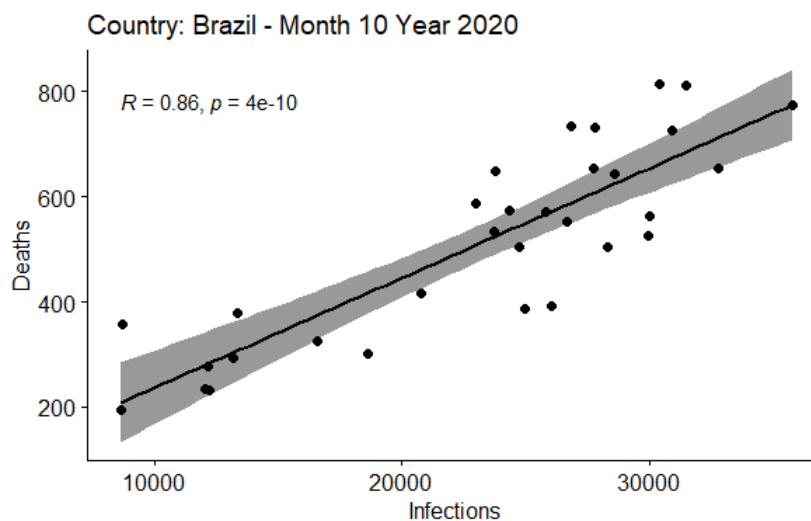
- Hiện thực trong R:

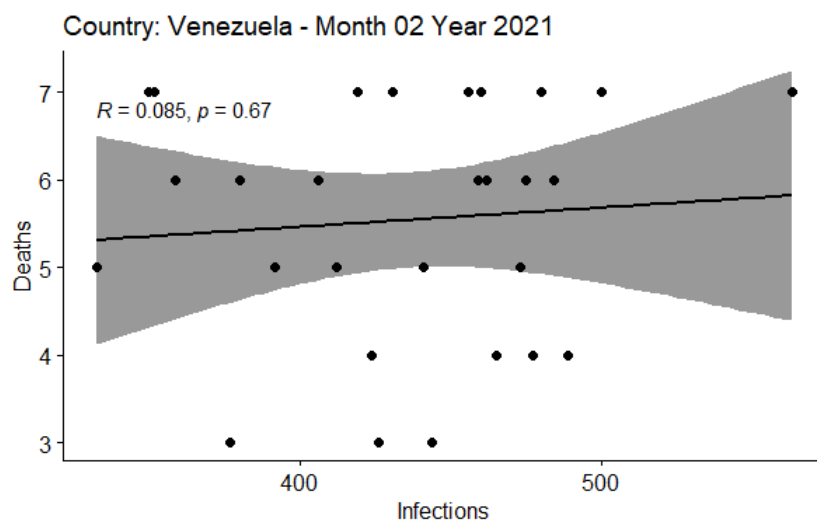
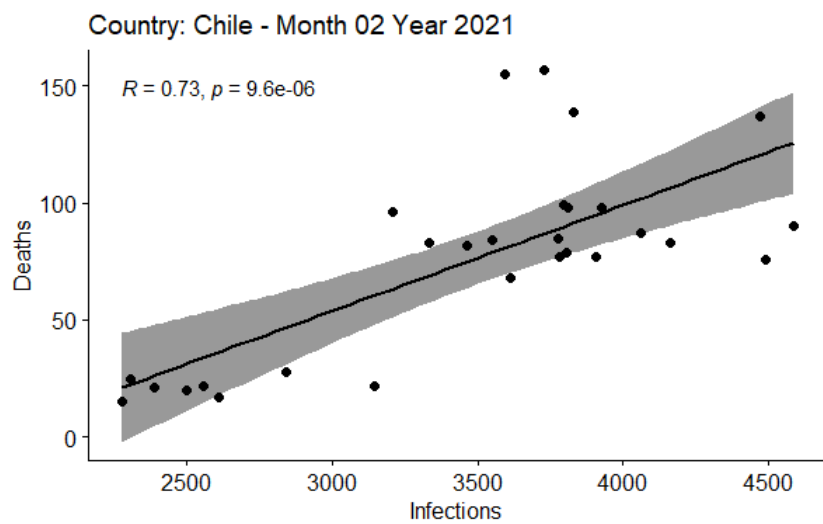
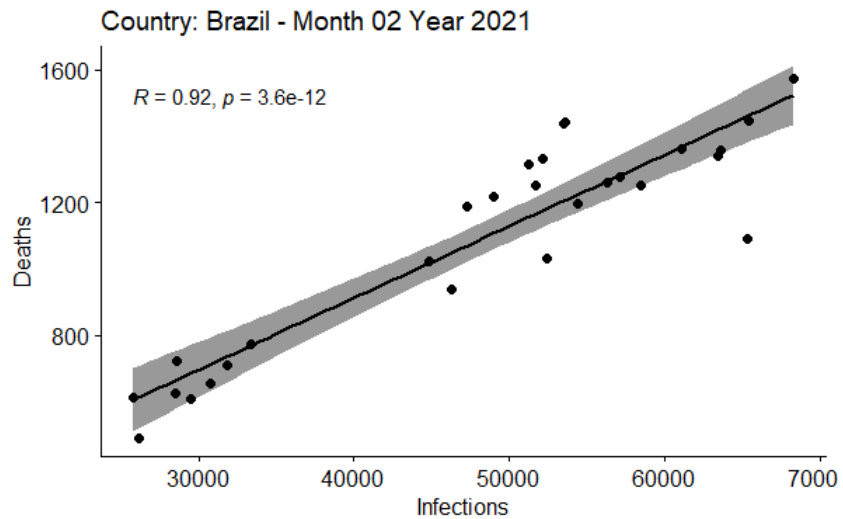
```
1 library(tidyverse)
2 library(dplyr)
3 library(datasets)
4 library("ggplot2")
5
6 covid = read.csv("owid-covid-data.csv", header = TRUE)
7
8 covid$new_cases <- abs(covid$new_cases)
9 covid$new_deaths <- abs(covid$new_deaths)
10
11 covid$date <- as.Date(covid$date, "%m/%d/%Y")
12 covid$month <- as.numeric(format(covid$date, "%m"))
13 covid$year <- as.numeric(format(covid$date, "%Y"))
14
15 covid <- subset(covid, !is.na(continent))
16 covid$new_cases[is.na(covid$new_cases)] <- 0
17 covid$new_deaths[is.na(covid$new_deaths)] <- 0
18 year <- unique(format(covid$date, format = "%Y"))
19 lcountry <- c("Brazil", "Chile", "Venezuela")
20 lmonth <- c("02", "03", "06", "10")
21
22
23
24 for (i in year){
25   for (k in lmonth){
26     for (c in lcountry){
27       temp = subset(covid, format(covid$date, format = "%Y") == i & format(covid$date, format =
28         "%m") == k & covid$location == c)
29       print(ggscatter(temp, x = "new_cases", y = "new_deaths",
30         add = "reg.line", conf.int = TRUE,
31         cor.coef = TRUE, cor.method = "pearson",
32         xlab = "Infections", ylab = "Deaths") + labs(title = paste("Country:", c,
33           "- Month", k, "Year", i)))
34     }
35   }
36 }
```

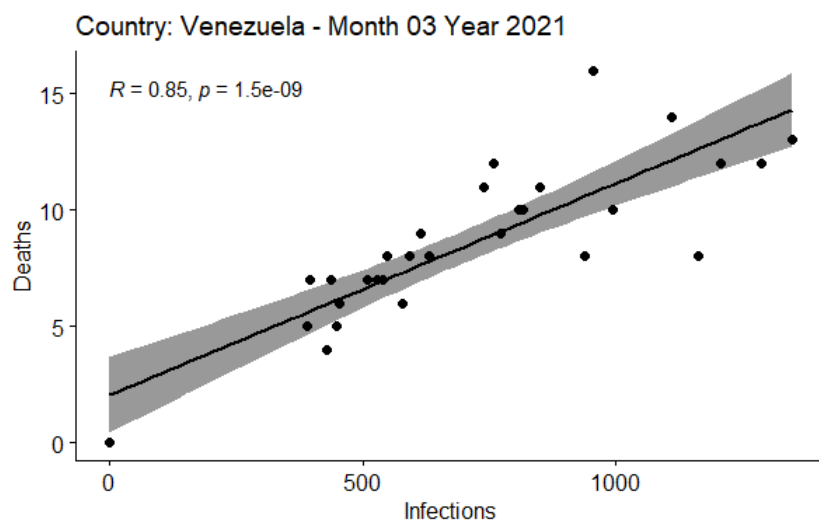
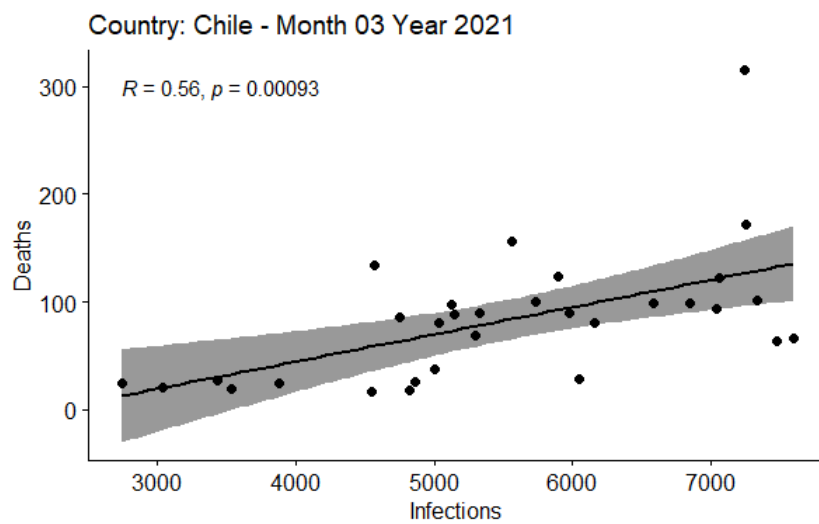
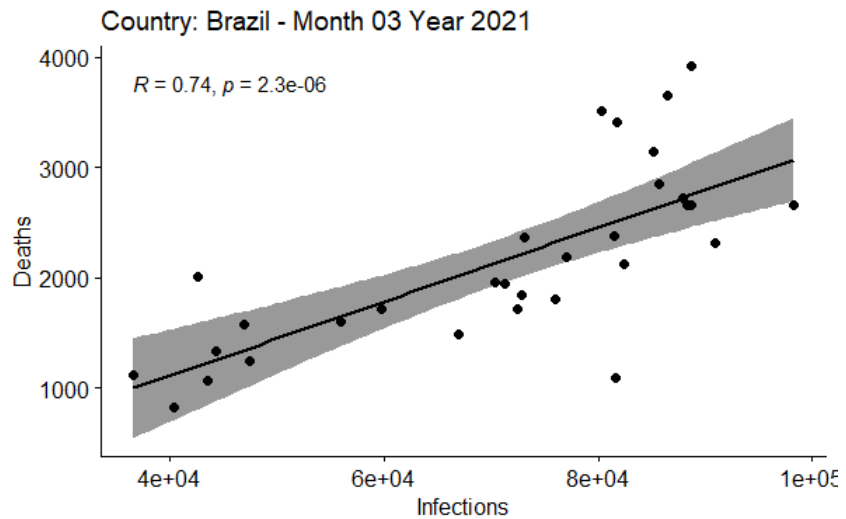
- Kết quả

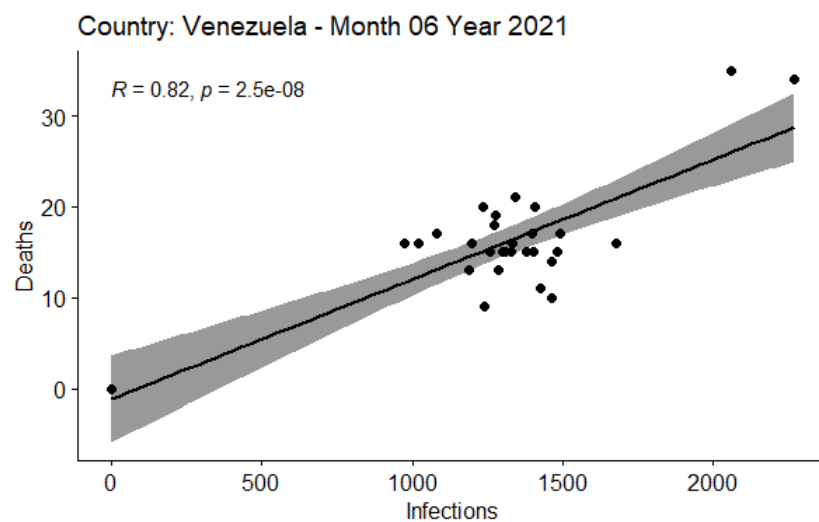
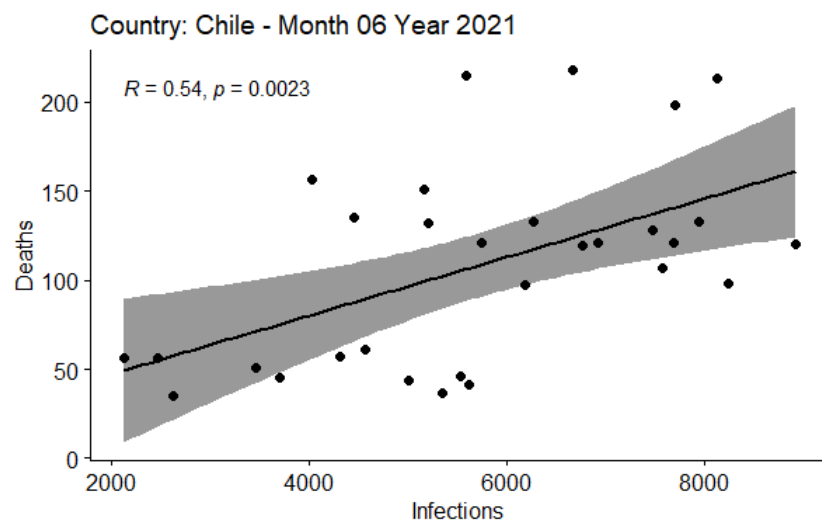
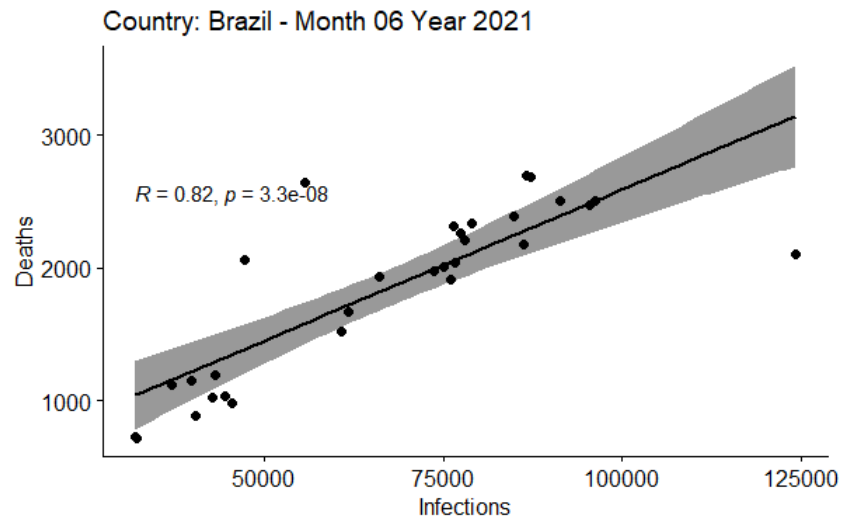


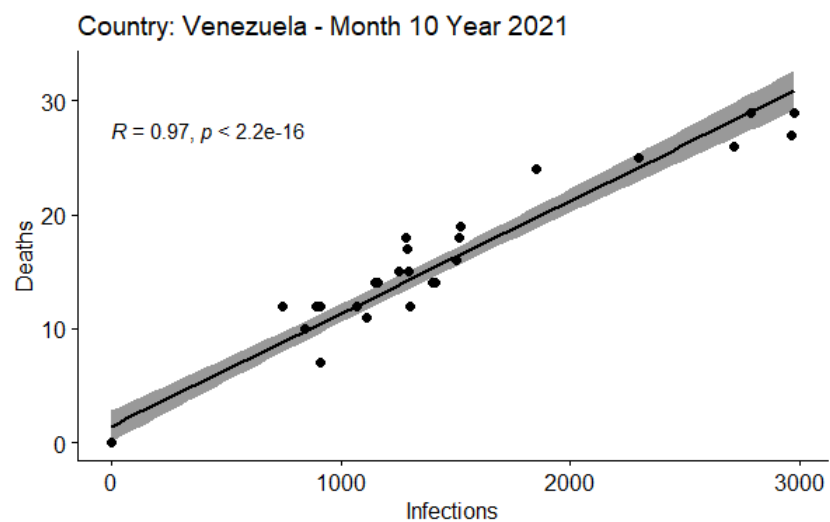
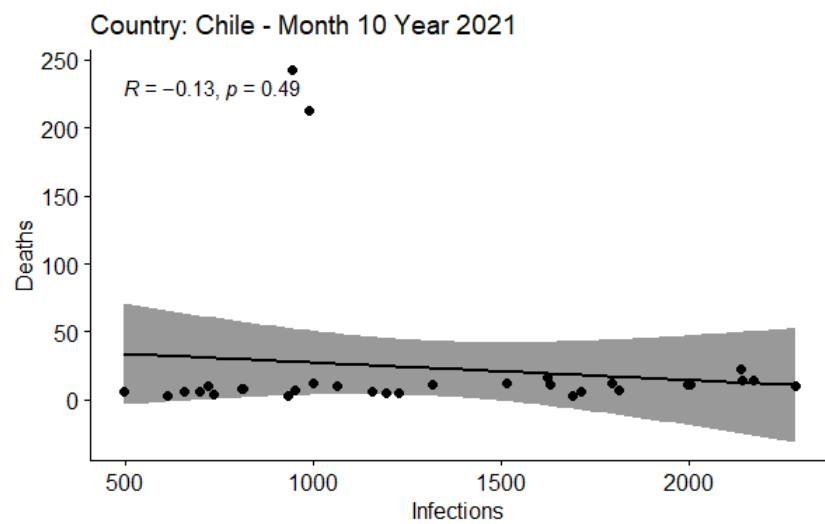
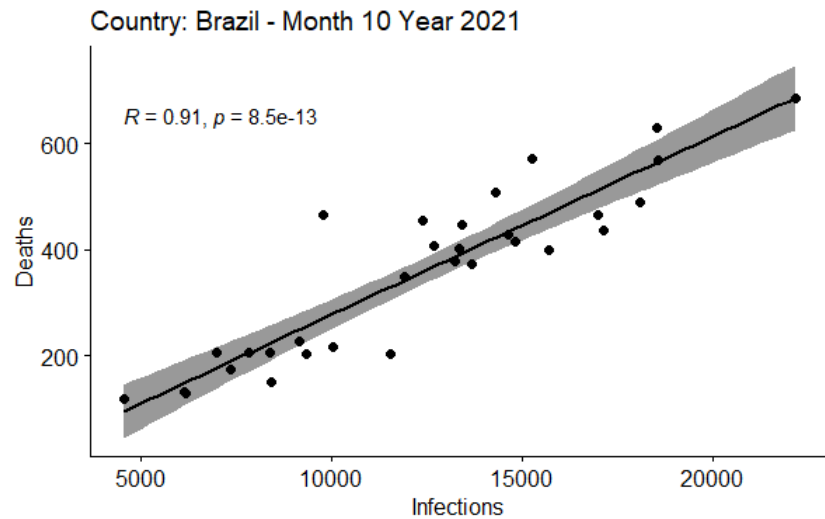


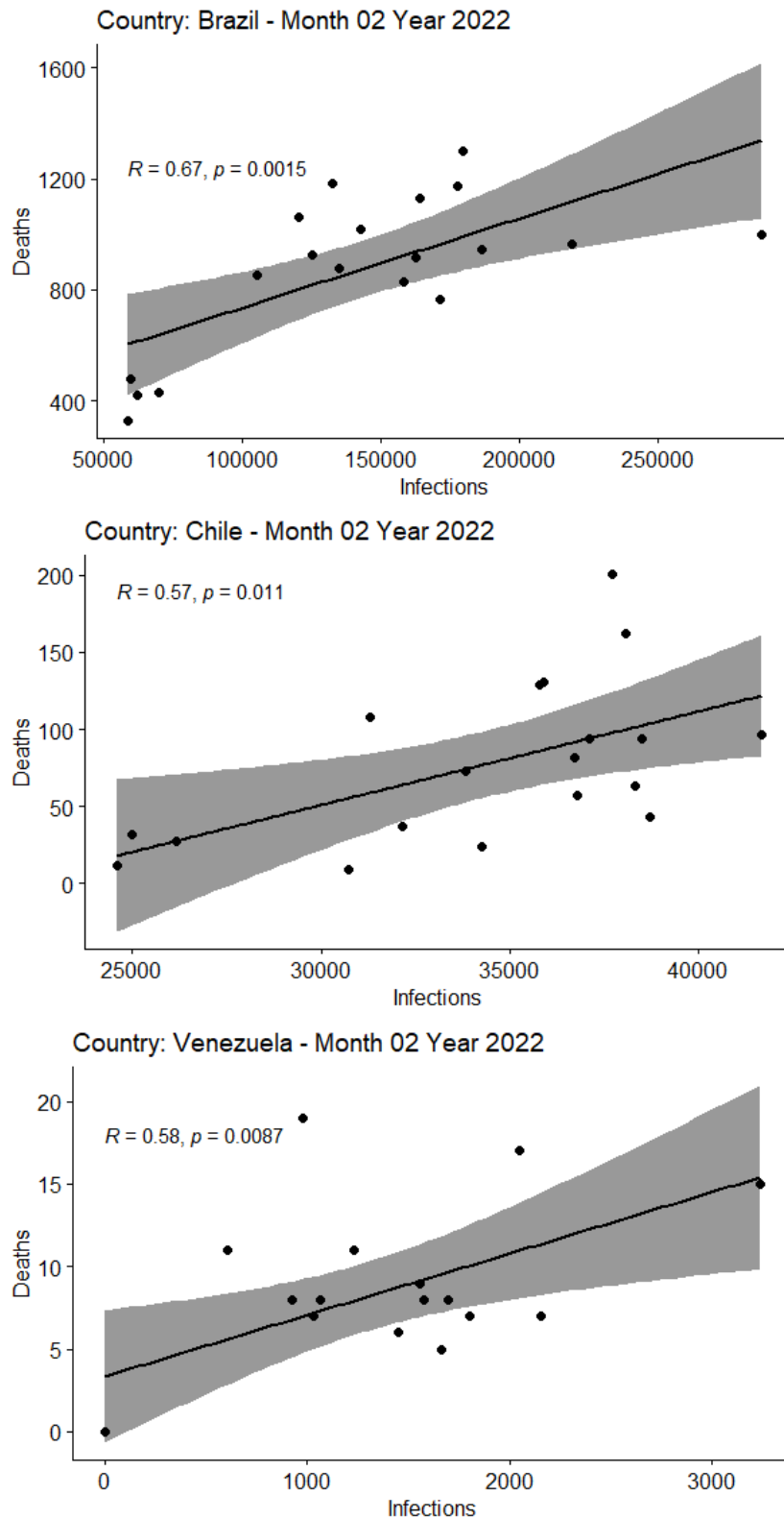












3) Xét tương quan trong mỗi tháng theo trung bình 7 ngày gần nhất

• Hiện thực trong R:

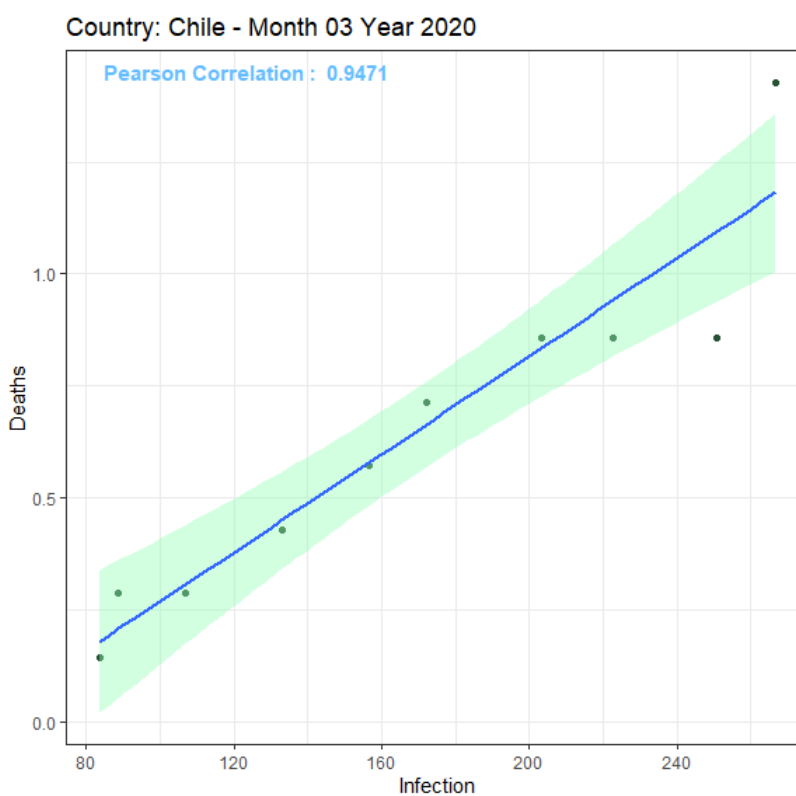
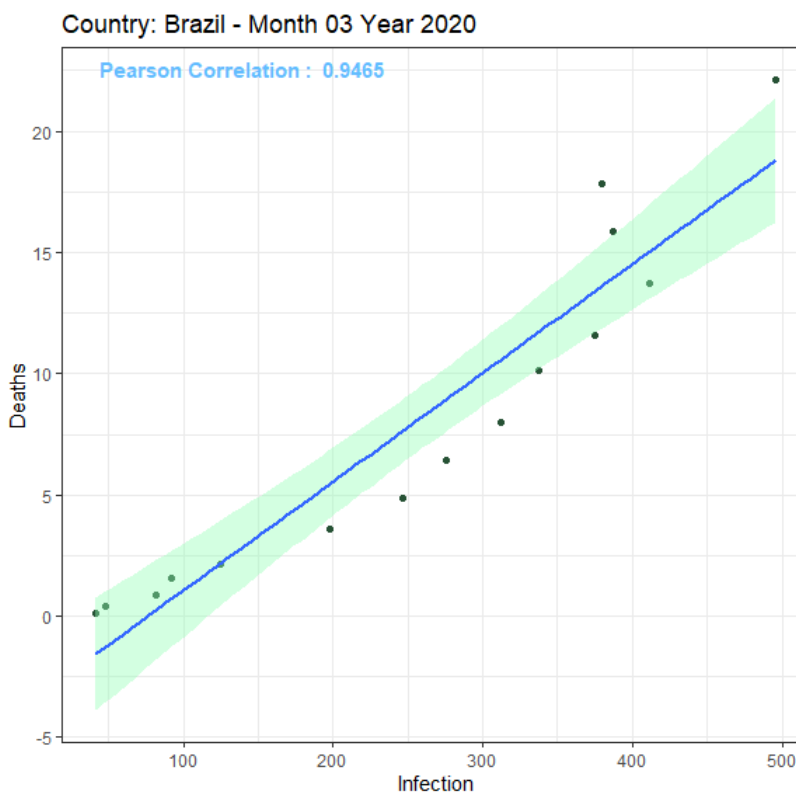
```
1 library(tidyverse)
2 library(dplyr)
3 library(datasets)
4 library("ggplot2")
```

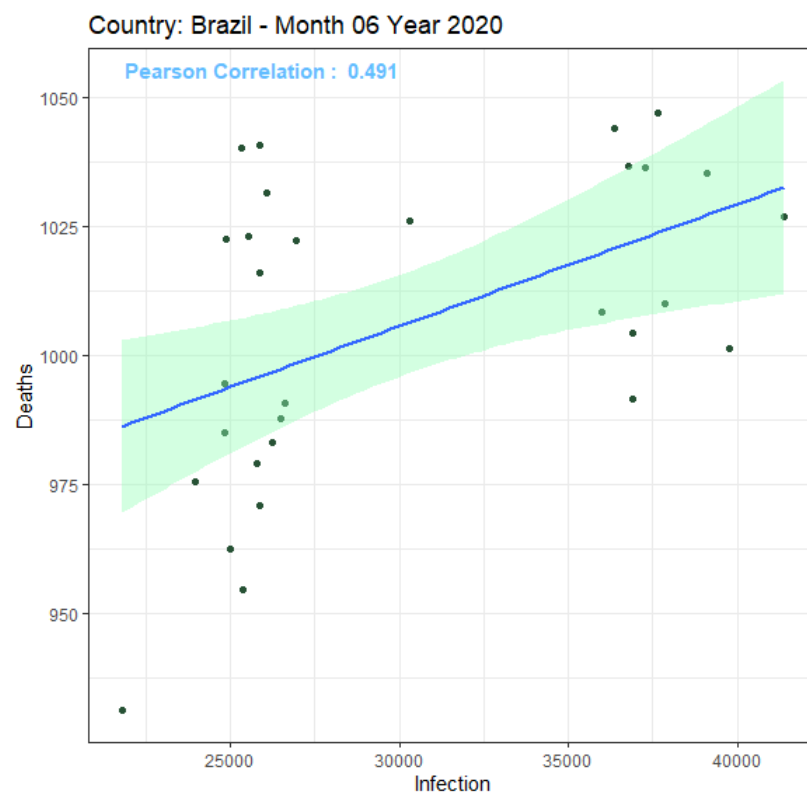
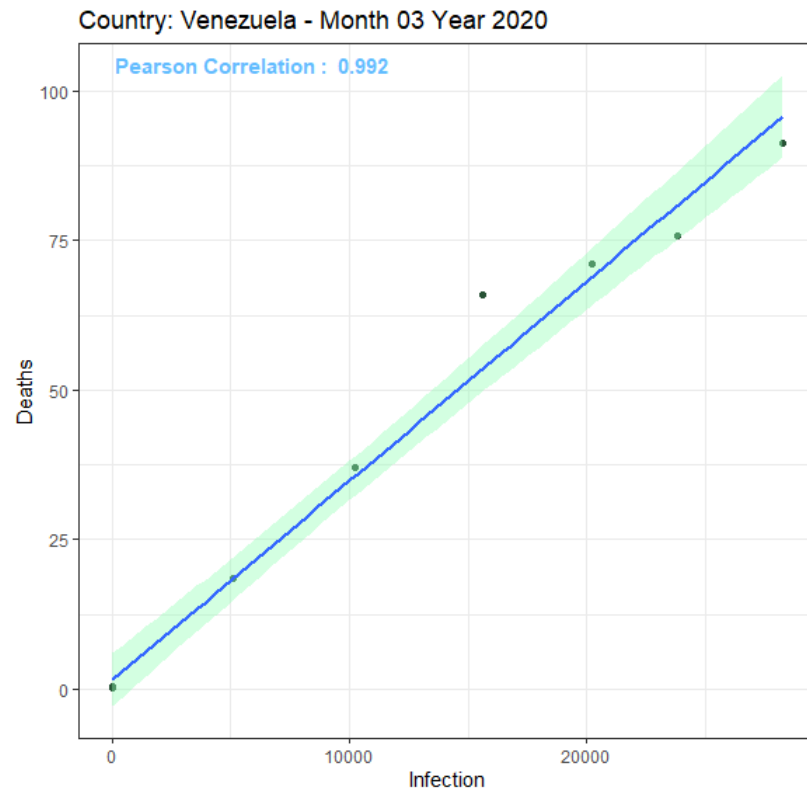
```

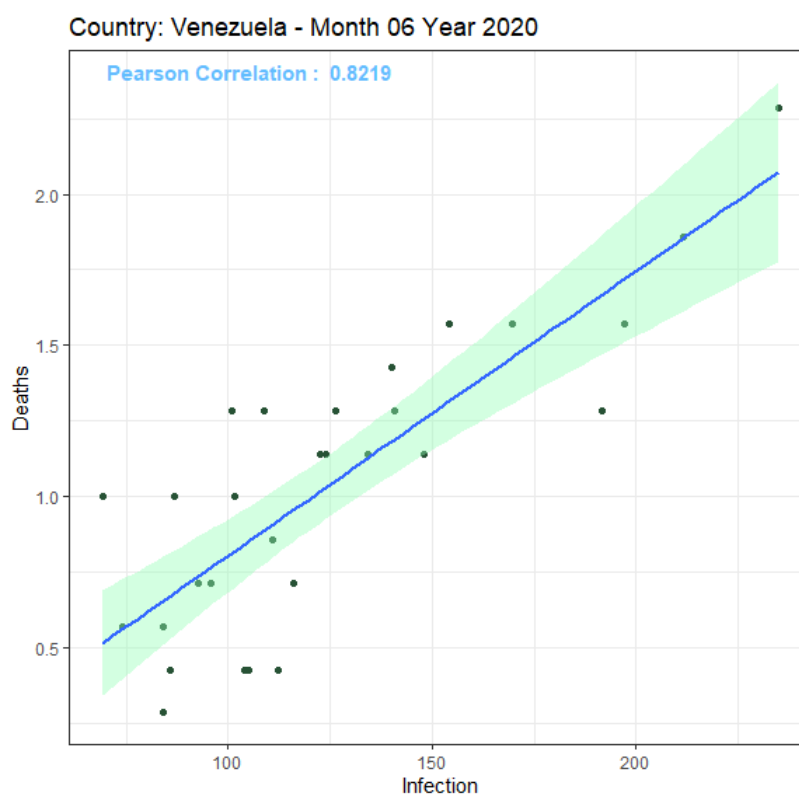
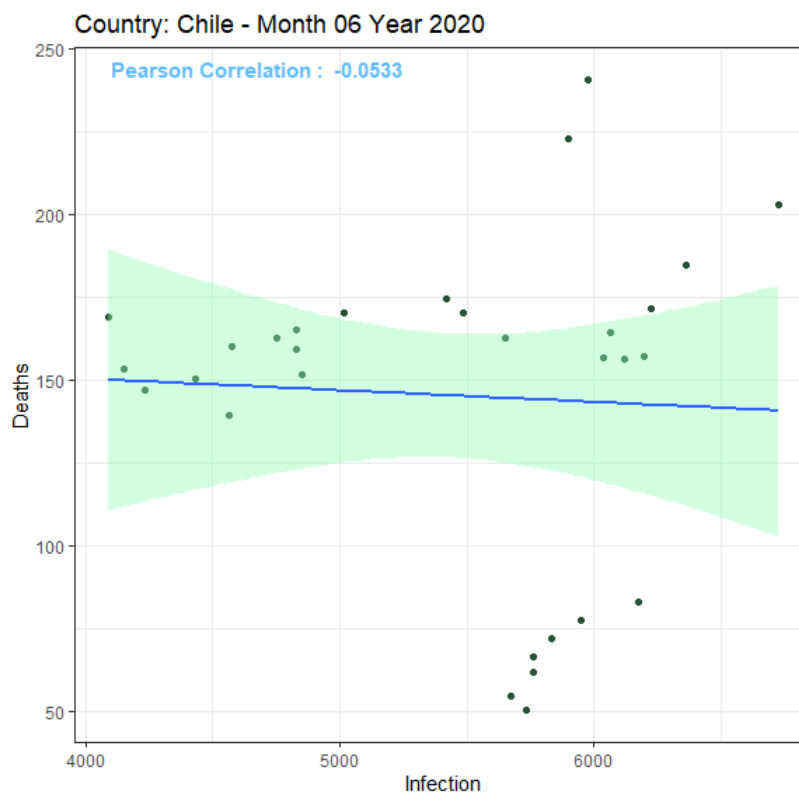
5 library(grid)
6 library(gridExtra)
7 covid = read.csv("C:/Users/Admin/Downloads/covidData.csv")
8
9 covid$new_cases <- abs(covid$new_cases)
10 covid$new_deaths <- abs(covid$new_deaths)
11
12 covid$date <- as.Date(covid$date, "%m/%d/%Y")
13 covid$month <- as.numeric(format(covid$date, '%m'))
14 covid$year <- as.numeric(format(covid$date, '%Y'))
15
16 covid <- subset(covid, !is.na(continent))
17 covid$new_cases[is.na(covid$new_cases)] <- 0
18 covid$new_deaths[is.na(covid$new_deaths)] <- 0
19 year <- unique(format(covid$date, format= "%Y"))
20 lcountry <- c("Brazil", "Chile", "Venezuela")
21 lmonth <- c("02", "03", "06", "10")
22
23 covid[["Infection"]] = NA
24 covid[["Deaths"]] = NA
25 covid = covid %>% filter(iso_code== "BRA" | iso_code == "CHL" | iso_code == "VEN")
26 for(j in 1:nrow(covid)) {
27   if(j <= 7) {
28     covid$Infection[[j]] = mean(covid$new_cases[1:j])
29     covid$Deaths[[j]] = mean(covid$new_deaths[1:j])
30   }
31   else{
32     covid$Infection[[j]] = sum( covid$new_cases[(j-6):j] )/7
33     covid$Deaths[[j]] = sum( covid$new_deaths[(j-6):j] )/7
34   }
35 }
36
37 for (i in year){
38   for (k in lmonth){
39     if(i==2022 & k > 2) break
40     for (c in lcountry){
41       #temp = covid %>% subset( location == c) %>% filter()
42
43       temp = subset(covid, format(covid$date, format= "%Y") == i & format(covid$date, format=
44 "%m") == k & covid$location == c)
45       temp = temp %>% filter(!Infection==0 | !Deaths ==0)
46
47       temp = temp %>% filter(!Infection ==0) %>% filter(!Deaths ==0)
48       #remove Outliers
49
50       grob1 = grobTree(textGrob(paste("Pearson Correlation : ", round(cor(temp$Infection, temp
51 $Deaths ), 4) ), x = 0.05, y = 0.97, hjust = 0, gp = gpar(col = "#63BDFE", fontsize = 11,
52 fontface = "bold")))
53       print(
54         ggplot(temp, aes(x =Infection, y=Deaths) ) + geom_point(color = "#275235")
55         + geom_smooth(formula = y ~ x, method = "lm", fill = "#91FFB4")+annotation_custom(grob1)+
56         theme_bw() +labs(title = paste("Country:",c,"- Month",k,"Year",i) )
57       )
58       cat(paste("Country:",c,"- Month",k,"Year",i))
59     }
60   }
61 }

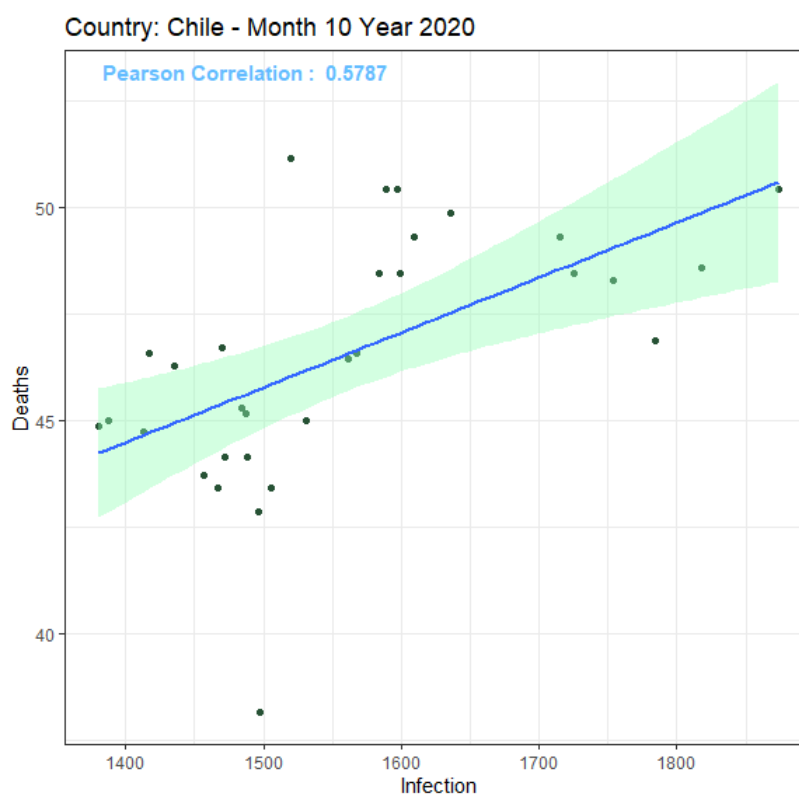
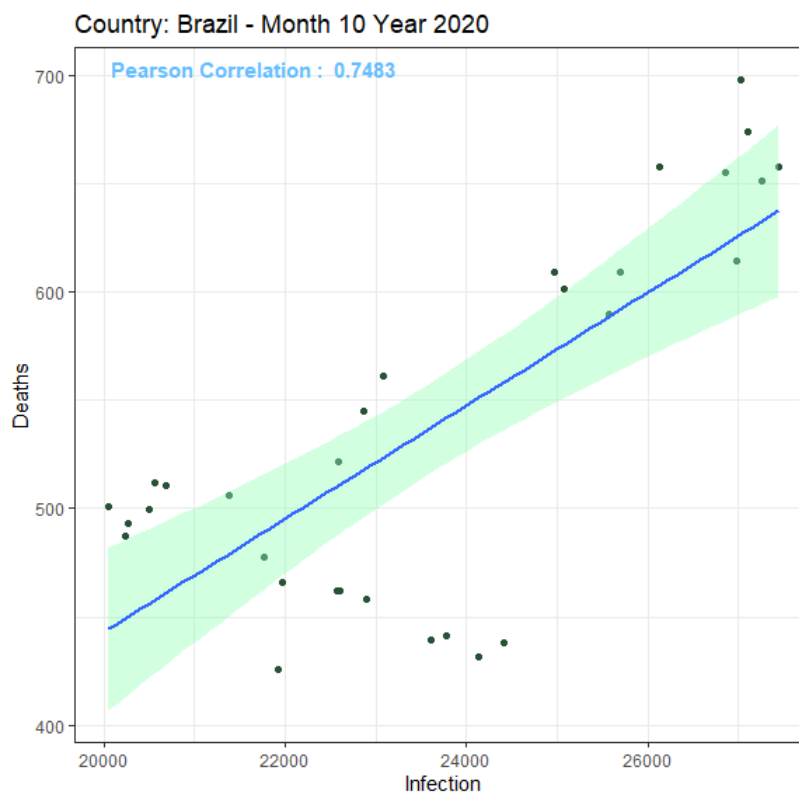
```

- Kết quả

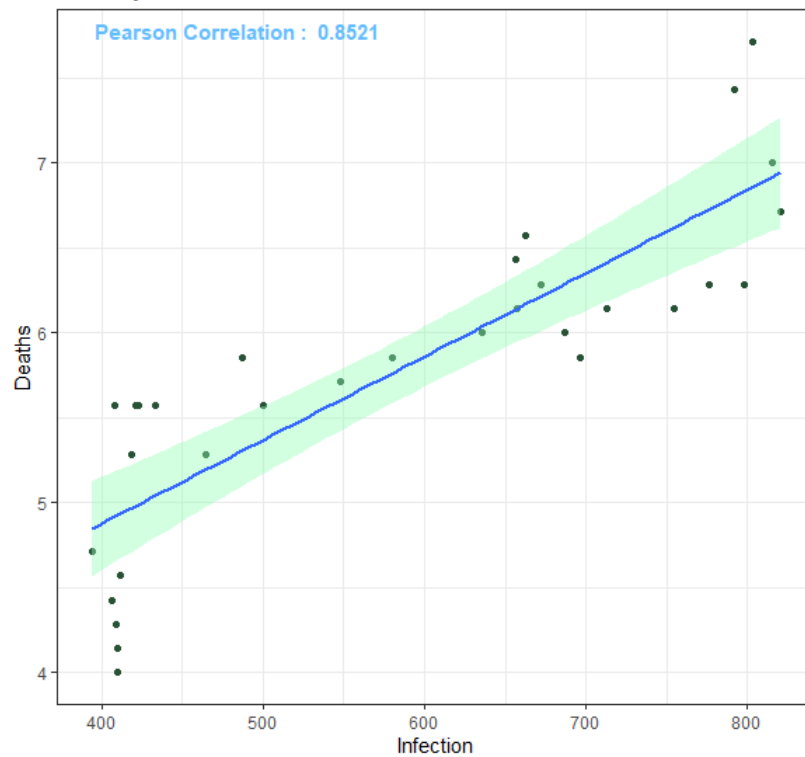




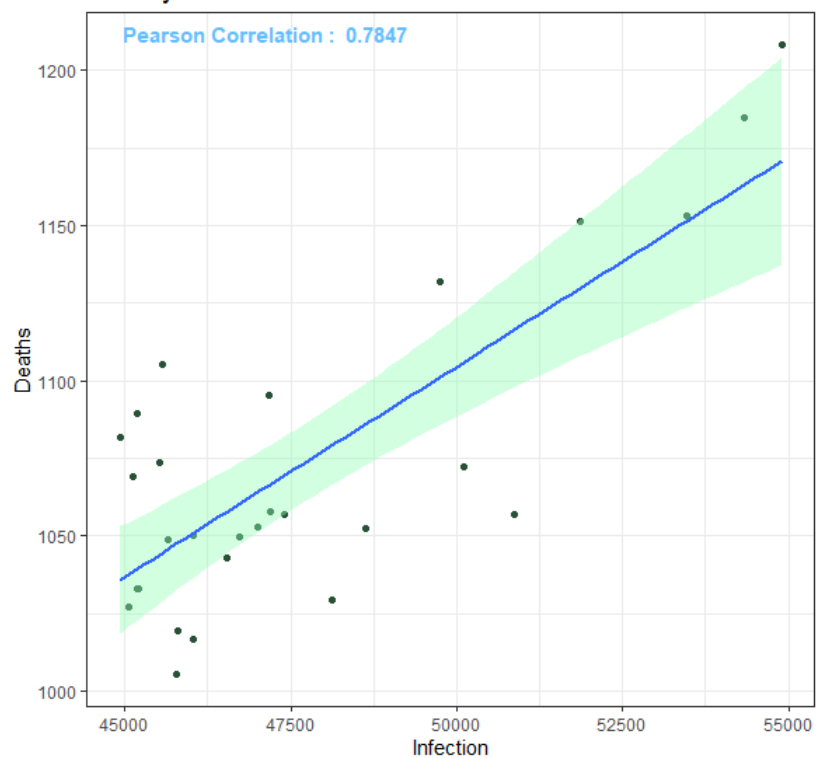


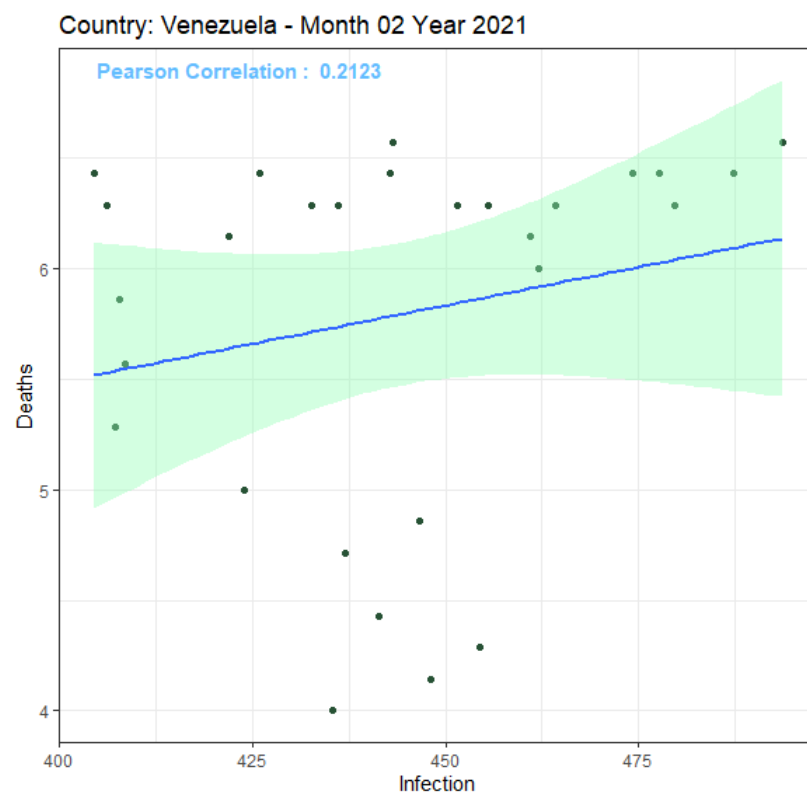
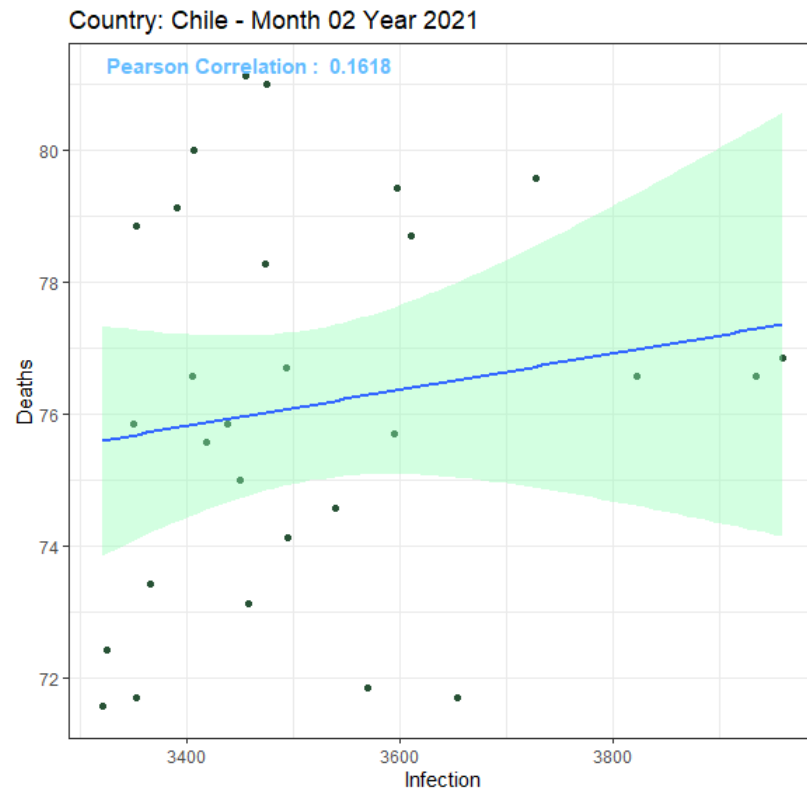


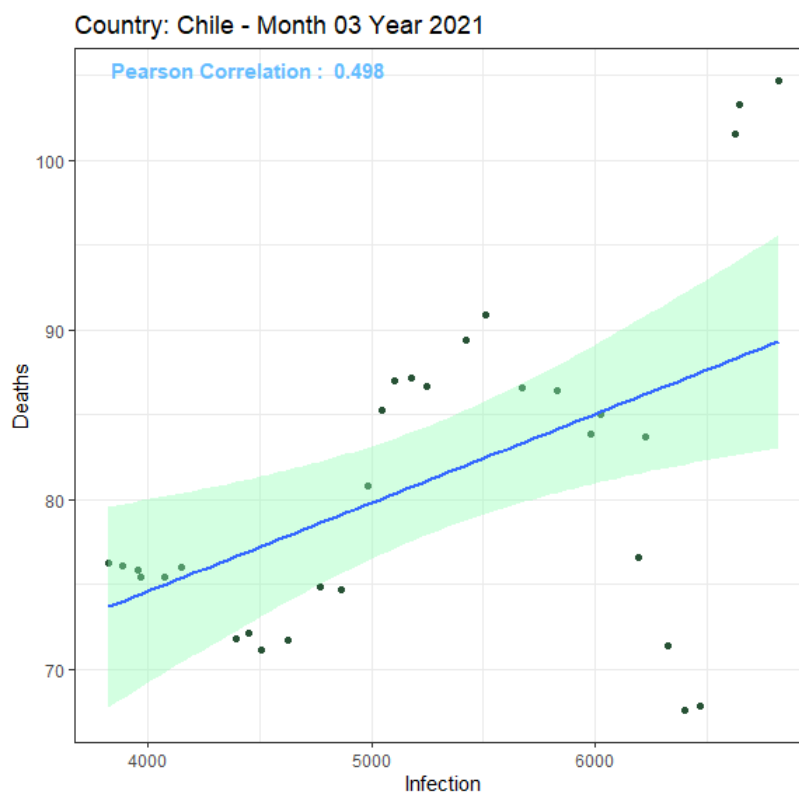
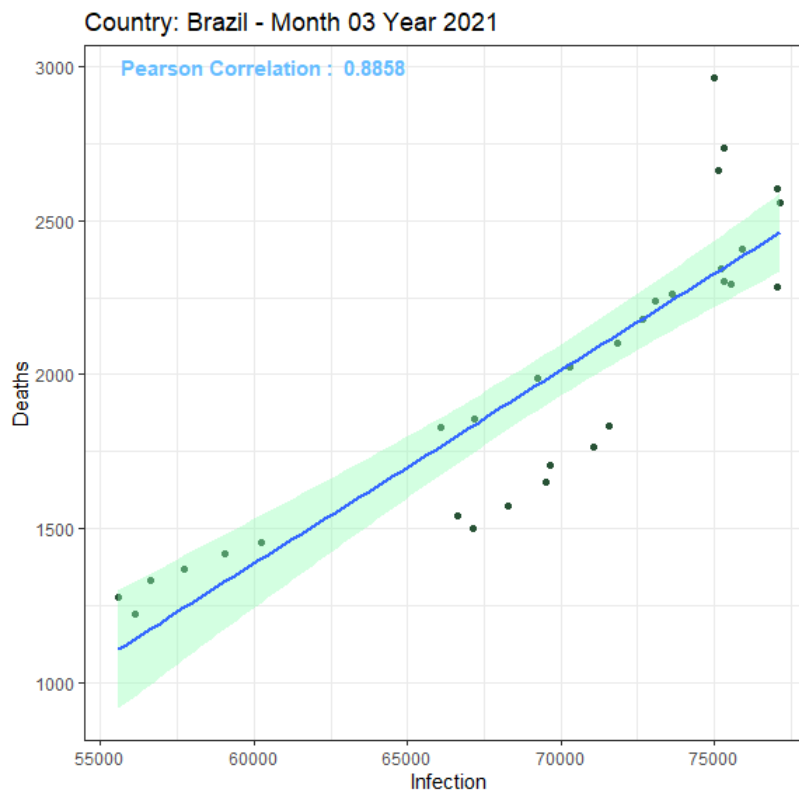
Country: Venezuela - Month 10 Year 2020

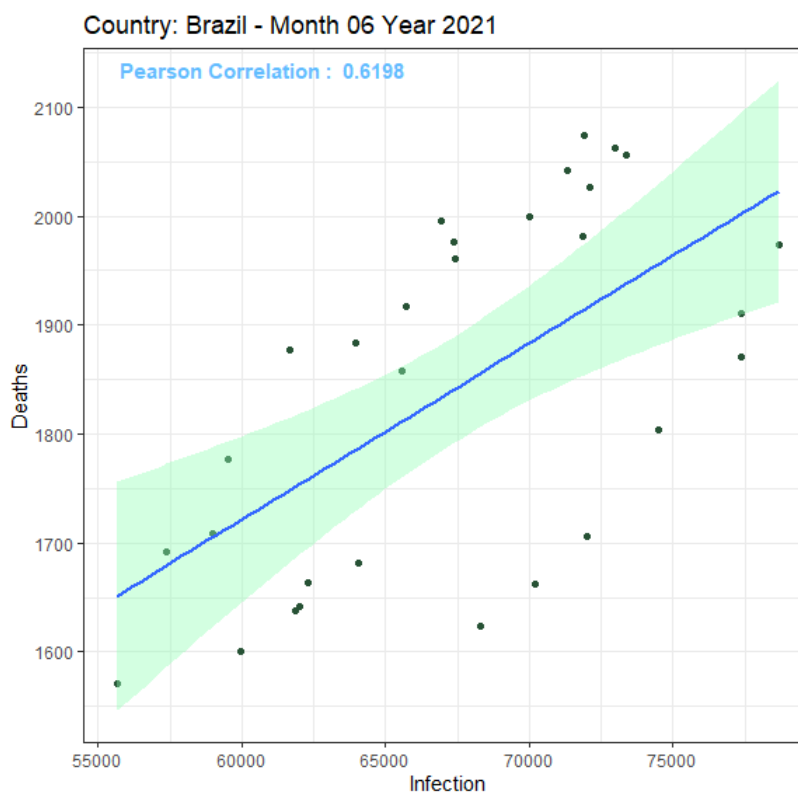
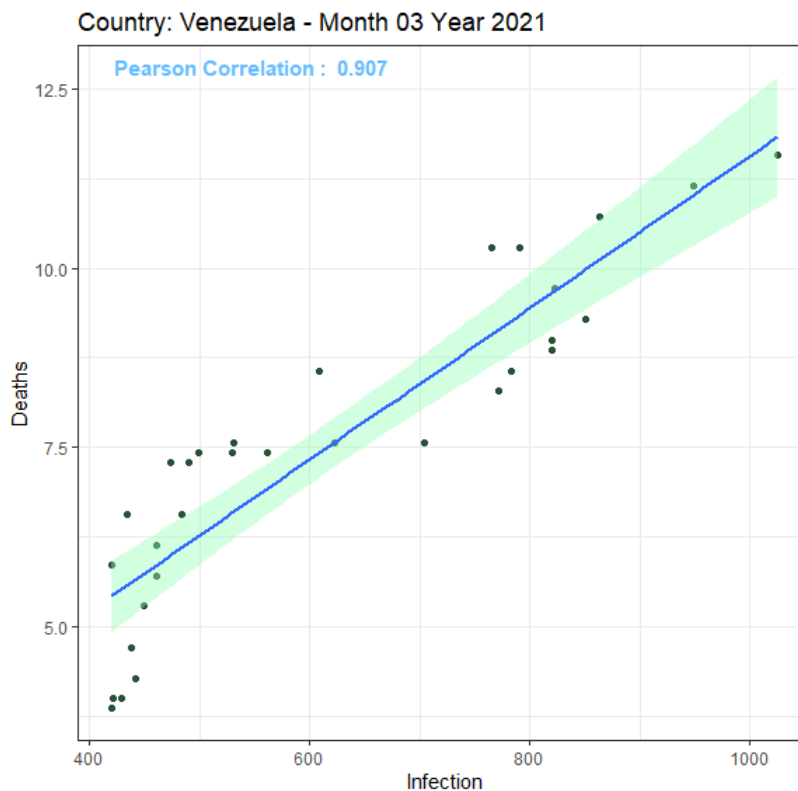


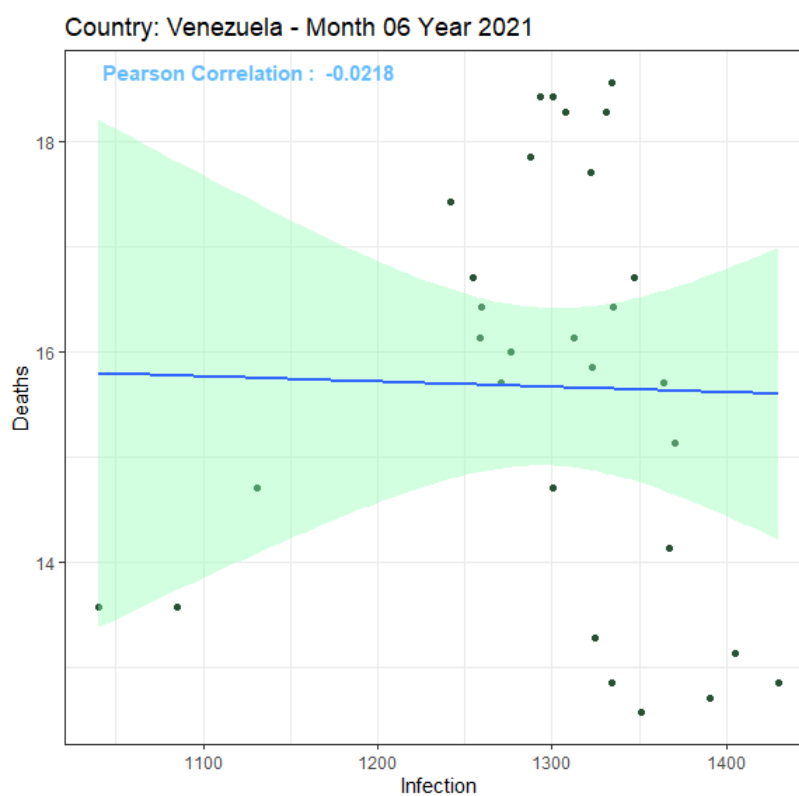
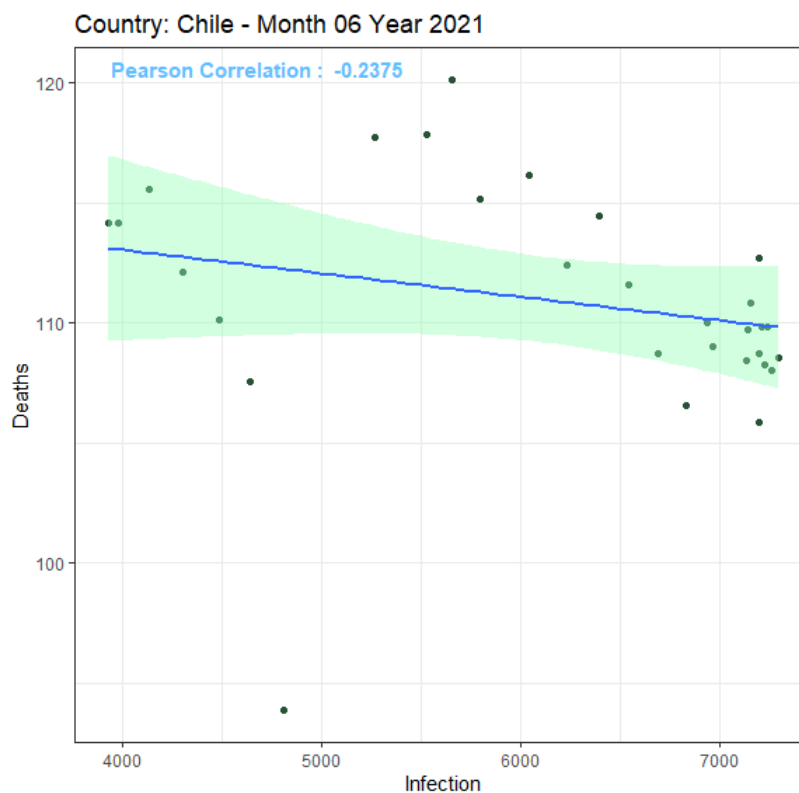
Country: Brazil - Month 02 Year 2021

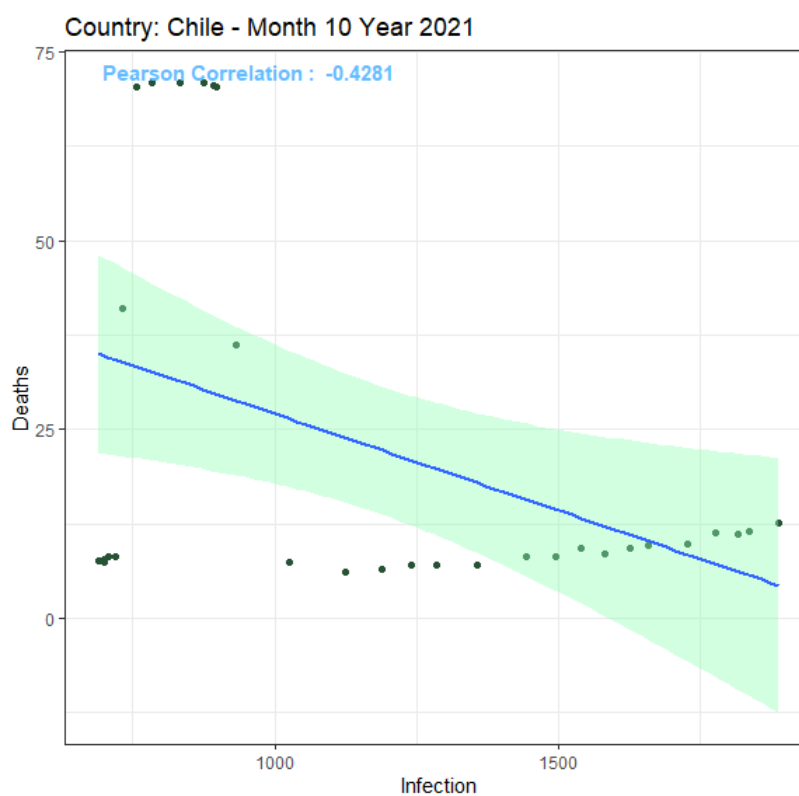
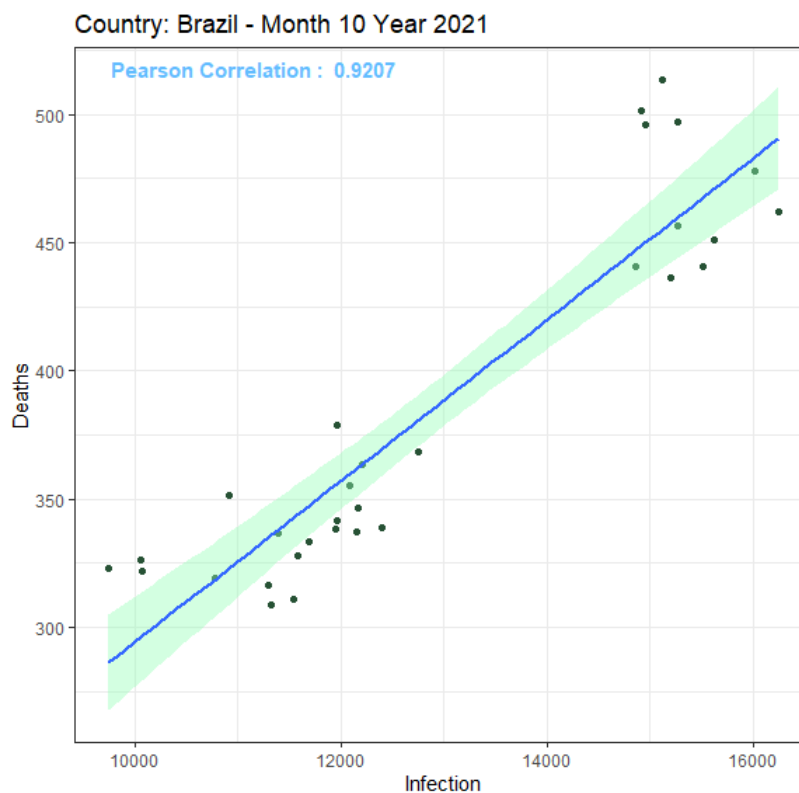


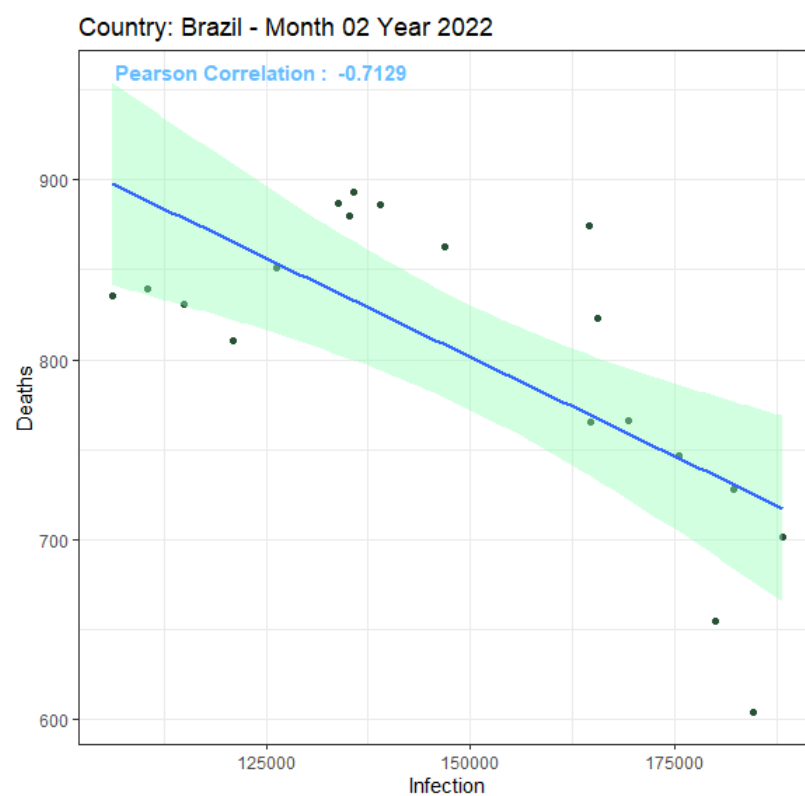
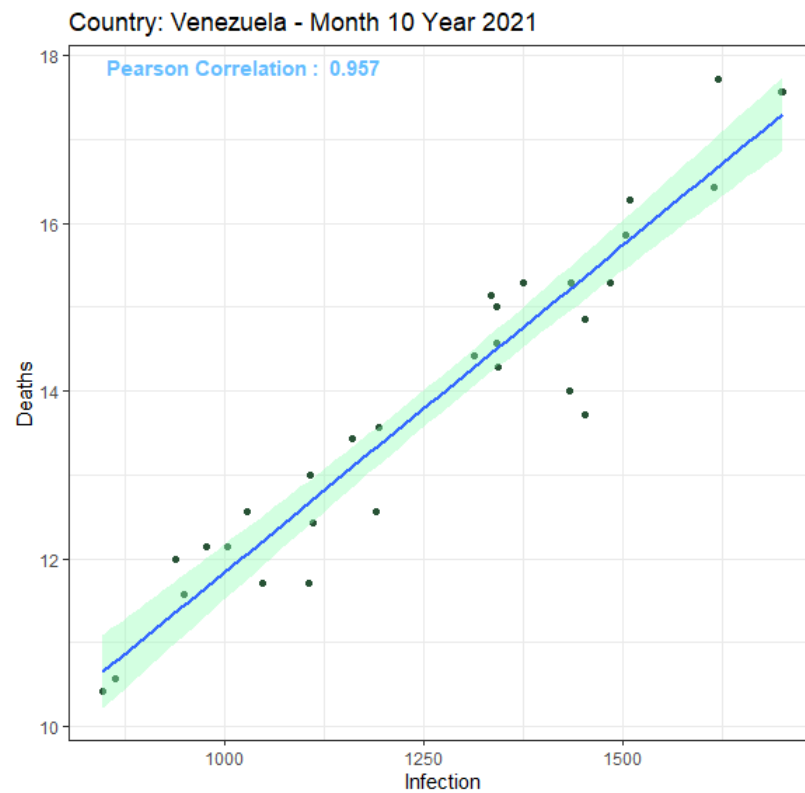


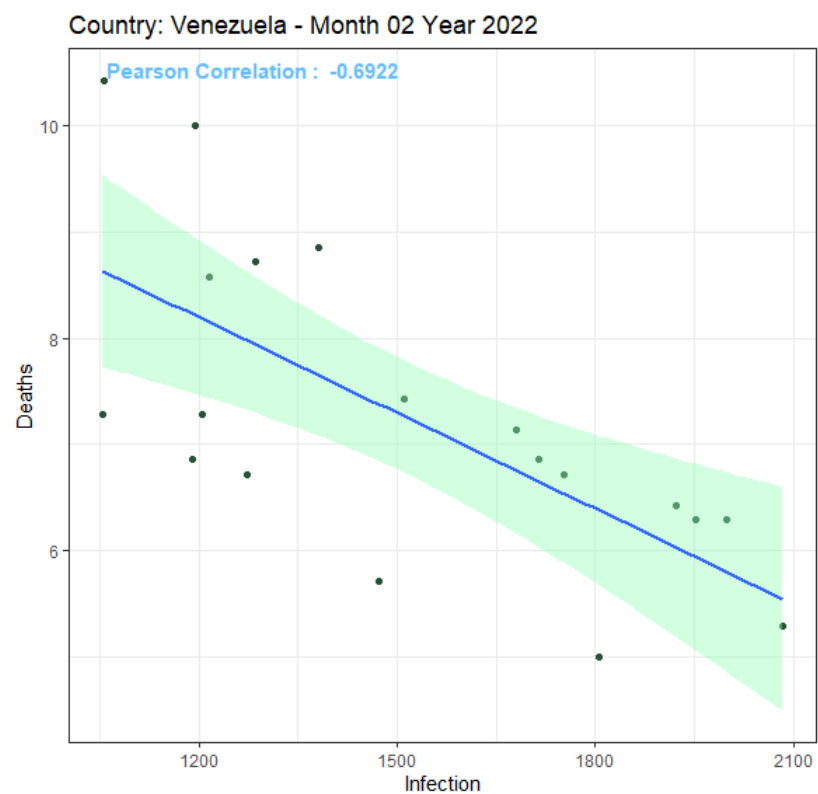
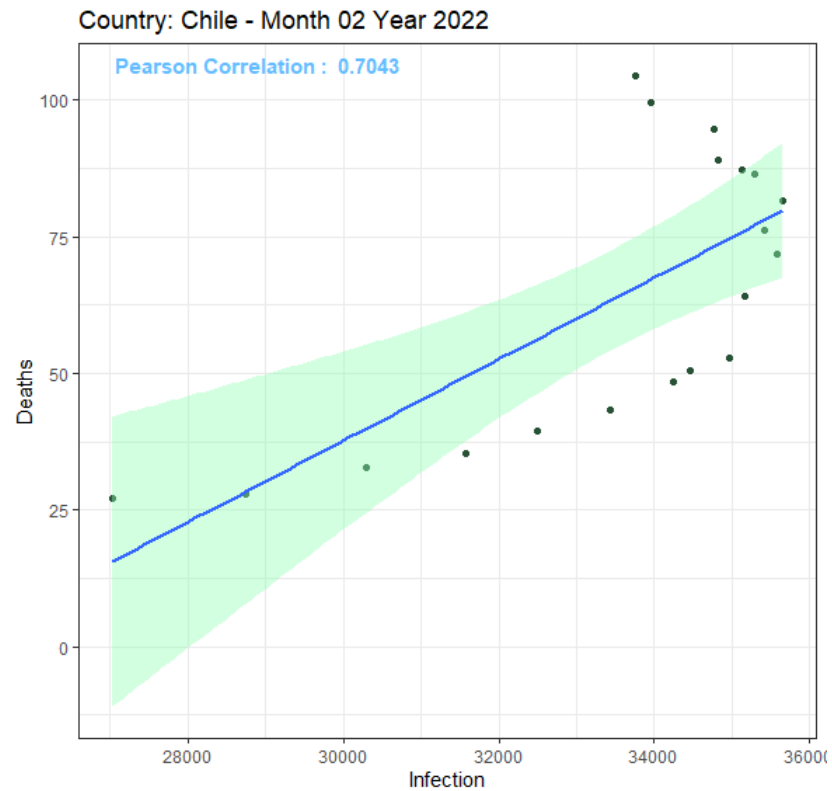






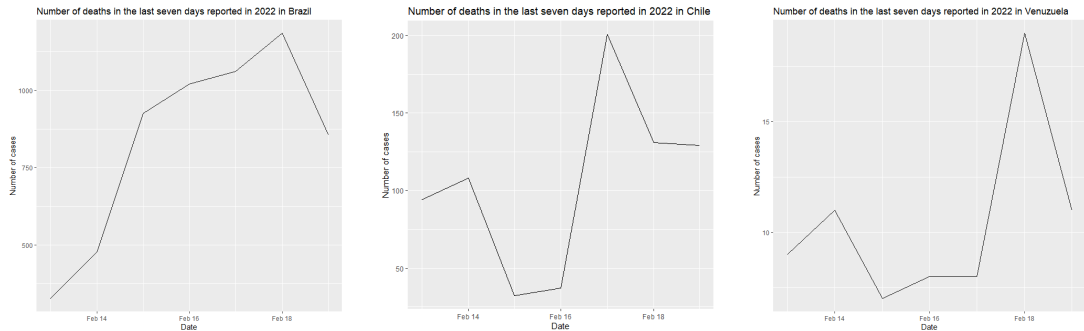






x) Nhóm câu hỏi riêng

- 2) So sánh tình trạng tử vong của các quốc gia trong 7 ngày cuối của năm cuối cùng
Sử dụng biểu đồ ở câu iv phần 4



- So sánh tình trạng tử vong:

Nhìn chung, tình trạng tử vong của cả ba nước đều đã đỉnh rồi bắt đầu giảm mạnh vào những ngày cuối. Trong đó, Brazil là nước có số ca tử vong lớn nhất, và rất lớn so với 2 nước còn lại, ngược lại số ca tử vong ở Venezuela trong 7 ngày cuối được cung cấp số liệu của năm 2022 là thấp nhất.

Ở Brazil, vào 7 ngày cuối được thu thập dữ liệu, số ca tử vong có xung hướng tăng mạnh đến khi đạt đỉnh số ca khoảng 1180 vào ngày 18/2 thì giảm mạnh vào ngày cuối -19/2. Biểu đồ của Chile và Venezuela chia sẻ xu thế giống nhau khi mà số ca tử vong dao động mạnh trước khi tăng nhanh đến điểm trong nửa sau của thời gian thực hiện khảo sát bằng biểu đồ. Sau khi chạm đỉnh với khoảng 200 ca và 19 ca ở lần lượt Chile và Venezuela, thì số ca tử vong bắt đầu giảm.

- 3) Cho biết các khoảng thời gian nào mà tỉ lệ tử vong tích lũy giảm mạnh nhưng tỉ lệ nhiễm bệnh tích lũy tăng mạnh hoặc ngược lại cho các quốc gia.

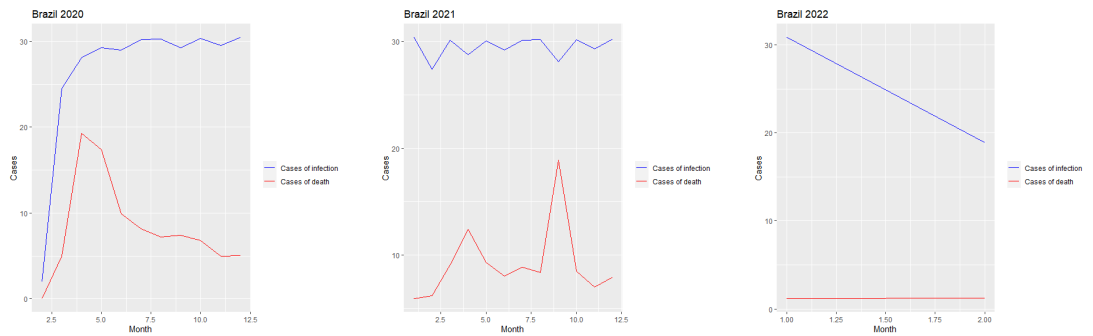
- Đối với Brazil

```
1 library(tidyverse)
2 library(datasets)
3 library(dplyr)
4 library("ggplot2")
5 library(lubridate)
6 covid = read.csv("C:/Users/Admin/Downloads/covidData.csv")
7 covid$date = as.Date(covid$date, format = "%m/%d/%Y")
8 covid = covid %>% filter(!continent == '') %>% mutate(new_cases = abs(new_cases), new_deaths =
9   abs(new_deaths))
10 covid[is.na(covid)] = 0
11 brazil = subset(covid, iso_code == "BRA")
12 chile = subset(covid, iso_code == "CHL")
13 venezuela = subset(covid, iso_code == "VEN")
14
15 #Doi voi brazil
16 brazil = brazil %>% mutate( totalcases = new_cases + new_deaths) %>% filter(!totalcases == 0)
17 brazil = brazil %>% mutate(new_cases = new_cases/totalcases , new_deaths = new_deaths/
18   totalcases)
19
20 brazil = brazil %>% mutate(month = as.numeric(format(as.Date(brazil$date), "%m")), year =
21   as.numeric(format(as.Date(brazil$date), "%Y")))
22 bra = aggregate(cbind(brazil$new_cases, brazil$new_deaths), by = list(brazil$month, brazil$
23   year), sum)
24 bra$V2 = bra$V2 * 10
25 ggplot(data=bra%>%filter(Group.2 == 2020), aes(x= Group.1)) + geom_line(aes(y= V1,color= "
26   Cases of infection")) +
27   geom_line(aes(y= V2,color = "Cases of death")) + xlab("Month") + ylab("Cases") + ggtitle("
28   Brazil 2020") +
29   scale_colour_manual("",
30     breaks = c("Cases of infection", "Cases of death"),
31     values = c("blue", "red"))
32
33 ggplot(data=bra%>%filter(Group.2 == 2021), aes(x= Group.1)) + geom_line(aes(y= V1,color= "
34   Cases of infection")) +
35   geom_line(aes(y= V2,color = "Cases of death")) + xlab("Month") + ylab("Cases") + ggtitle("
36   Brazil 2021") +
37   scale_colour_manual("",
```

```

31     breaks = c("Cases of infection", "Cases of death"),
32     values = c("blue", "red"))
33 ggplot(data=bra%>%filter(Group.2 == 2022), aes(x= Group.1)) + geom_line(aes(y= V1,color= "
34     Cases of infection")) +
35     geom_line(aes(y= V2,color = "Cases of death")) + xlab("Month") + ylab("Cases") + ggtitle("
36     Brazil 2022") +
37     scale_colour_manual("",
38     breaks = c("Cases of infection", "Cases of death"),
39     values = c("red", "green"))

```



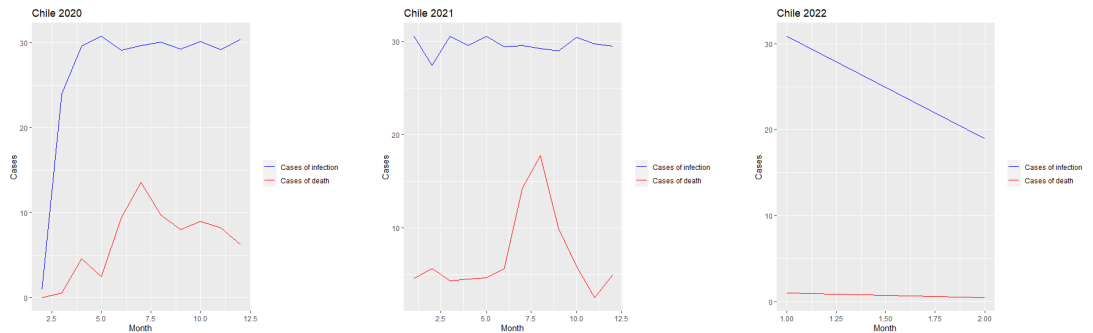
- Nhận xét từ đồ thị của Brazil:

Khoảng thời gian nào mà tỉ lệ tử vong tích lũy giảm mạnh nhưng tỉ lệ nhiễm bệnh tích lũy tăng mạnh ở Brazil: Vào khoảng tháng 4 năm 2020 đến hết năm 2020 và khoảng tháng 9 đến tháng 10 năm 2021.

Khoảng thời gian nào mà tỉ lệ nhiễm bệnh tích lũy giảm mạnh nhưng tỉ lệ tử vong tích lũy tăng mạnh ở Brazil: Vào khoảng tháng 3 đến tháng 4 và khoảng tháng 8 đến tháng 9 năm 2021.

Thực hiện đoạn code tương tự với Chile và Venezuela

- Đối với Chile

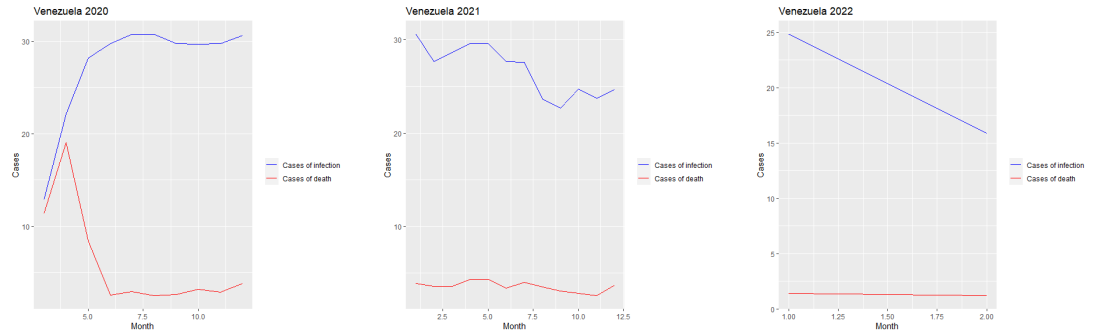


- Nhận xét từ đồ thị của Chile:

Khoảng thời gian nào mà tỉ lệ tử vong tích lũy giảm mạnh nhưng tỉ lệ nhiễm bệnh tích lũy tăng mạnh ở Chile: tháng 4 đến tháng 5 và tháng 7 đến tháng 9 năm 2020. Ngoài ra tỉ lệ tử vong tích lũy giảm mạnh nhất vào tháng 8 đến tháng 11 năm 2021.

Khoảng thời gian nào mà tỉ lệ nhiễm bệnh tích lũy giảm mạnh nhưng tỉ lệ tử vong tích lũy tăng mạnh ở Chile: tháng 5 đến tháng 6 năm 2020 và tháng 1 đến tháng 2 năm 2021.

- Đối với Venezuela



- Nhận xét từ đồ thị của Venezuela:

Khoảng thời gian nào mà tỉ lệ tử vong tích lũy giảm mạnh nhưng tỉ lệ nhiễm bệnh tích lũy tăng mạnh ở Venezuela: tháng 4 đến tháng 6 năm 2020.

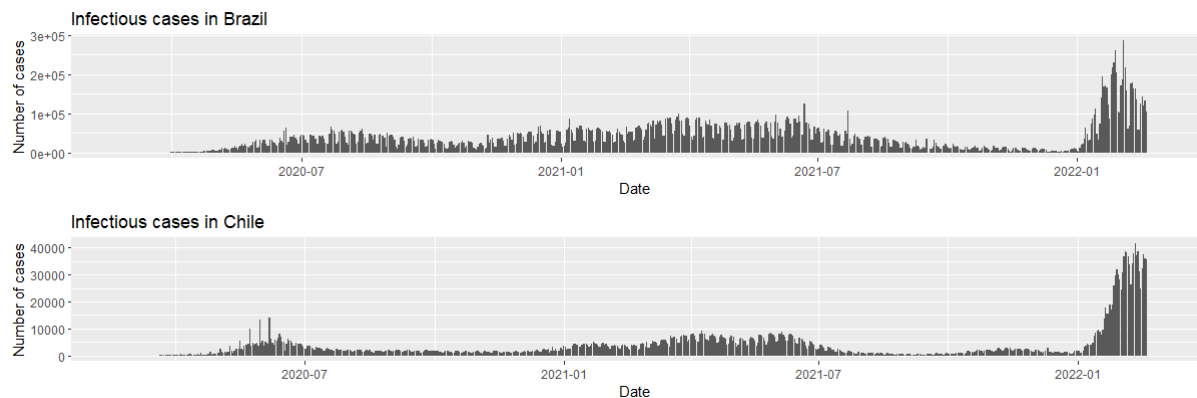
Khoảng thời gian nào mà tỉ lệ nhiễm bệnh tích lũy giảm mạnh nhưng tỉ lệ tử vong tích lũy tăng mạnh ở Venezuela: tỉ lệ nhiễm bệnh tích lũy có xu hướng giảm kể từ lần tăng đạt đỉnh vào giữa tháng 7 năm 2020 trong khi tỉ lệ tử vong tích lũy lại dao động không đáng kể

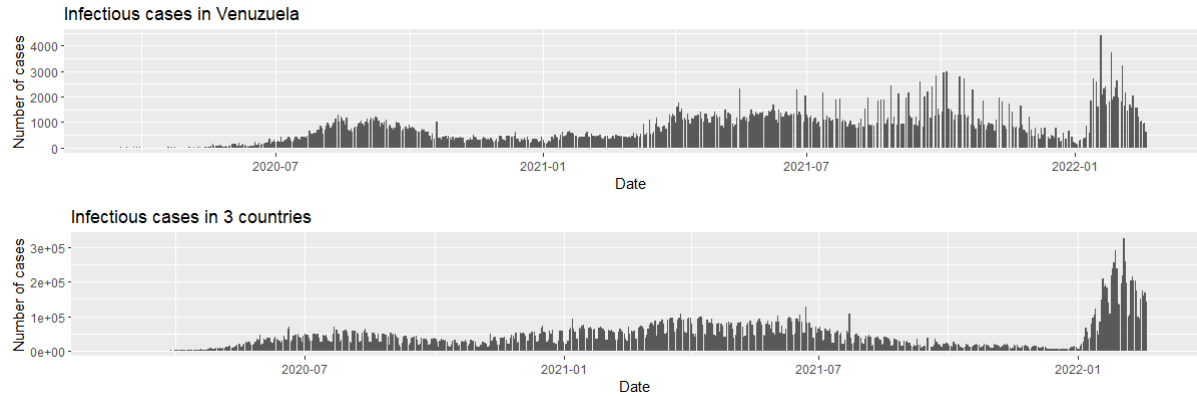
- 6) Khoảng thời bùng phát nhiễm bệnh lớn nhất giữa các quốc gia có chồng lên nhau không, Cho biết khoảng thời gian giao nhau đó?

- Hiện thực trong R:

```
1 brazil = subset(covid, iso_code == "BRA")
2 bra = brazil[order(brazil$date),]
3 ggplot(data = bra, aes(x = date, y = new_cases)) + geom_col() + labs(title = "Infectious cases
  in Brazil", x = "Date", y = "Number of cases")
4
5 chile = subset(covid, iso_code == "CHL")
6 chi = chile[order(chile$date),]
7 ggplot(data = chi, aes(x = date, y = new_cases)) + geom_col() + labs(title = "Infectious cases
  in Chile", x = "Date", y = "Number of cases")
8
9 venezuela = subset(covid, iso_code == "VEN")
10 ven = venezuela[order(ven$date),]
11 ggplot(data = ven, aes(x = date, y = new_cases)) + geom_col() + labs(title = "Infectious cases
  in Venezuela", x = "Date", y = "Number of cases")
12
13 total <- rbind(brazil, venezuela, chile)
14 tot = total[order(total$date),]
15 ggplot(data = tot, aes(x = date, y = new_cases)) + geom_col() + labs(title = "Infectious cases
  in 3 countries", x = "Date", y = "Number of cases")
16 }
```

- Kết quả





- Nhận xét: Dựa vào các biểu đồ, khoảng thời gian bùng phát nhiễm bệnh lớn nhất ở cả 3 quốc gia Brazil, Chile, Venezuela có sự giao nhau vào khoảng tháng 2 năm 2022.
- 10) Hãy mô tả mối quan hệ tuyến tính giữa nhiễm bệnh và tử vong bằng cách đo độ kết hợp của mối quan hệ dùng correlation r (correlation coefficient) và hướng kết hợp.

- Hiện thực trong R:

```

1 library(tidyverse)
2 library(dplyr)
3 library(datasets)
4 library("ggplot2")
5 library("ggpubr")
6 setwd("E:/BTL_CTRR")
7 covid = read.csv("owid-covid-data.csv", header = TRUE)
8 covid = covid %>% filter(!continent == '') %>% mutate(new_cases = abs(new_cases), new_deaths =
9   abs(new_deaths))
10 brazil = subset(covid, iso_code == "BRA")
11 chile = subset(covid, iso_code == "CHL")
12 venezuela = subset(covid, iso_code == "VEN")
13 brazil[is.na(brazil)] = 0
14 chile[is.na(chile)] = 0
15 venezuela[is.na(venezuela)] = 0
16 Countries = c("Brazil", "Chile", "Venezuela")
17
18 #remove outliers from new cases col
19 q1 = quantile(brazil$new_cases, 0.25)
20 q3 = quantile(brazil$new_cases, 0.75)
21 iqr = IQR(brazil$new_cases)
22 brazil = brazil %>% subset(brazil$new_cases > (q1 - 1.5*iqr) & brazil$new_cases < (q3 + 1.5*iqr)
23   )
24 q1 = quantile(brazil$new_deaths, 0.25)
25 q3 = quantile(brazil$new_deaths, 0.75)
26 iqr = IQR(brazil$new_deaths)
27 brazil = brazil %>% subset(brazil$new_deaths > (q1 - 1.5*iqr) & brazil$new_deaths < (q3 + 1.5*
28   iqr) )
29
30 ggscatter(brazil, x = "new_cases", y = "new_deaths",
31   add = "reg.line", conf.int = TRUE,
32   cor.coef = TRUE, cor.method = "pearson",
33   xlab = "Infections", ylab = "Deaths", title = "brazil")
34
35 q1 = quantile(chile$new_cases, 0.25)
36 q3 = quantile(chile$new_cases, 0.75)
37 iqr = IQR(chile$new_cases)
38 chile = chile %>% subset(chile$new_cases > (q1 - 1.5*iqr) & chile$new_cases < (q3 + 1.5*iqr) )
39
40 q1 = quantile(chile$new_deaths, 0.25)
41 q3 = quantile(chile$new_deaths, 0.75)
42 iqr = IQR(chile$new_deaths)
43 chile = chile %>% subset(chile$new_deaths > (q1 - 1.5*iqr) & chile$new_deaths < (q3 + 1.5*iqr)
44   )

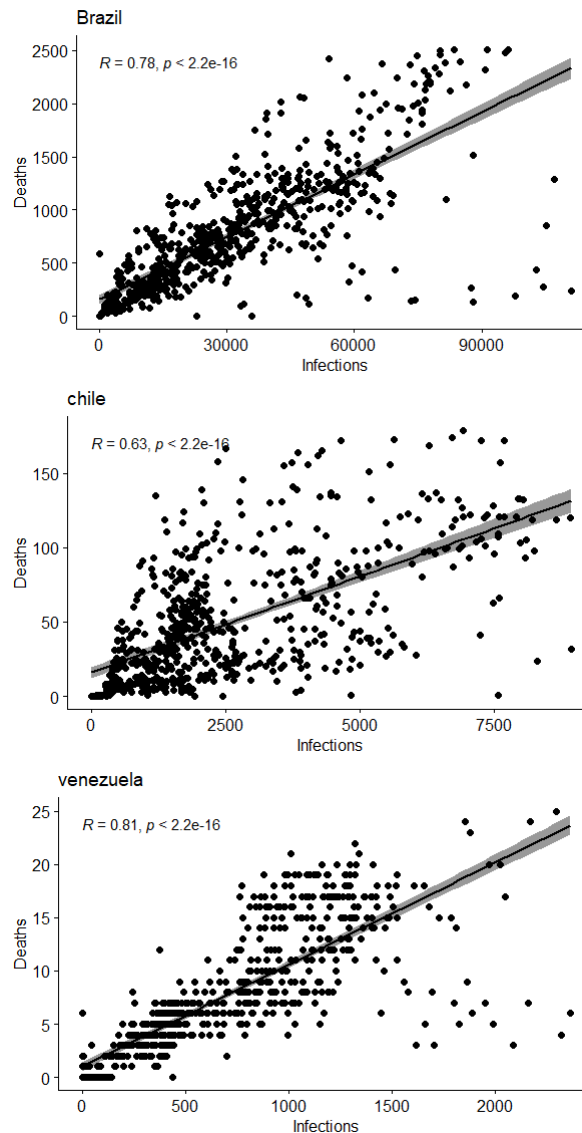
```

```

45 ggscatter(chile, x = "new_cases", y = "new_deaths",
46           add = "reg.line", conf.int = TRUE,
47           cor.coef = TRUE, cor.method = "pearson",
48           xlab = "Infections", ylab = "Deaths", title = "chile")
49
50 #remove outliers from new cases col
51 q1 = quantile(venezuela$new_cases, 0.25)
52 q3 = quantile(venezuela$new_cases, 0.75)
53 iqr = IQR(venezuela$new_cases)
54 venezuela = venezuela %>% subset(venezuela$new_cases > (q1 - 1.5*iqr) & venezuela$new_cases < (
55                               q3 + 1.5*iqr) )
56
57 q1 = quantile(venezuela$new_deaths, 0.25)
58 q3 = quantile(venezuela$new_deaths, 0.75)
59 iqr = IQR(venezuela$new_deaths)
60 venezuela = venezuela %>% subset(venezuela$new_deaths > (q1 - 1.5*iqr) & venezuela$new_deaths <
61                               (q3 + 1.5*iqr) )
62
63 ggscatter(venezuela, x = "new_cases", y = "new_deaths",
64           add = "reg.line", conf.int = TRUE,
65           cor.coef = TRUE, cor.method = "pearson",
66           xlab = "Infections", ylab = "Deaths", title = "venezuela")
67 }

```

- Kết quả



- Nhận xét:
 - Đối với cả 3 biểu đồ tương ứng với 3 quốc gia Brazil, Chile, Venezuela, hệ số tương quan Pearson(R) đều dương (>0), nghĩa là số ca nhiễm(infections) và số ca tử vong(deaths) ở 3 quốc gia có mối tương quan tuyến tính với nhau. Nói cách khác, số ca nhiễm(infections) tăng thì số ca tử vong(deaths) cũng tăng.
 - Ta thấy, hệ số tương quan Pearson(R) ứng với Venezuela là 0.81 lớn nhất trong 3 nước nên mối tương quan tuyến tính giữa số ca nhiễm và số ca tử vong được thu thập tại Venezuela là chặt chẽ nhất.
 - Ngược lại, hệ số tương quan Pearson(R) ứng với Chile là 0.63 bé nhất nên mối tương quan tuyến tính giữa số ca nhiễm và số ca tử vong được thu thập tại đây ít chặt chẽ hơn 2 nước còn lại.

Tài liệu

- [Dal] Dalgaard, P. *Introductory Statistics with R*. Springer 2008.
- [K-Z] Kenett, R. S. and Zacks, S. *Modern Industrial Statistics: with applications in R, MINITAB and JMP*, 2nd ed., John Wiley and Sons, 2014.
- [Ker] Kerns, G. J. *Introduction to Probability and Statistics Using R*, 2nd ed., CRC 2015.
- [4] CSETI. *Thống kê mô tả trong nghiên cứu - Các đại lượng về độ phân tán*. <<http://thongke.cesti.gov.vn/dich-vu-thong-ke/tai-lieu-phan-tich-thong-ke/845-thong-ke-mo-ta-trong-nghien-cuu-dai-luong-do-phan-tan>>.
- [5] Sololearn. *R course*, <<https://www.sololearn.com/learning/1147>>.
- [6] Đoàn Quỳnh, Nguyễn Huy Doan, Nguyễn Xuân Liêm, Đặng Hùng Thắng, Trần Văn Vương. *Đại Số 10 nâng cao*, tái bản lần thứ 14. NXB Giáo dục Việt Nam 2020.
- [7] Luanvan2S. *Lý thuyết về hệ số tương quan Pearson - Phân tích tương quan trong SPSS*, <<https://luanvan2s.com/he-so-tuong-quan-pearson-trong-spss-bid61.html>>.