

ĐẠI HỌC ĐÀ NẴNG
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC CÔNG NGHỆ TIỀN TIẾN



BÁO CÁO MÔN HỌC
PHƯƠNG PHÁP THỐNG KÊ
VÀ
PHÂN TÍCH DỮ LIỆU

GVHD: TS. Lê Dân
SVTH: Nguyễn Đình Mẫn
Lớp: 18PFIEV3
Nhóm: 18.87

Đà Nẵng, 4/2021

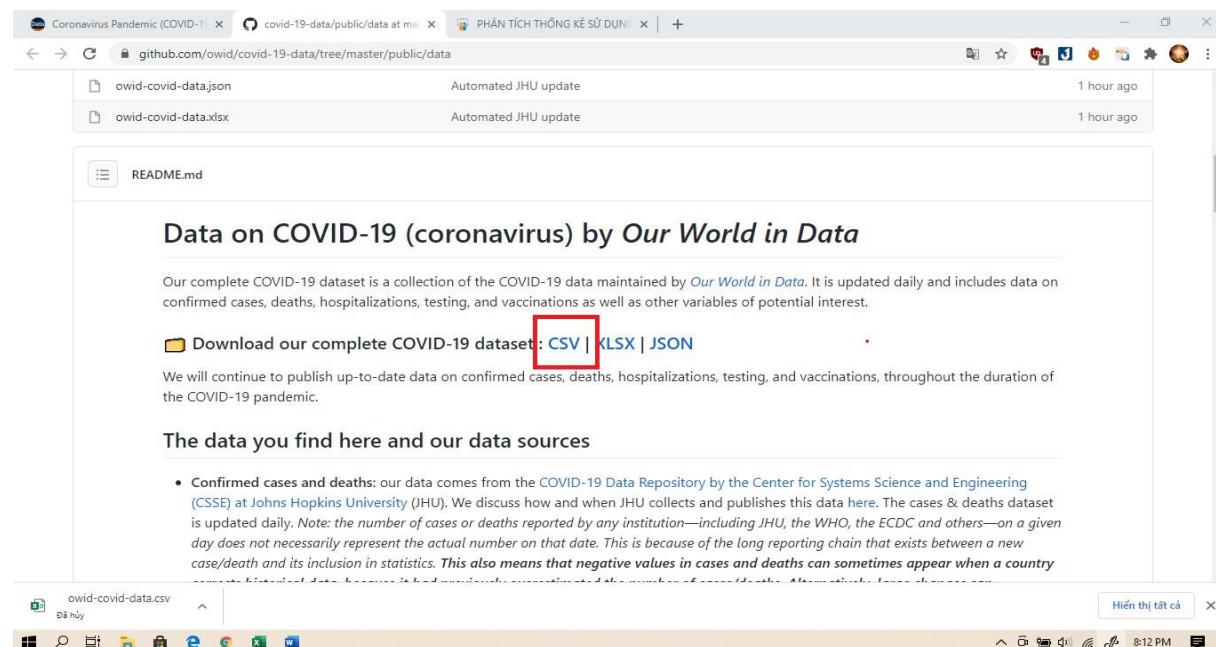
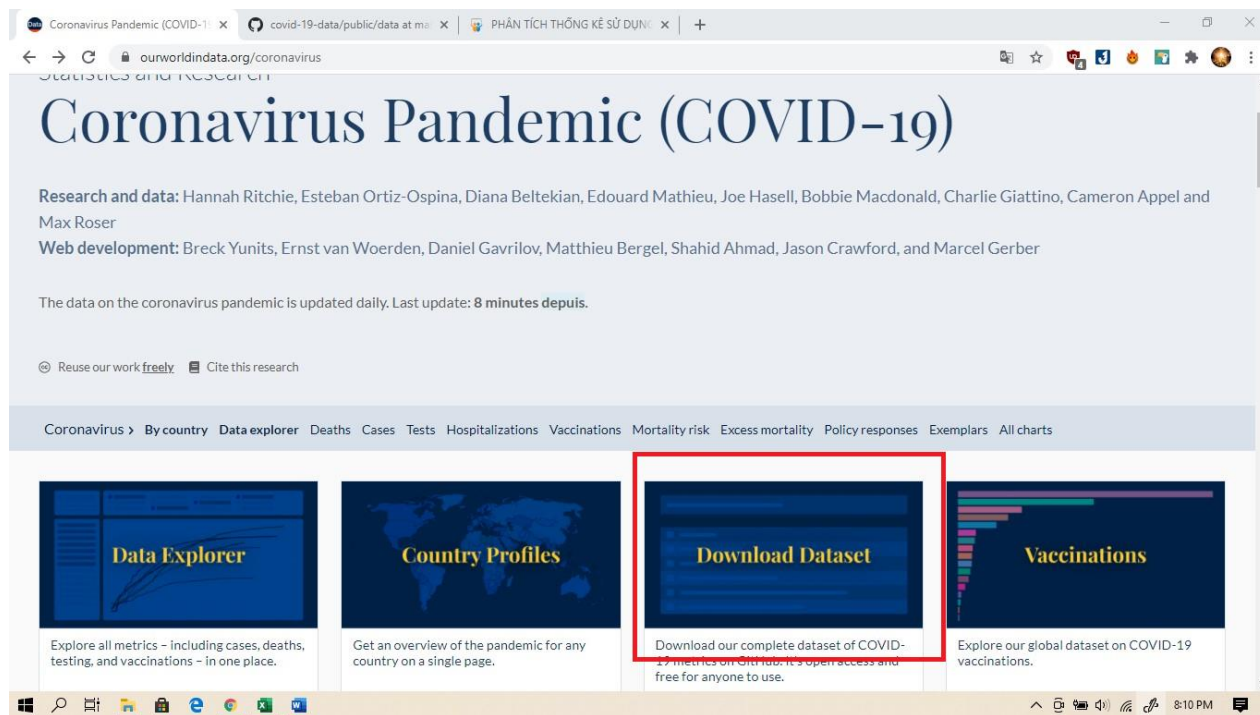
PHỤ LỤC

I. PHƯƠNG PHÁP DOWNLOAD VÀ LƯU TRỮ FILE DỮ LIỆU	3
1. Download dữ liệu từ trang web.....	3
2. Thực hiện xử lý file có đuôi .csv bằng Excel	4
II. PHƯƠNG PHÁP XỬ LÝ DỮ LIỆU TRÊN PHẦN MỀM SPSS.....	5
1. Cách import dữ liệu vào SPSS	5
2. Làm sạch dữ liệu bằng SPSS:	7
3. Mô tả dữ liệu và phân tích kết quả.....	9
3.1 Mô tả dữ liệu bằng Bảng trong SPSS.....	9
3.2 Mô tả dữ liệu bằng Đồ thị	11
3.3 Mô tả đồ thị bằng đại lượng thống kê.....	14
4. Kiểm định.....	15
a. Kiểm định trung bình tổng thể	15
b. Kiểm định 1 mẫu	16
c. Kiểm định mẫu cặp	17

I. PHƯƠNG PHÁP DOWNLOAD VÀ LƯU TRỮ FILE DỮ LIỆU

1. Download dữ liệu từ trang web

- Dữ liệu được lấy từ trang web: <https://ourworldindata.org/coronavirus>
- Chọn **Download Dataset** và chọn file định dạng .csv để tải về
- Nguồn dữ liệu đã được kiểm định, có mức độ tin cậy cao: Đây là dự án của [Global Change Data Lab](#) hợp tác với Đại học Oxford (UK)



2. Thực hiện xử lý file có đuôi .csv bằng Excel

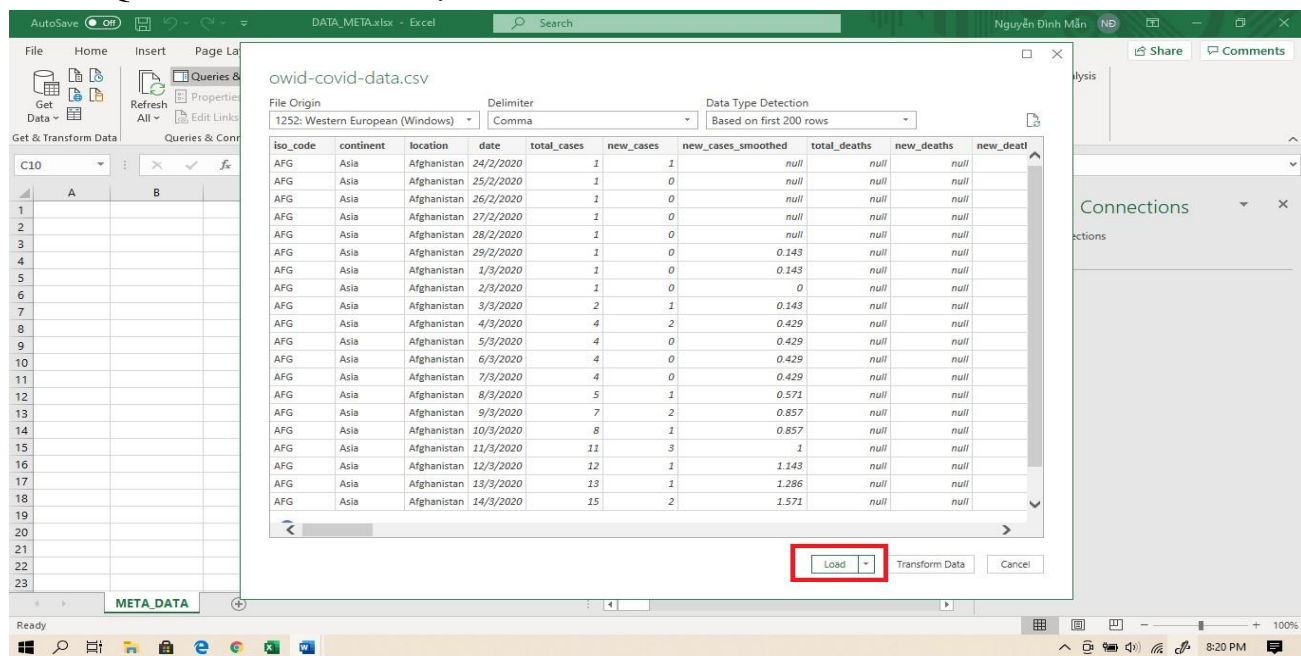
(Việc xử lý được thực hiện trên Excel - Office 365)

Bước 1: Mở trang Excel mới

Bước 2: Chọn **Data > Get Data > From File > From Text/Csv**

Bước 3: Chọn file có đuôi *.csv muốn import vào Excel, sau đó nhấn **Load**.

Quá trình load file dữ liệu:



- Dữ liệu sau khi Import từ file *.csv

The screenshot shows the Excel spreadsheet with the imported data. The data is displayed in a table with columns for iso_code, continent, location, date, total_cases, new_cases, new_cases_smoothed, total_deaths, new_deaths, and new_deaths_smoothed. The data is sorted by date, showing the progression of the pandemic from February 24, 2020, to March 16, 2020.

iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed
AFG	Asia	Afghanistan	24/2/2020	1	1	0.143	0	0	0
AFG	Asia	Afghanistan	25/2/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	26/2/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	27/2/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	28/2/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	29/2/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	1/3/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	2/3/2020	1	0	0.143	0	0	0
AFG	Asia	Afghanistan	3/3/2020	2	1	0.143	0	0	0
AFG	Asia	Afghanistan	4/3/2020	4	2	0.429	0	0	0
AFG	Asia	Afghanistan	5/3/2020	4	0	0.429	0	0	0
AFG	Asia	Afghanistan	6/3/2020	4	0	0.429	0	0	0
AFG	Asia	Afghanistan	7/3/2020	4	0	0.429	0	0	0
AFG	Asia	Afghanistan	8/3/2020	5	1	0.571	0	0	0
AFG	Asia	Afghanistan	9/3/2020	7	2	0.857	0	0	0
AFG	Asia	Afghanistan	10/3/2020	8	1	0.857	0	0	0
AFG	Asia	Afghanistan	11/3/2020	11	3	1	0	0	0
AFG	Asia	Afghanistan	12/3/2020	12	1	1.143	0	0	0
AFG	Asia	Afghanistan	13/3/2020	13	1	1.286	0	0	0
AFG	Asia	Afghanistan	14/3/2020	15	2	1.571	0	0	0
AFG	Asia	Afghanistan	15/3/2020	16	1	1.571	0	0	0
AFG	Asia	Afghanistan	16/3/2020	18	2	1.571	0	0	0

Bước 4: Thực hiện tổ chức lại dữ liệu

Dữ liệu sau khi được tổ chức lại:

AutoSaveOff

DATA_META.xlsx - Excel

Search

Nguyễn Đình SảnND

Share

Comments

FileHomeInsertPage LayoutFormulasDataReviewViewHelpPower PivotTable Design

Calibri11A^A

B

I

U

Font

Wrap Text

Merge & Center

Alignment

General

\$ % ∞

Number

Conditional Formatting

Format as Table

Cell Styles

Insert

Delete

Format

Cells

Sort & Find & Filter & Select

Editing

Ideas

Sensitivity

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

Ideas

Sensitivity

113

<

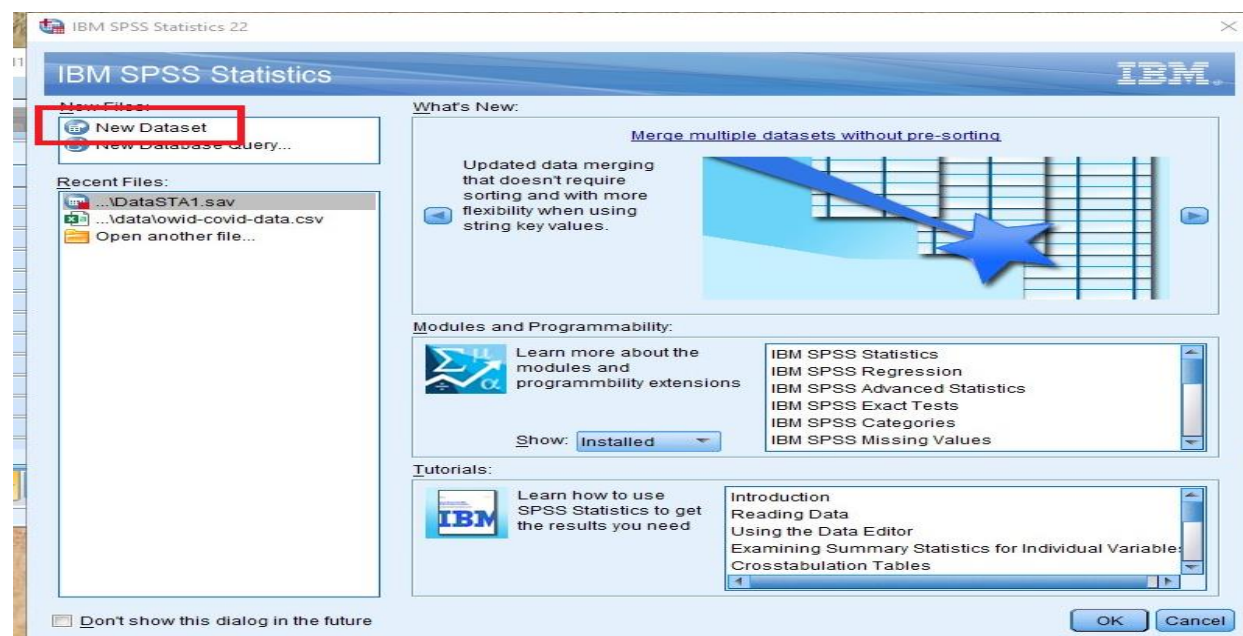
Bước 5: Lưu file với tên META_DATE.xlsx

II. PHƯƠNG PHÁP XỬ LÝ DỮ LIỆU TRÊN PHẦN MỀM SPSS

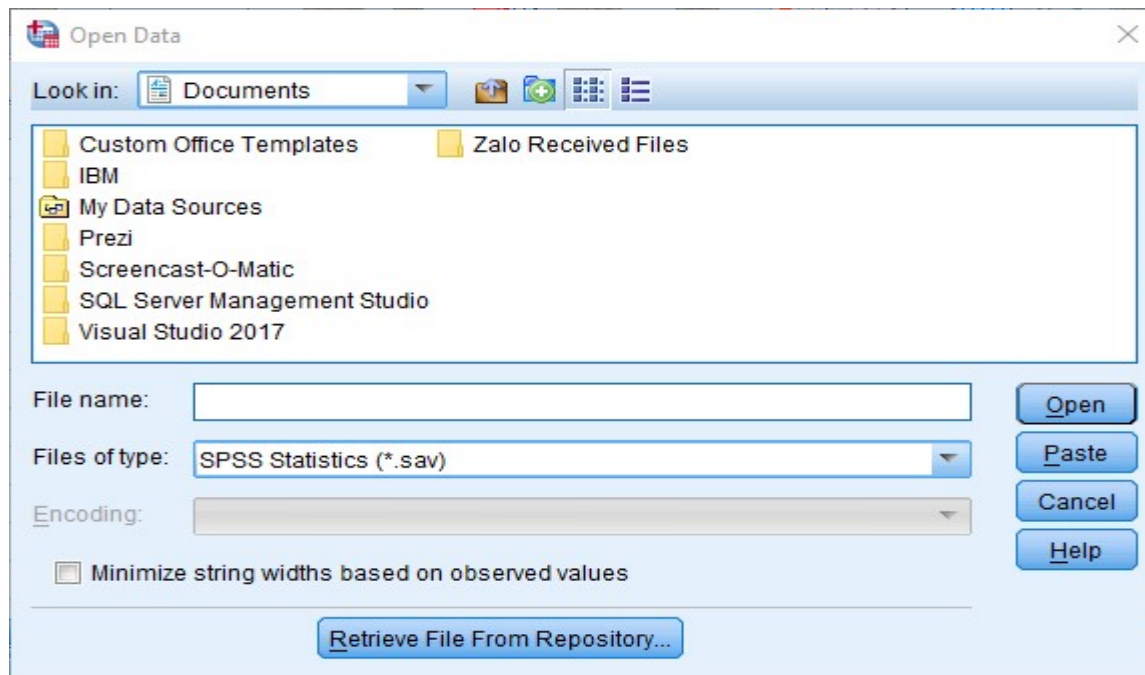
1. Cách import dữ liệu vào SPSS

Bước 1: Mở phần mềm IBM SPSS Statistic 22

Bước 2: Chọn New Dataset



Bước 3: Chọn **File > Open > Data**, xuất hiện hộp thoại **Open Data**



Bước 4:

- Chọn thư mục lưu file dữ liệu
- Điều chỉnh **File of type** thành định dạng **Excel**
- Chọn file META_DATA.xlsx. Sau đó nhấn **Open**
- File dữ liệu sẽ được hiển thị như sau:

SPSS Statistics Data Editor - IBM SPSS Statistics Data Editor

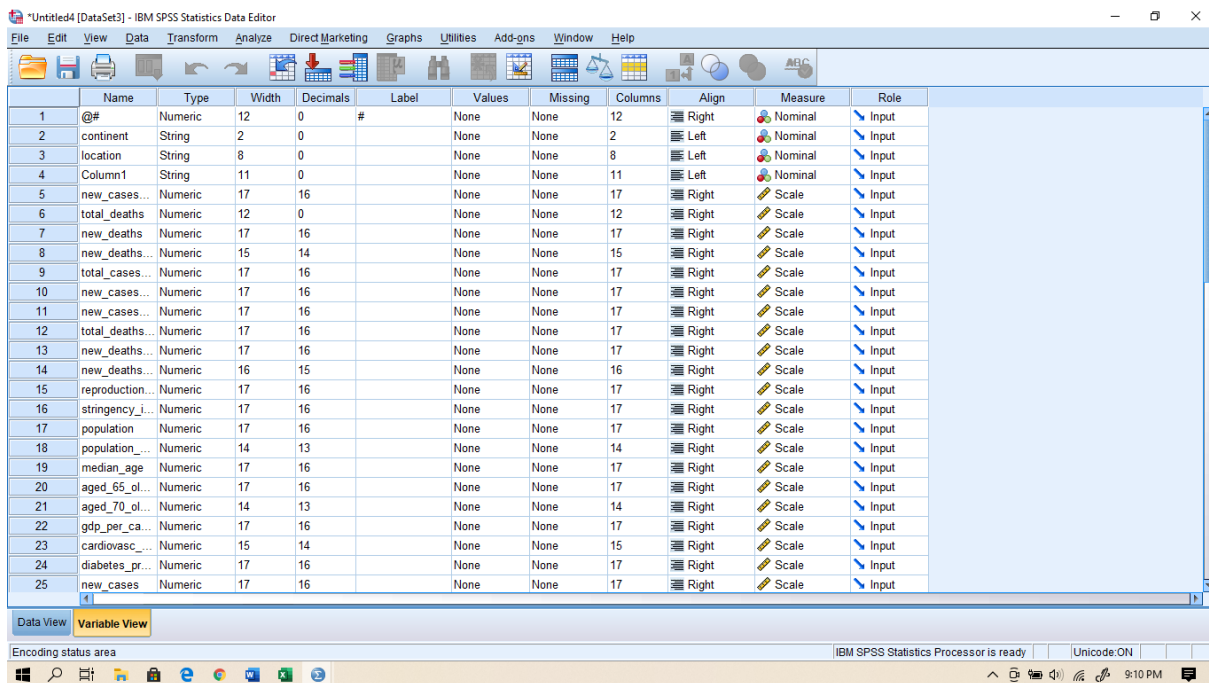
File Edit View Data Transform Analyze Direct Marketing Graphs Utilities Add-ons Window Help

Visible: 30 of 30 Variables

	@#	cont in t	location	Column1	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	total_cases_per_million	new_cases_per_million	n
1	1	SA	ARG	7	2290457.7090000000000000	8185904	55611.0000000000000000	55217.71299999997000	7177991.5630000001000000	51389.9959999999800000	
2	.	.	BOL	25	268419.8710000001600000	2193810	12211.0000000000000000	12148.0070000000000000	3646014.2419999999000000	23159.9660000000040000	
3	.	.	BRA	26	12381895.86500000200000	48456240	313866.0000000000000000	306651.00499999990000	8299663.7160000000000000	59153.4100000000000000	
4	.	.	CHL	34	962790.7090000000600000	4122773	23070.0000000000000000	22624.8590000000000000	8419062.6449999980000000	51499.9620000000360000	
5	.	.	COL	40	2367972.294000002000000	9474759	63079.0000000000000000	62665.15599999996000	6723891.8519999990000000	46966.2769999999500000	
6	.	.	ECU	54	320392.4329999999600000	3409876	16746.0000000000000000	16631.5780000000000000	2928522.2639999995000000	18427.8630000000000000	
7	.	.	FLK	62	
8	.	.	GUY	80	9981.1400000000000000	32938	229.0000000000000000	223.1600000000000000	1602860.0579999993000000	12944.9840000000060000	
9	.	.	PER	159	1508473.8479999998000000	9496060	51635.0000000000000000	51021.8600000000000000	7680632.2550000030000000	46497.8849999999950000	
10	.	.	PRY	164	204575.1450000001000000	433376	4113.0000000000000000	3962.013000000000100	2991303.3630000002000000	29502.1620000000040000	
11	.	.	SUR	183	9087.8580000000000000	29539	177.0000000000000000	177.020000000000004	2546956.0100000002000000	15519.0399999999800000	
12	.	.	URY	203	92541.5700000000000000	58229	928.0000000000000000	877.8920000000000000	1519868.7160000000000000	28667.7820000000070000	
13	2	OC	VEN	208	155008.0029999999400000	218470	1577.0000000000000000	1543.437000000000060	844608.2210000000600000	5554.3569999999950000	
14	.	.	AUS	10	7634449.0000000000000000	29296	29249.1370000000000000	212936.00000000000000	909.0000000000000000	909.0130000000006000	
15	.	.	FJI	61	12783.0000000000000000	67	65.88100000000000700	457.0000000000000000	2.0000000000000000	2.0020000000000000	
16	.	.	FSM	65	68.0000000000000000	1	2860000000000000000	.	.	.0000000000000000	
17	.	.	MHL	124	560.0000000000000000	4	2.574000000000000000	.	.	.0000000000000000	
18	.	.	NZL	144	677131.0000000000000000	2495	2481.5790000000000600	8224.0000000000000000	26.0000000000000000	26.0090000000000360	
19	.	.	PNG	161	213870.0000000000000000	5349	4546.2909999999990000	2293.0000000000000000	49.0000000000000000	41.0129999999999900	
20	.	.	SLB	175	2616.0000000000000000	18	16.28900000000000120	.	.	.0000000000000000	
21	.	.	VUT	210	184.0000000000000000	3	2.288000000000000000	.	.	.0000000000000000	
22	3	NA	WSM	211	298.0000000000000000	3	2.2879999999999994	.	.	.0000000000000000	

Data View Variable View

- Chuyển sang Variable View để xem thuộc tính của các biến và các tổ chức dữ liệu



Trong bảng Variable View cho thấy có 30 biến trong file dữ liệu được import vào.

Bước 5: Chọn **File > Save As** với tên DATA_SPSS.sav

2. Làm sạch dữ liệu bằng SPSS:

Do trong quá trình thực hiện cập nhật dữ liệu, các số liệu có thể bị nhập sai hoặc tại thời điểm nhập dữ liệu, một số thông tin vẫn chưa có nên chưa thể nhập. Ví dụ như các biến vẫn chưa nhập dữ liệu, nhập số âm,...

Do vậy trước khi dữ liệu được phân tích, cần làm sạch để hạn chế sai số trong quá trình phân tích số liệu.

Cách thực hiện:

Chọn: **Analyze > Descriptive Statistics > Frequencies**

Chọn các biến muốn làm sạch, nhấn OK

Sau khi xử lý, bảng tần số xuất hiện như sau

Variable	Value
new_cases_smoothed	3544957
total_deaths	4122773
new_deaths	4351796
new_deaths_smoothed	4477916
total_cases_per_million	4615295
new_cases_per_million	8185904
new_cases_smoothed_per_million	9474759
total_deaths_per_million	9496060
new_deaths_per_million	12095855
new_deaths_smoothed_per_million	30331024
reproduction_rate	48456240
Total	192044473
Missing System	16
Total	207

Nhận xét: ta thấy có lỗi Missing, tức là lỗi hàng trống trong dữ liệu nhập vào.

- Hoặc đối với biến continent:

Statistics

continent

N	Valid	207
	Missing	0

continent

	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	201	97.1	97.1	97.1
AF	1	.5	.5	97.6
AS	1	.5	.5	98.1
EU	1	.5	.5	98.6
NA	1	.5	.5	99.0
OC	1	.5	.5	99.5
SA	1	.5	.5	100.0
Total	207	100.0	100.0	

Nhận xét: ta thấy lỗi Missing = 0, điều này cho biết dữ liệu nhập vào không sai sót.

Ngoài cách lọc tần số nói trên, ta cũng có thể lọc dữ liệu bằng cách dùng bảng kết hợp hay xử lý trực tiếp trên excel.

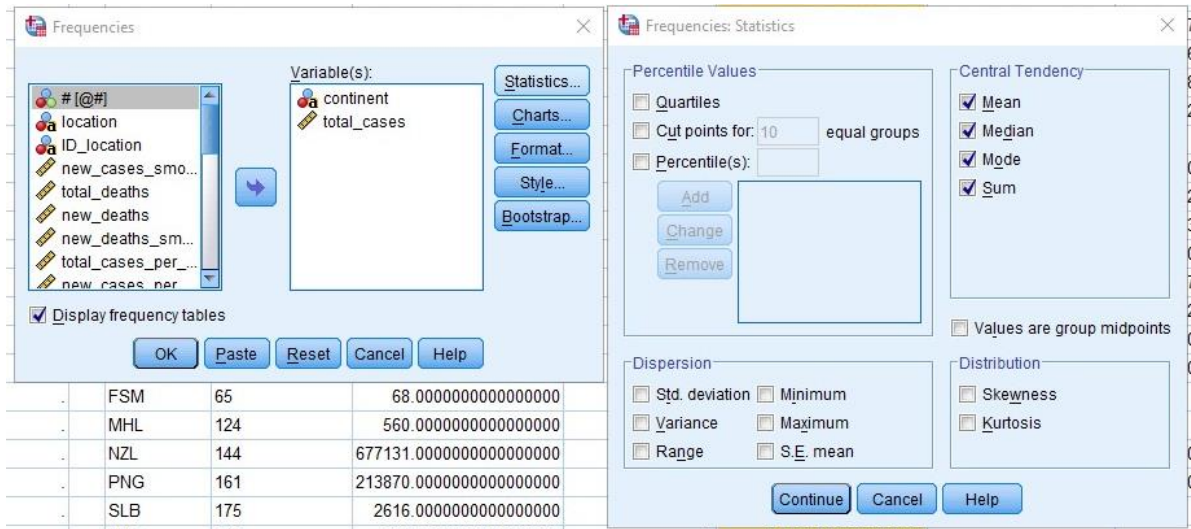
3. Mô tả dữ liệu và phân tích kết quả

3.1 Mô tả dữ liệu bằng Bảng trong SPSS

a. Mô tả một biến

- Cách thực hiện mô tả theo tần số xuất hiện:

- Chọn **Analyze > Descriptive Statistics > Frequencies**
- Chọn biến cần mô tả
- Nhấn vào nút Statistics để chọn các thông số
- Nhấn OK



- Sau khi thực hiện xong, của sổ Output sẽ xuất hiện:

➔ Frequencies

Statistics			
		continent	total_cases
N	Valid	207	191
	Missing	0	16
Mean			31340772.16
Median			2801.340000
Mode			895.400000 ^a
Sum			5986087482

a. Multiple modes exist. The smallest value is shown

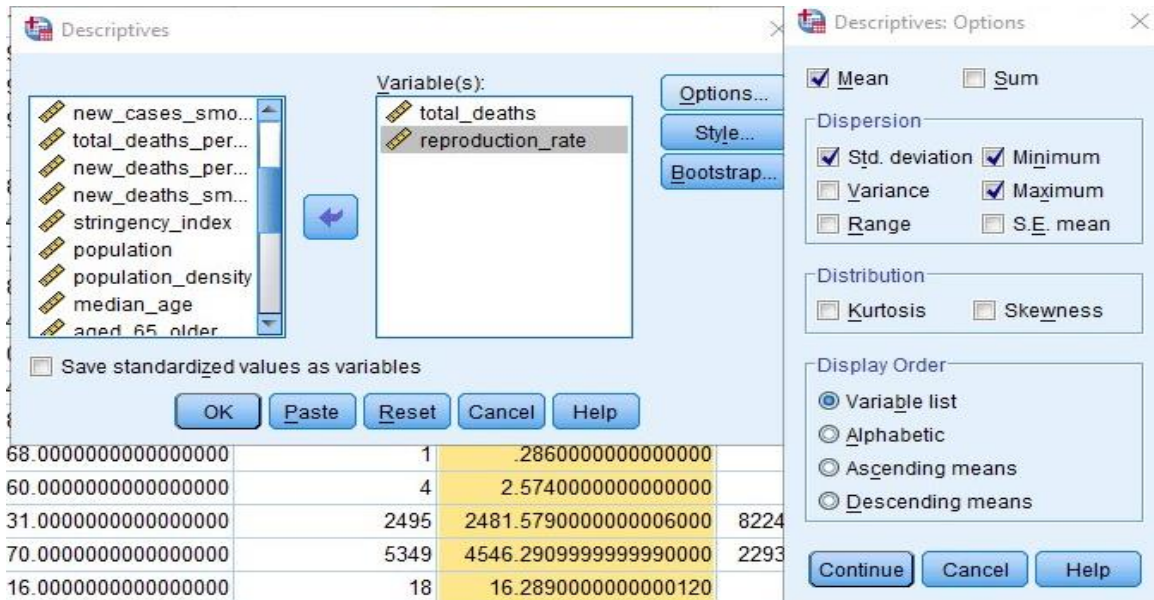
Frequency Table

continent				
	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	201	97.1	97.1	97.1
AF	1	.5	.5	97.6
AS	1	.5	.5	98.1
EU	1	.5	.5	98.6
NA	1	.5	.5	99.0
OC	1	.5	.5	99.5
SA	1	.5	.5	100.0
Total	207	100.0	100.0	

Nhận xét:

Bảng Frequencies, thể hiện thông số của các biến và các lỗi có thể có trong dữ liệu của biến.
Bảng Frequency Table thể hiện các giá trị tần số xuất hiện của các giá trị trong biến.

- **Cách thực hiện mô tả giá trị trung bình Mean**
 - Chọn **Analyze > Descriptive Statistics > Descriptives**
 - Chọn biến cần mô tả
 - Nhấn vào nút Options để chọn các thông số
 - Nhấn OK



- Bảng Output

Descriptive Statistics					
	N	Minimum	Maximum	Mean	Std. Deviation
new_cases_smoothed	191	68.00000000	1.41337E+10	147996406.6	1072579454
total_deaths	191	1	192044473	2010936.89	14498178.28
new_deaths	191	.2860000000	105176151.9	1101320.963	7950851.273
new_deaths_smoothed	181	71.00000000	345373335.7	3816279.952	26686499.89
total_cases_per_million	181	1.000000000	56629851.90	625744.2199	4393747.446
new_cases_per_million	191	.0000000000	2620323.245	27437.93974	194391.4896
new_cases_smoothed_per_million	191	591.1920000	561585494.4	5880476.381	40603955.80
total_deaths_per_million	191	6.728000000	6331683.960	66300.35560	458226.4863
new_deaths_per_million	191	2.484000000	4614989.362	48324.49594	333698.6269
new_deaths_smoothed_per_million	181	4.247000000	10580257.84	116908.9264	786895.1781
reproduction_rate	181	.3510000000	84719.72400	936.1295470	6290.175764
stringency_index	191	.0000000000	414946.9750	4344.994503	30626.34567
Valid N (listwise)	181				

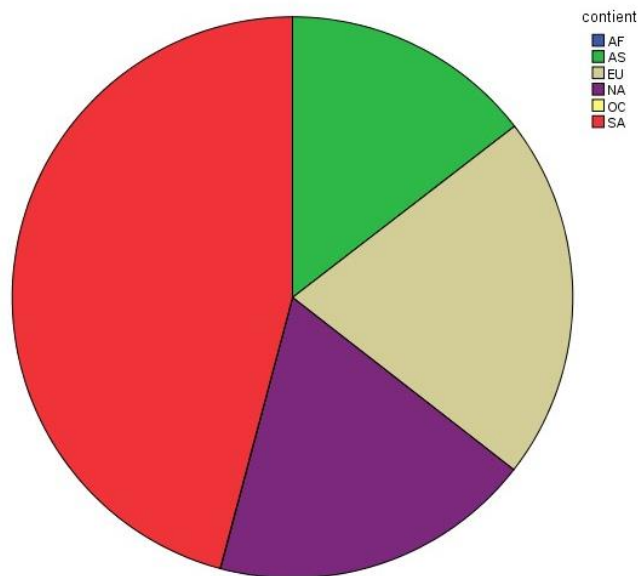
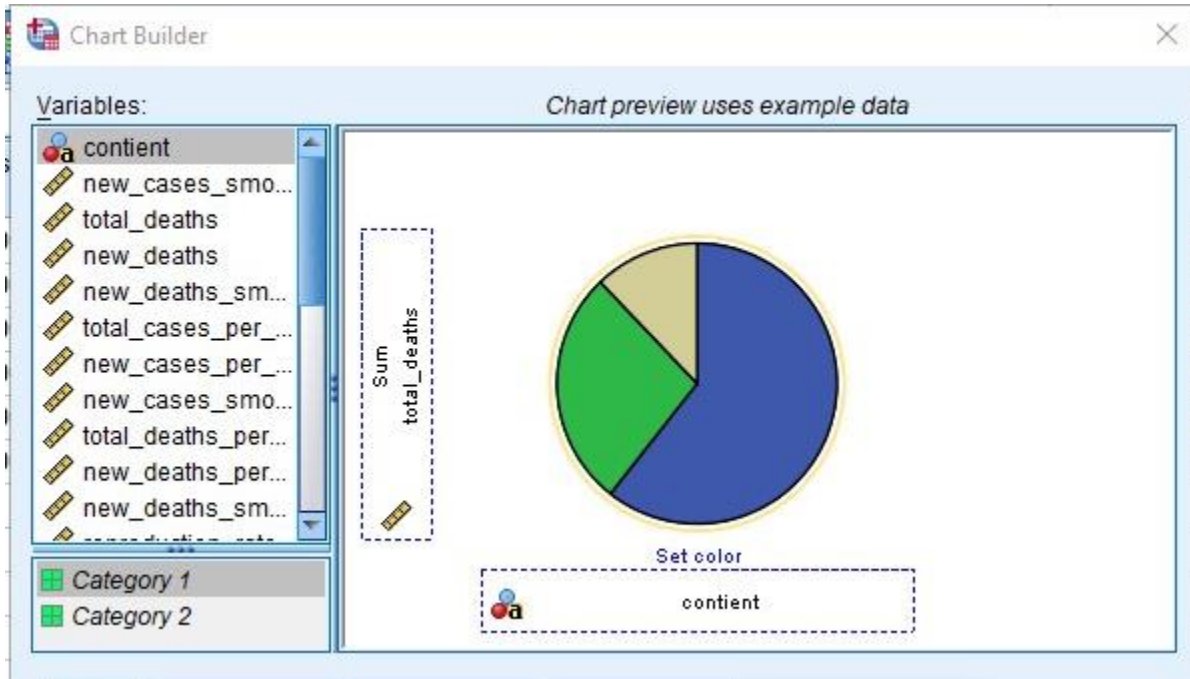
Bảng thể hiện các giá trị trung bình của các biến được chọn.

3.2 Mô tả dữ liệu bằng Đồ thị

- Cách thực hiện:

- Chọn **Graphs > Chart Builder**
- Chọn loại biểu đồ phù hợp với dữ liệu
- Chọn các biến cho biểu đồ

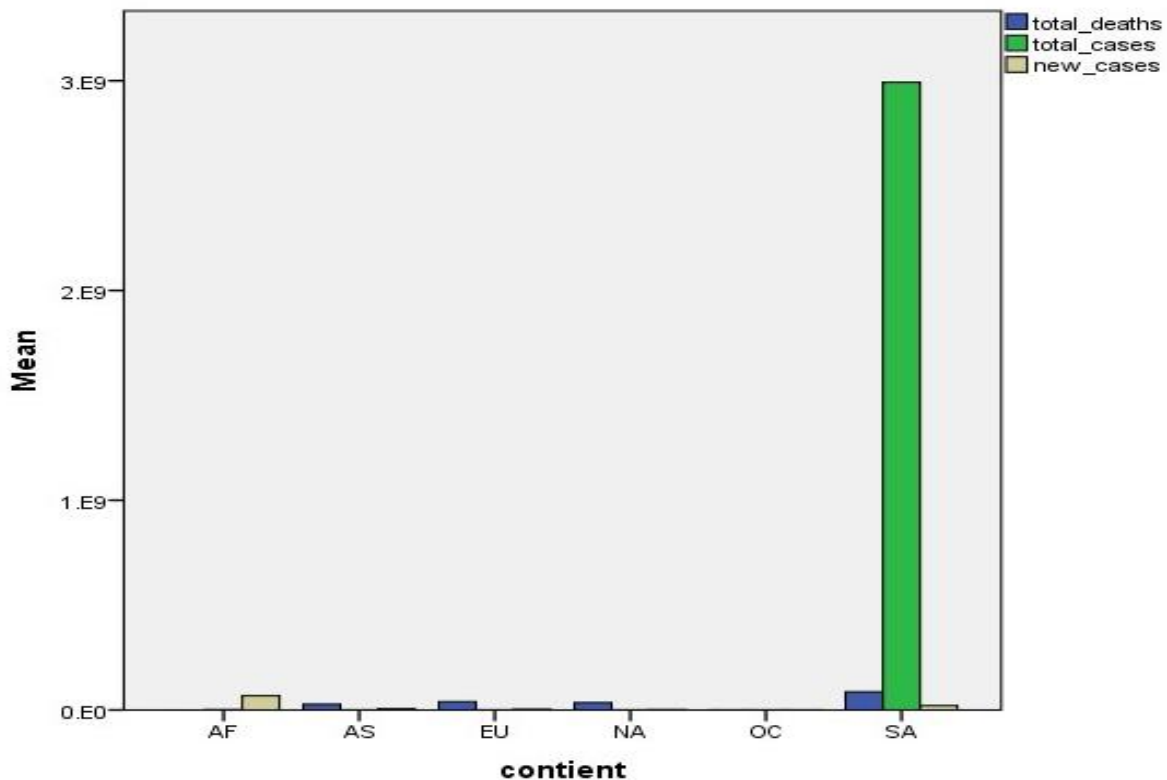
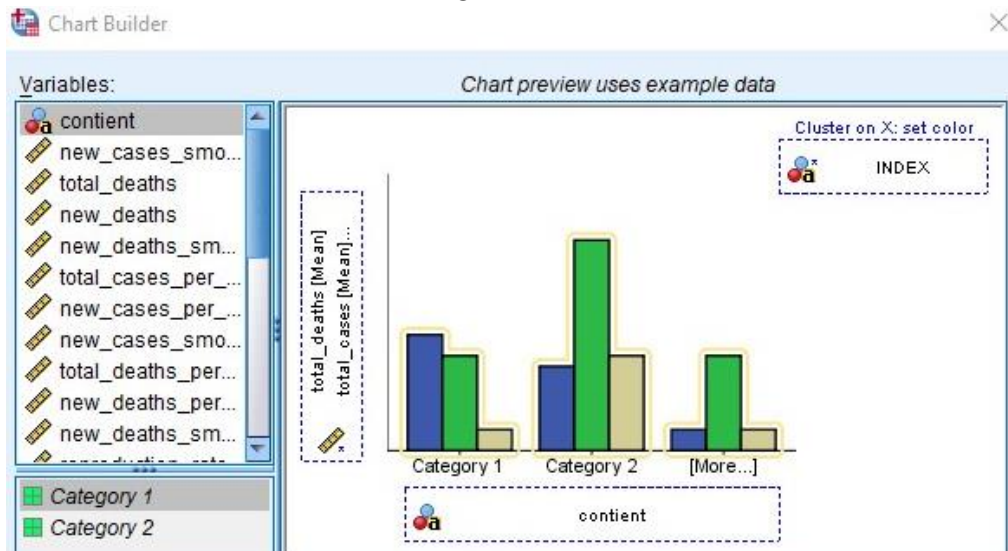
- Đồ thị thể hiện tổng số ca đã tử vong do covid-19 trên 6 lục địa



Nhận xét:

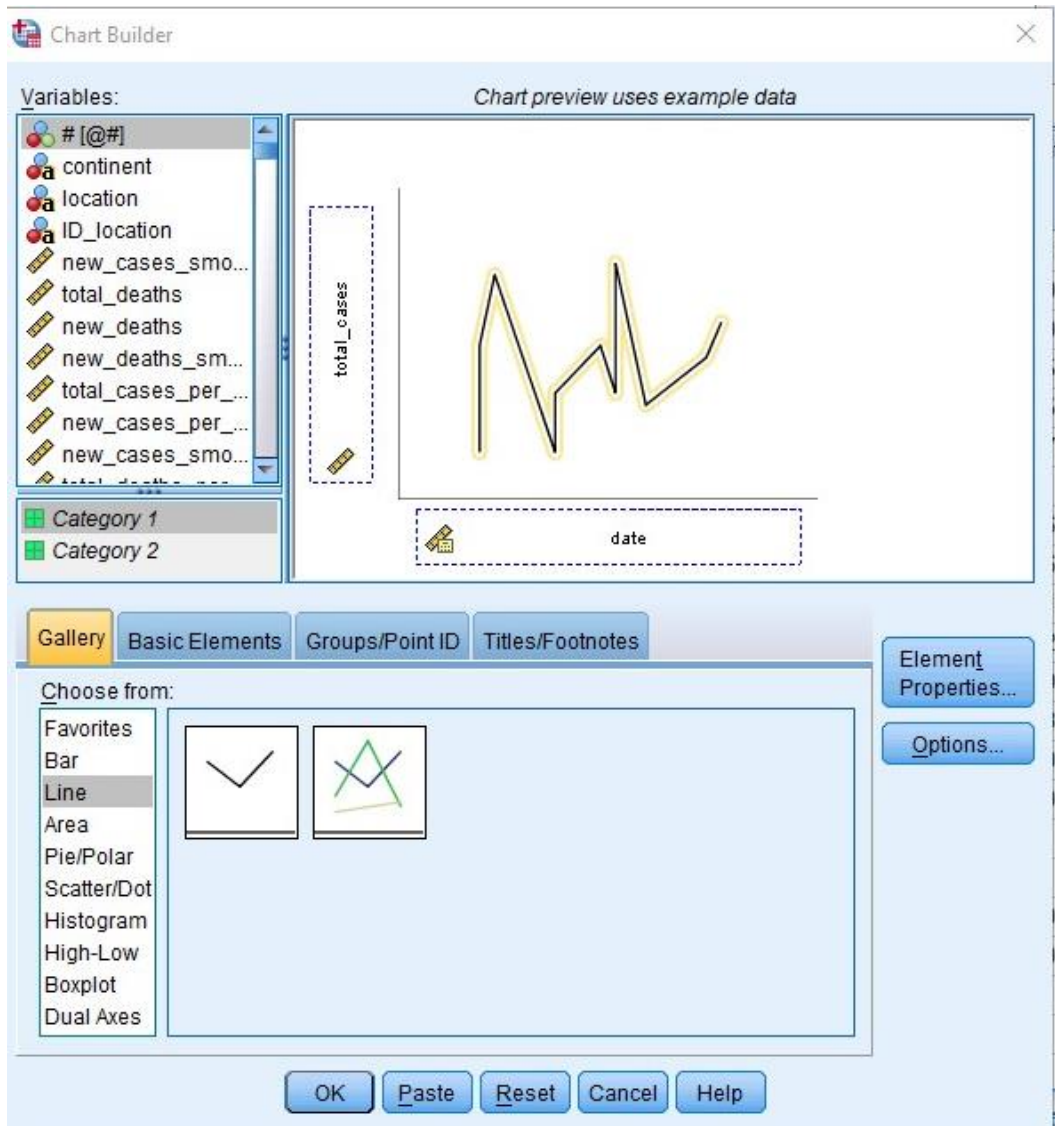
- Từ đồ thị trên ta thấy rằng, tổng số ca tử vong do covid 19 tập trung chủ yếu ở vùng SA(Bắc Mỹ)
- Hai vùng OC(Châu Đại Dương) và AF (Châu Phi) có số lượng người tử vong ít nhất.

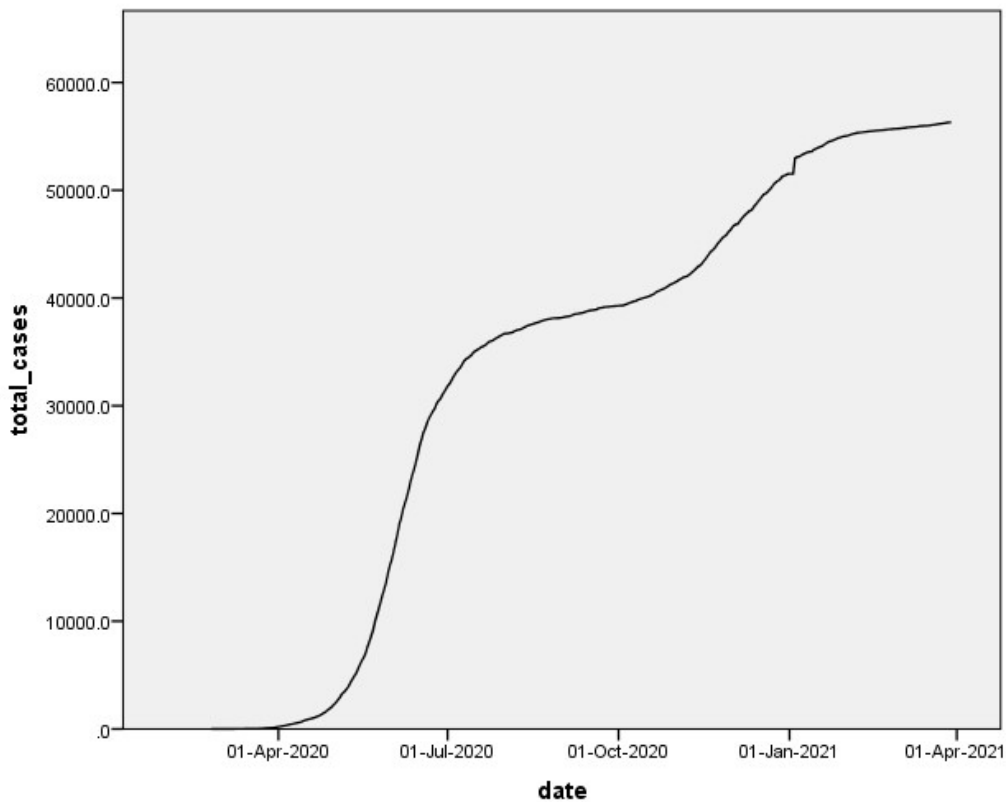
- Đồ thị thể hiện số ca mắc,tử vong và ca mắc mới do Covid-19



Nhận xét:

- Ta thấy tăng số ca mắc tập trung chủ yếu ở vùng SA(Nam Mỹ) Bên cạnh đó số ca tử vong do Covid-19 cũng nhiều nhất ở khu vực này.
 - Số ca mắc mới nhiều nhất ở khu vực AF (Châu Phi).
- Đồ thị thể hiện số ca mắc theo ngày của Châu Á





Nhận xét:

- Từ đồ thị ta thấy rằng, số ca mắc của Asia tăng trong khoảng thời gian từ tháng 2/2020 đến 4/2021
- Cả chỉ số tăng nhanh trong khoảng thời gian này.

3.3 Mô tả đồ thị bằng đại lượng thống kê

Dữ liệu sử dụng được thu thập từ ngày 7/2/2020-31/3/2021

- Bài này thực hiện trên biến Total_case (numeric) trong bảng Continents.
- Các đại lượng thống kê dùng cho bài này: Mean(), Std.deviation(), Minimum, Maximum.
- **Cách thực hiện:**
 - Chọn **Analyze > Descriptive Statistics > Descriptives**
 - Chọn Options, sau đó chọn các đại lượng cần tính.
 - Nhấn OK.

- Bảng Output

Descriptives

Descriptive Statistics					
	N	Minimum	Maximum	Mean	Std. Deviation
total_cases	6	29695.60000	2992518115	499090694.7	1221525129
Valid N (listwise)	6				

Phân tích:

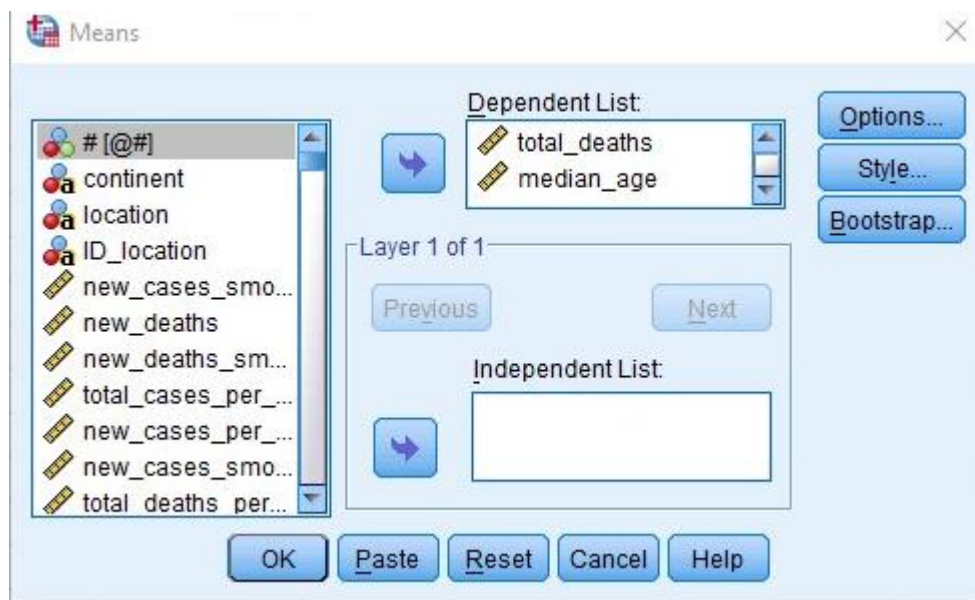
- Tổng số ca mắc ít nhất trong 6 khu vực trong khoảng thời gian thu thập số liệu là 29695 ca.
- Số ca mắc nhiều nhất được ghi nhận trên 6 khu vực trong khoảng thời gian thu thập là 499090694 ca.
- Độ lệch chuẩn của dữ liệu sau khi xử lý là 1221525129 ca.

4. Kiểm định

a. Kiểm định trung bình tổng thể

Cách thực hiện:

- Chọn **Analyze > Compare Means > Means**
- Đưa biến Total_deaths, median_age vào Dependent List
- Nhấn OK



- Kết quả Output:

Case Processing Summary

	Cases					
	Included		Excluded		Total	
	N	Percent	N	Percent	N	Percent
total_deaths	191	92.3%	16	7.7%	207	100.0%
median_age	205	99.0%	2	1.0%	207	100.0%

Report

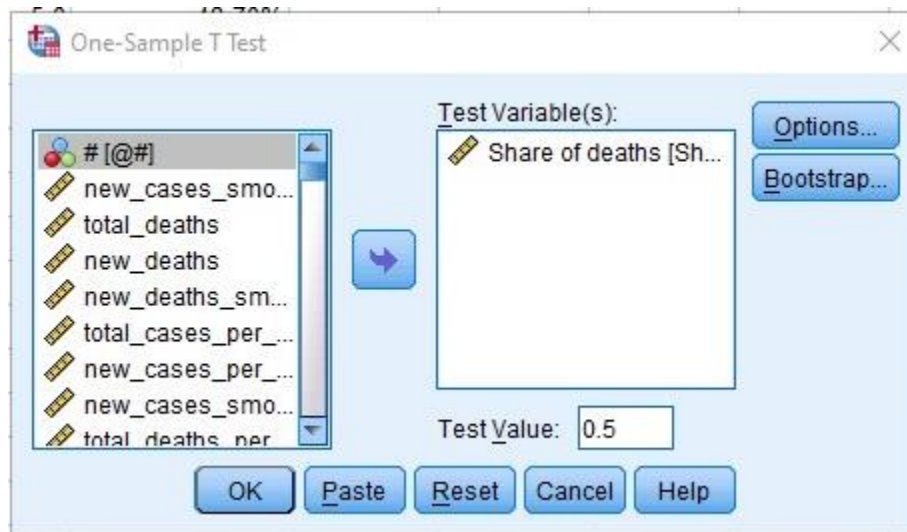
	total_deaths	median_age
Mean	2010936.89	2.95920E+10
N	191	205
Std. Deviation	14498178.28	2.19605E+11

b. Kiểm định 1 mẫu

Giả thuyết: Khả năng mắc bệnh Covid ở các độ tuổi là 50%

- **Cách thực hiện:**

- Chọn **Analyze > Compare Means > One-Sample T Test**
- Đưa biến Share of death vào Test Variables
- Nhập Test Value là 50%



- Bảng Output

One-Sample Statistics

	N	Mean	Std. Deviation	Std. Error Mean
Share of deaths	6	33.3267%	36.99912%	15.10483%

One-Sample Test

	Test Value = 0.5					
	t	df	Sig. (2-tailed)	Mean Difference	95% Confidence Interval of the Difference	
					Lower	Upper
Share of deaths	2.173	5	.082	32.82667%	-6.0015%	71.6549%

Phân tích: Quan sát trong bảng One-Sample Test, ta thấy giá trị Sig.(2-tailed) bằng 0.082, lớn hơn 0.05 . Vậy ta *chưa đủ cơ sở để bác bỏ giả thuyết* ban đầu, Khả năng mắc bệnh Covid ở các độ tuổi là 50%

c. Kiểm định mẫu cặp

Giả thuyết: Không có sự khác nhau giữa số ca mắc mới và tổng số ca tử vong do Covid

- **Cách thực hiện:**

- Chọn **Analyze > Compare Means > Paired Sample T Test.**

- Đưa biến new_case vào Pairs_Variables- variable1
- Đưa biến total_case vào Pairs_Variables-variable2
- Thực hiện kiểm định với độ tin cậy 95%(mức ý nghĩa 5%)



- Kết quả Output

Paired Samples Statistics

		Mean	N	Std. Deviation	Std. Error Mean
Pair 1	total_cases	2636.035909	44	784.9102138	118.3296670
	total_deaths	893160.727	44	1306648.799	196984.7182

Paired Samples Correlations

		N	Correlation	Sig.
Pair 1	total_cases & total_deaths	44	-.103	.507

Paired Samples Test

		Paired Differences					t	df	Sig. (2-tailed)
		Mean	Std. Deviation	Std. Error Mean	95% Confidence Interval of the Difference				
					Lower	Upper			
Pair 1	total_cases - total_deaths	-890524.691	1306729.600	196996.8993	-1287806.80	-493242.581	-4.521	43	.000

Phân tích: Từ bảng Pair Samples Test, ta thấy rằng chỉ số Sig.(2-tailed) = 0.00 chỉ số này bé hơn 0.05. Vì vậy ta có thể **bác bỏ giả thuyết trên**