# Data Science Project- Dashboard

# Title of Dashboard- "The Covid-19"

# Submitted by Jasmine Kaur

# Roll No- 102183047, Sub-group- 3co25

Original dataset before preprocessing:



| | Country/Region | Continent | Population | TotalCases | NewCases | TotalDeaths | NewDeaths | TotalRecovered | NewRecovered | ActiveCases | Serious,Critical | Tot Cases/1M Dea |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Country/Region | Continent | Population | TotalCases | NewCases | TotalDeaths | NewDeaths | TotalRecovered | NewRecovered | ActiveCases | Serious,Critical | Tot Cases/1M Dea |
| 2 | USA | North America | 331198130 | 5032179 | | 162804 | | 2576668 | | 2292707 | 18296 | 15194 |
| 3 | Brazil | South America | 212710692 | 2917562 | | 98644 | | 2047660 | | 771258 | 8318 | 13716 |
| 4 | India | Asia | 1381344997 | 2025409 | | 41638 | | 1377384 | | 606387 | 8944 | 1466 |
| 5 | Russia | Europe | 145940924 | 871894 | | 14606 | | 676357 | | 180931 | 2300 | 5974 |
| 6 | South Africa | Africa | 59381566 | 538184 | | 9604 | | 387316 | | 141264 | 539 | 9063 |
| 7 | Mexico | North America | 129066160 | 462690 | 6590 | 50517 | 819 | 308848 | 4140 | 103325 | 3987 | 3585 |
| 8 | Peru | South America | 33016319 | 455409 | | 20424 | | 310337 | | 124648 | 1426 | 13793 |
| 9 | Chile | South America | 19132514 | 366671 | | 9889 | | 340168 | | 16614 | 1358 | 19165 |
| 10 | Colombia | South America | 50936262 | 357710 | | 11939 | | 192355 | | 153416 | 1493 | 7023 |
| 11 | Spain | Europe | 46756648 | 354530 | | 28500 | | | | | 617 | 7582 |
| 12 | Iran | Asia | 84097623 | 320117 | | 17976 | | 277463 | | 24678 | 4156 | 3806 |
| 13 | UK | Europe | 67922029 | 308134 | | 46413 | | | | | 73 | 4537 |
| 14 | Saudi Arabia | Asia | 34865919 | 284226 | | 3055 | | 247089 | | 34082 | 1915 | 8152 |
| 15 | Pakistan | Asia | 221295851 | 281863 | | 6035 | | 256058 | | 19770 | 809 | 1274 |
| 16 | Bangladesh | Asia | 164851401 | 249651 | | 3306 | | 143824 | | 102521 | | 1514 |



| | vered | NewRecovered | ActiveCases | Serious,Critical | Tot Cases/1M pop | Deaths/1M pop | TotalTests | Tests/1M pop | WHO Region | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | vered | NewRecovered | ActiveCases | Serious,Critical | Tot Cases/1M pop | Deaths/1M pop | TotalTests | Tests/1M pop | WHO Region | | | |
| 2 | '6668 | | 2292707 | 18296 | 15194 | 492 | 63139605 | 190640 | Americas | | | |
| 3 | l7660 | | 771258 | 8318 | 13716 | 464 | 13206188 | 62085 | Americas | | | |
| 4 | '7384 | | 606387 | 8944 | 1466 | 30 | 22149351 | 16035 | South-EastAsia | | | |
| 5 | '6357 | | 180931 | 2300 | 5974 | 100 | 29716907 | 203623 | Europe | | | |
| 6 | l7316 | | 141264 | 539 | 9063 | 162 | 3149807 | 53044 | Africa | | | |
| 7 | l8848 | 4140 | 103325 | 3987 | 3585 | 391 | 1056915 | 8189 | Americas | | | |
| 8 | .0337 | | 124648 | 1426 | 13793 | 619 | 2493429 | 75521 | Americas | | | |
| 9 | l0168 | | 16614 | 1358 | 19165 | 517 | 1760615 | 92022 | Americas | | | |
| 10 | l2355 | | 153416 | 1493 | 7023 | 234 | 1801835 | 35374 | Americas | | | |
| 11 | | | | 617 | 7582 | 610 | 7064329 | 151087 | Europe | | | |
| 12 | '7463 | | 24678 | 4156 | 3806 | 214 | 2612763 | 31068 | EasternMediterranean | | | |
| 13 | | | | 73 | 4537 | 683 | 17515234 | 257873 | Europe | | | |
| 14 | l7089 | | 34082 | 1915 | 8152 | 88 | 3635705 | 104277 | EasternMediterranean | | | |
| 15 | i6058 | | 19770 | 809 | 1274 | 27 | 2058872 | 9304 | EasternMediterranean | | | |
| 16 | l3824 | | 102521 | | 1514 | 20 | 1225124 | 7432 | South-EastAsia | | | |

**R code for Data PreProcessing:**

```r
dataset<-read.csv(choose.files())

library(dplyr)

library(WriteXLS)

library('writexl')

install.packages('data.table')


install.packages('WriteXLS')

head(dataset)

ncol(dataset)

nrow(dataset)

colnames(dataset)

data.frame(dataset)

# we are beginning with data preprocessing



# 1) checking null values in the dataset


sum(is.na(dataset))

sum(is.na(dataset$NewCases))

sum(is.na(dataset$NewDeaths))

sum(is.na(dataset$NewRecovered))

#df2<-data.frame(dataset[, c("Country.Region","NewCases", "NewDeaths","NewRecovered")])

# these are features with high no of na values therefore we will split them into different dataset
```

sum(is.na(dataset$Population)) # since there is only one row with na value we can drop the row

dataset[-(is.na(dataset$Population)==TRUE),]


sum(is.na(dataset$TotalRecovered)) #replacing the nas with mean value

dataset$TotalRecovered[is.na(dataset$TotalRecovered)] <- mean(dataset$TotalRecovered, na.rm =TRUE)


sum(is.na(dataset$Deaths.1M.pop)) #replacing the nas with median value

dataset$Deaths.1M.pop[is.na(dataset$Deaths.1M.pop)] <- mean(dataset$Deaths.1M.pop, na.rm =TRUE)



# 2) checking anomlous (inf) values in dataset


sum(is.infinite(dataset$Deaths.1M.pop))


#replacing inf values with mean

dataset$Deaths.1M.pop[which(!is.finite(dataset$Deaths.1M.pop))]<-variable

#mean(dataset$Deaths...100.Recovered)

variable=mean(dataset$Deaths...100.Recovered) #saved 39.47385

#variable=39.47385





# 4) since the name of continents are containing '/' we have to process them

num<-grep("/",dataset$Continent,value=FALSE)

for ( i in num)

```
{

  print(dataset$Continent[i])

}

#since the slash is only in one continent

for ( i in num)

{ dataset$Continent[i]="Australia"

}
```



```
# 5) we want to know the total no of cases so far has been registered per total population of each continent

df5<-dataset %>%

  group_by(Continent)%>%

  summarise_at("Population",sum)


tot_cases<-dataset %>%

  group_by(Continent)%>%

  summarise_at("TotalCases",sum)
```
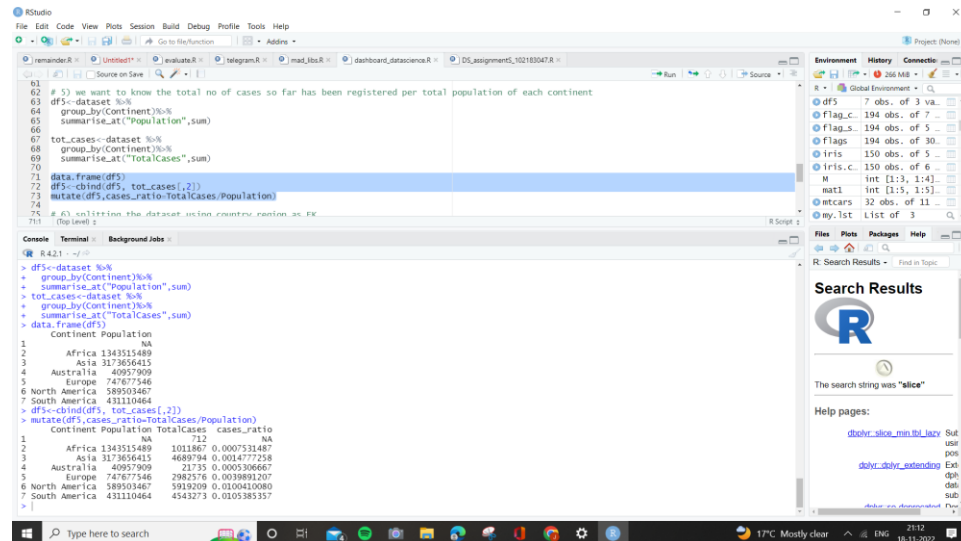
data.frame(df5)

df5<-cbind(df5, tot_cases[,2])

mutate(df5,cases_ratio=TotalCases/Population)



# 6) splitting the dataset using country.region as FK

df1<-data.frame(dataset[, c("Country.Region","TotalCases","TotalDeaths","TotalRecovered","ActiveCases")])

df1


df2<-data.frame(dataset[, c("Country.Region","NewCases", "NewDeaths","NewRecovered")])


df3<-data.frame(dataset[,c("Country.Region","Tot.Cases.1M.pop","Deaths.1M.pop","Tests.1M.pop")])


df4<-df5
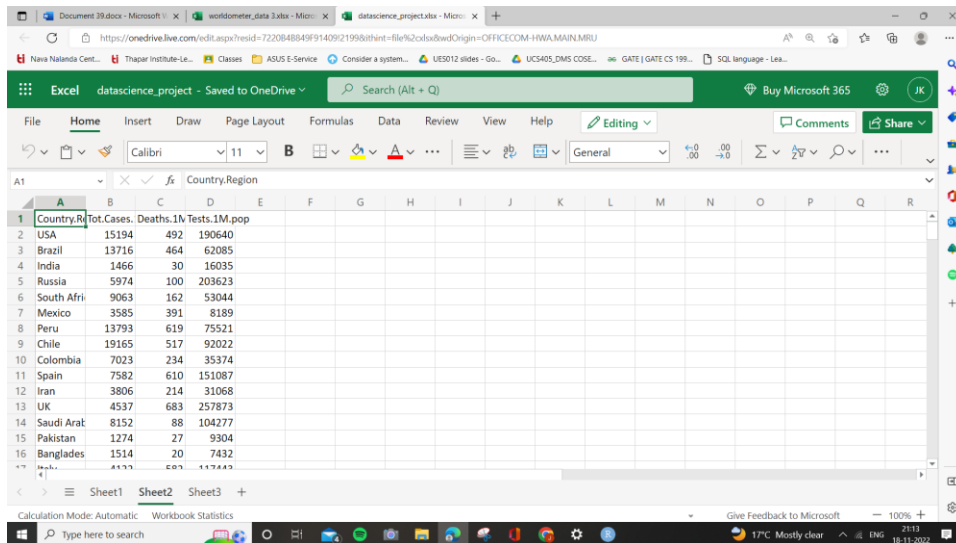
 #7) preprocessing of df2

df2

sum(is.na(df2)) #dropping

# 8) now are four datasets are preprocessed for tableau representation, so we will convert them into xlsx file

write.csv(df4, file="C:\\Users\\jaskirat singh\\Downloads\\datascienceproj\\dataset4.csv", row.names = FALSE)

write.csv(df1, file="C:\\Users\\jaskirat singh\\Downloads\\datascienceproj\\dataset1.csv", row.names = FALSE)

write.csv(df3, file="C:\\Users\\jaskirat singh\\Downloads\\datascienceproj\\dataset3.csv", row.names = FALSE)
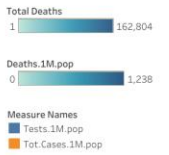
The final preprocessed data is:

## The dashboard:

### Covid-19 Data Visualization During Year( 2020-22)

Country-wise total no. of deaths,cases and recoveries



© 2022 Mapbox © OpenStreetMap

Country wise total deaths, tests conducted per million



Total Deaths

1 [            ] 162,804

Deaths.1M.pop

0 [            ] 1,238

Measure Names
■ Tests.1M.pop
■ Tot.Cases.1M.pop

Continent-wise covid cases along with population



Country-wise total no of tests conducted throught along with the total cases cases observed