

# COVID-19 Data Visualization

## Using Tableau and Spark

Team 1:

*Jacequai De Magalhaes, Brandon Cho, Jack Nguyen, Newyork Her, and Aaron Schomer*

# MVP

- Created a Spark Application, written in Scala, that uses the CLI for user interaction.
- Users can create an account and login to the program.
  - Accounts persists after sessions.
  - Passwords are automatically encrypted using **bCrypt API**.
  - Input validation for correct/present username and password.
  - Basic and Admin privileges. Admins have access to user management and log viewing.
- System will log each action done during sessions within a Spark SQL table.
  - Also set logging level to error to hide unnecessary log information during runtime
- Users are able to run queries on COVID data files provided by the program.
- Produced a .jar file for the project, allowing for the application to be ran in Ubuntu and AWS.
- Utilized Tableau and Zeppelin to visualize results and analyze trends.

# Stretch Goals

- Finding a relation between COVID data and other kinds of data, i.e., GDP and Population.
- Use Zeppelin for data visualization.
- Using AWS or Azure to run our .jar file.
- Visualizing Spatial Data in Tableau.

Here is our Demo!

# AWS EC2

The screenshot shows the AWS EC2 Management console interface. The browser address bar indicates the URL is `us-east-1.console.aws.amazon.com/ec2/v2/home?region=us-east-1#Instances:instanceState=running`. The left sidebar has a 'New EC2 Experience' toggle and sections for EC2 Dashboard, EC2 Global View, Events, Tags, Limits, Instances (selected), Images, and Feedback. The Instances section is expanded, showing 'Instances (1) Info'. A search bar at the top of the main content area contains the filter 'Instance state = running'. Below it is a table with one row, where the 'Name' column value 'project2' is highlighted with a red circle. The table columns are Name, Instance ID, Instance state, Instance type, and Status check. The instance details are: Name: project2, Instance ID: i-0b74a35668b7ea97c, Instance state: Running, Instance type: t2.large, Status check: 2/2 checks passed. A modal window titled 'Select an instance' is open at the bottom.

Name	Instance ID	Instance state	Instance type	Status check
project2	i-0b74a35668b7ea97c	Running	t2.large	2/2 checks passed

Select an instance

# AWS EC2

The screenshot shows the AWS EC2 Instance Details page for an instance named **i-0b74a35668b7ea97c (project2)**. The instance is currently **Running**.

**Instance summary for i-0b74a35668b7ea97c (project2) [Info](#)**

Updated less than a minute ago

Value	Description	Value	Description
Instance ID	Public IPv4 address	i-0b74a35668b7ea97c (project2)	3.91.215.172   <a href="#">open address</a>
IPv6 address	Instance state	-	Running
Hostname type	Private IP DNS name (IPv4 only)	IP name: ip-172-31-20-91.ec2.internal	ip-172-31-20-91.ec2.internal
Instance type	Elastic IP addresses	t2.large	-
AWS Compute Optimizer finding	IAM Role	Opt-in to AWS Compute Optimizer for recommendations.   <a href="#">Learn more</a>	-
Auto Scaling Group name	Subnet ID	-	subnet-063c66e1faf4a1297

**EC2 Dashboard** | **EC2 Global View** | **Events** | **Tags** | **Limits**

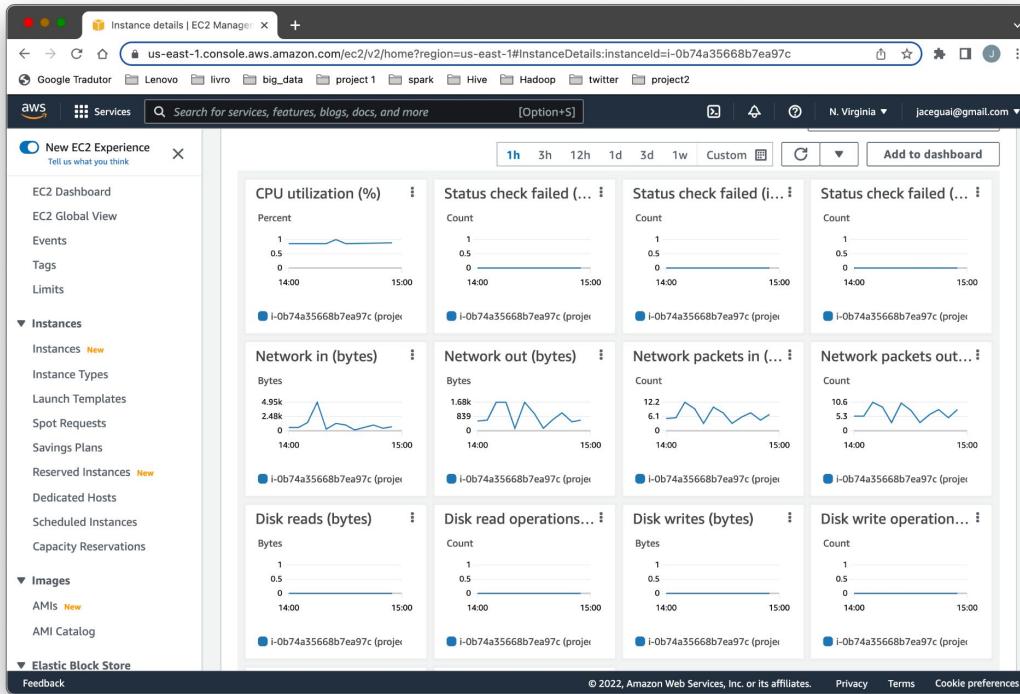
**Instances** [New](#) | **Instance Types** | **Launch Templates** | **Spot Requests** | **Savings Plans** | **Reserved Instances** [New](#) | **Dedicated Hosts** | **Scheduled Instances** | **Capacity Reservations**

**Images** [AMIs](#) [New](#)

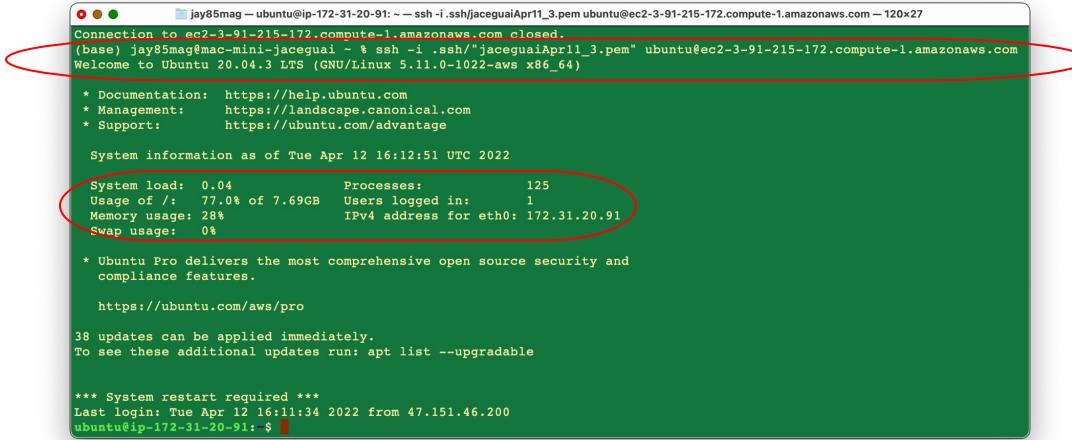
[Feedback](#) | [Privacy](#) | [Terms](#) | [Cookie preferences](#)

© 2022, Amazon Web Services, Inc. or its affiliates.

# AWS EC2



# AWS EC2



jay85mag — ubuntu@ip-172-31-20-91: ~ ssh -i .ssh/jacegualApr11\_3.pem ubuntu@ec2-3-91-215-172.compute-1.amazonaws.com - 120x27

Connection to ec2-3-91-215-172.compute-1.amazonaws.com closed.

(base) jay85mag@mac-mini-jacegual ~ % ssh -i .ssh/"jacegualApr11\_3.pem" ubuntu@ec2-3-91-215-172.compute-1.amazonaws.com

Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.11.0-1022-aws x86\_64)

\* Documentation: <https://help.ubuntu.com>  
\* Management: <https://landscape.canonical.com>  
\* Support: <https://ubuntu.com/advantage>

System information as of Tue Apr 12 16:12:51 UTC 2022

System load:	Processes:
0.04	125

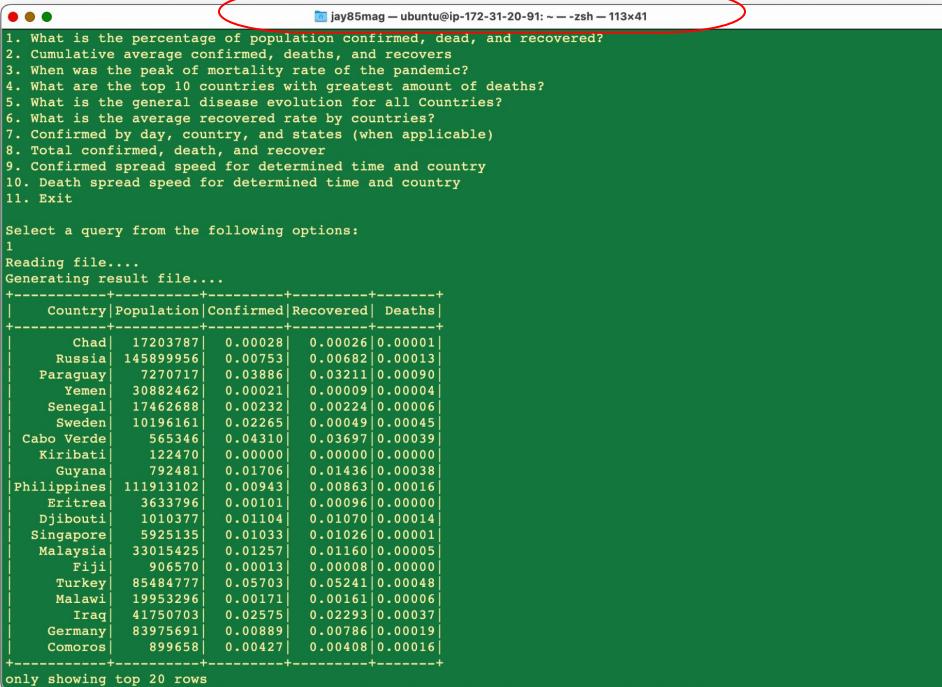
Usage of /: 77.0% of 7.69GB Users logged in: 1  
Memory usage: 28% IPv4 address for eth0: 172.31.20.91  
Swap usage: 0%

\* Ubuntu Pro delivers the most comprehensive open source security and compliance features.  
<https://ubuntu.com/aws/pro>

38 updates can be applied immediately.  
To see these additional updates run: apt list --upgradable

\*\*\* System restart required \*\*\*  
Last login: Tue Apr 12 16:11:34 2022 from 47.151.46.200  
ubuntu@ip-172-31-20-91: \$

# AWS EC2



jay85mag — ubuntu@ip-172-31-20-91: ~ - zsh - 113x41

1. What is the percentage of population confirmed, dead, and recovered?
2. Cumulative average confirmed, deaths, and recoveries
3. When was the peak of mortality rate of the pandemic?
4. What are the top 10 countries with greatest amount of deaths?
5. What is the general disease evolution for all Countries?
6. What is the average recovered rate by countries?
7. Confirmed by day, country, and states (when applicable)
8. Total confirmed, death, and recover
9. Confirmed spread speed for determined time and country
10. Death spread speed for determined time and country
11. Exit

Select a query from the following options:

1

Reading file....

Generating result file....

Country	Population	Confirmed	Recovered	Deaths
Chad	17203787	0.00028	0.00026	0.00001
Russia	145899956	0.00753	0.00682	0.00013
Paraguay	7270717	0.03886	0.03211	0.00090
Yemen	30882462	0.00021	0.00009	0.00004
Senegal	17462688	0.00232	0.00224	0.00006
Sweden	10196161	0.02265	0.00049	0.00045
Cabo Verde	565346	0.04310	0.03697	0.00039
Kiribati	122470	0.00000	0.00000	0.00000
Guyana	792481	0.01706	0.01436	0.00038
Philippines	111913102	0.00943	0.00863	0.00016
Eritrea	3633796	0.00101	0.00096	0.00000
Djibouti	1010377	0.01104	0.01070	0.00014
Singapore	5925135	0.01033	0.01026	0.00001
Malaysia	33015425	0.01257	0.01160	0.00005
Fiji	906570	0.00013	0.00008	0.00000
Turkey	85484777	0.05703	0.05241	0.00048
Malawi	19953296	0.00171	0.00161	0.00006
Iraq	41750703	0.02575	0.02293	0.00037
Germany	83975691	0.00889	0.00786	0.00019
Comoros	899658	0.00427	0.00408	0.00016

only showing top 20 rows

# AWS EC2 – Connecting

```
root@DESKTOP-NC309U:~# ssh -i /home/brandon/jaceguaiApril11_3.pem ubuntu@ec2-3-91-215-172.compute-1.amazonaws.com
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.11.0-1022-aws x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/advantage

System information as of Tue Apr 12 16:15:16 UTC 2022

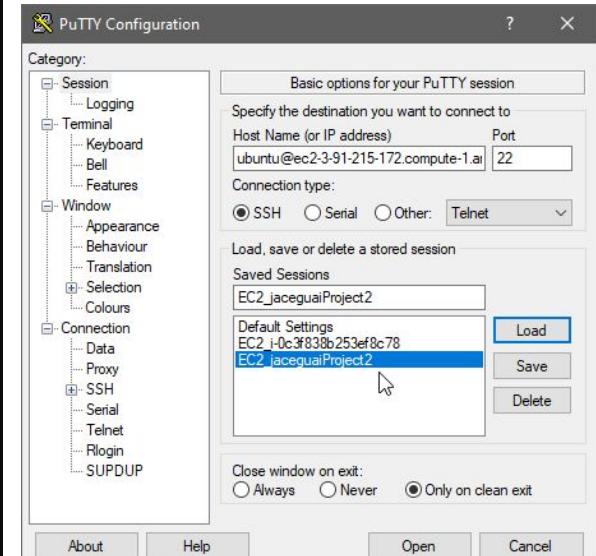
System load:  0.04      Processes:          124
Usage of /:   77.0% of 7.69GB  Users logged in:     0
Memory usage: 28%           IPv4 address for eth0: 172.31.20.91
Swap usage:   0%

* Ubuntu Pro delivers the most comprehensive open source security and
  compliance features.

  https://ubuntu.com/aws/pro

38 updates can be applied immediately.
To see these additional updates run: apt list --upgradable

*** System restart required ***
Last login: Tue Apr 12 16:15:07 2022 from 47.151.46.200
ubuntu@ip-172-31-20-91:~$ ls
apache-hive-3.1.2-bin apache-hive-3.1.2-bin.tar.gz derby.log hadoop hadoop-3.3.2.tar.gz metastore_db project2 resources spark spark-3.2.1-bin-hadoop3.2.tgz
ubuntu@ip-172-31-20-91:~$ jps
52210 Jps
6819 SecondaryNameNode
6579 DataNode
6422 NameNode
7164 NodeManager
7085 ResourceManager
ubuntu@ip-172-31-20-91:~$ sudo su - brandon
brandon@ip-172-31-20-91:~$ jps
52236 Jps
brandon@ip-172-31-20-91:~$ logout
ubuntu@ip-172-31-20-91:~$ |
```



# AWS EC2 - Connecting

```
ubuntu@ip-172-31-20-91:~  
Using username "ubuntu".  
Authenticating with public key "imported-openssh-key"  
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.11.0-1022-aws x86_64)  
* Documentation: https://help.ubuntu.com  
* Management: https://landscape.canonical.com  
* Support: https://ubuntu.com/advantage  
  
System information as of Tue Apr 12 16:38:39 UTC 2022  
  
System load: 0.0 Processes: 131  
Usage of /: 77.0% of 7.69GB Users logged in: 1  
Memory usage: 28% IPv4 address for eth0: 172.31.20.91  
Swap usage: 0%  
  
* Ubuntu Pro delivers the most comprehensive open source security and  
compliance features.  
  
https://ubuntu.com/aws/pro  
  
38 updates can be applied immediately.  
To see these additional updates run: apt list --upgradable  
  
*** System restart required ***  
Last login: Tue Apr 12 16:35:26 2022 from 47.151.46.200  
ubuntu@ip-172-31-20-91:~$  
  
ubuntu@ip-172-31-20-91:~  
Using username "ubuntu".  
Authenticating with public key "imported-openssh-key"  
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.11.0-1022-aws x86_64)  
* Documentation: https://help.ubuntu.com  
* Management: https://landscape.canonical.com  
* Support: https://ubuntu.com/advantage  
  
System information as of Tue Apr 12 16:39:02 UTC 2022  
  
System load: 0.0 Processes: 136  
Usage of /: 77.0% of 7.69GB Users logged in: 1  
Memory usage: 28% IPv4 address for eth0: 172.31.20.91  
Swap usage: 0%  
  
* Ubuntu Pro delivers the most comprehensive open source security and  
compliance features.  
  
https://ubuntu.com/aws/pro  
  
38 updates can be applied immediately.  
To see these additional updates run: apt list --upgradable  
  
*** System restart required ***  
Last login: Tue Apr 12 16:39:40 2022 from 47.151.46.200  
ubuntu@ip-172-31-20-91:~$ users  
ubuntu ubuntu ubuntu ubuntu ubuntu  
ubuntu@ip-172-31-20-91:~$ ls  
ubuntu ubuntu ubuntu ubuntu ubuntu  
ubuntu@ip-172-31-20-91:~$  
  
ubuntu@ip-172-31-20-91:~  
Using username "ubuntu".  
Authenticating with public key "imported-openssh-key"  
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.11.0-1022-aws x86_64)  
* Documentation: https://help.ubuntu.com  
* Management: https://landscape.canonical.com  
* Support: https://ubuntu.com/advantage  
  
System information as of Tue Apr 12 16:35:25 UTC 2022  
  
System load: 0.02 Processes: 130  
Usage of /: 77.0% of 7.69GB Users logged in: 1  
Memory usage: 28% IPv4 address for eth0: 172.31.20.91  
Swap usage: 0%  
  
* Ubuntu Pro delivers the most comprehensive open source security and  
compliance features.  
  
https://ubuntu.com/aws/pro  
  
38 updates can be applied immediately.  
To see these additional updates run: apt list --upgradable  
  
*** System restart required ***  
Last login: Tue Apr 12 16:35:26 2022 from 47.151.46.200  
ubuntu@ip-172-31-20-91:~$  
  
Administrator PowerShell x ubuntu@ip-172-31-20-91:~ + v  
brandon@DESKTOP-NOE399U:~$ chmod 700 /home/brandon/Keys/AWS  
brandon@DESKTOP-NOE399U:~$ ssh -i /home/brandon/Keys/AWS/jaceguaiApr11_3.pem ubuntu@ec2-3-91-215-172.compute-1.amazonaws.com  
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.11.0-1022-aws x86_64)  
  
* Documentation: https://help.ubuntu.com  
* Management: https://landscape.canonical.com  
* Support: https://ubuntu.com/advantage  
  
System information as of Tue Apr 12 16:39:54 UTC 2022  
  
System load: 0.0 Processes: 139  
Usage of /: 77.0% of 7.69GB Users logged in: 1  
Memory usage: 28% IPv4 address for eth0: 172.31.20.91  
Swap usage: 0%  
  
* Ubuntu Pro delivers the most comprehensive open source security and  
compliance features.  
  
https://ubuntu.com/aws/pro  
  
38 updates can be applied immediately.  
To see these additional updates run: apt list --upgradable  
  
*** System restart required ***  
Last login: Tue Apr 12 16:39:03 2022 from 47.151.46.200  
ubuntu@ip-172-31-20-91:~$ users  
ubuntu ubuntu ubuntu ubuntu ubuntu  
ubuntu@ip-172-31-20-91:~$ !
```

# Data Cleaning/Sanitization

```
1+ ##### Initial setup #####
2 # Checks current working directory
3 getwd()
4
5 # Sets working directory
6 # Change the path to where your files are located
7 setwd("C:\\Users\\brany\\Documents\\Revature\\Training\\Projects\\P2\\Data_Set2")
8
9
10+ ##### The files we are cleaning/changing #####
11 # Covid-19 data file
12 covid_19_data = read.csv("covid_19_data.csv", header = F)
13
14 # Times Series files
15 ts_c19_confirmed = read.csv("time_series_covid_19_confirmed.csv", header = F)
16 ts_c19_confirmed_us = read.csv("time_series_covid_19_confirmed_us.csv", header = F)
17 ts_c19_deaths = read.csv("time_series_covid_19_deaths.csv", header = F)
18 ts_c19_deaths_us = read.csv("time_series_covid_19_deaths_us.csv", header = F)
19 ts_c19_recovered = read.csv("time_series_covid_19_recovered.csv", header = F)
20
21
22
23
24+ ##### Functions #####
25 # Convert date to format that is compatible with Hive (ONLY FOR TIMES SERIES)
26 # Potential formats: M/D/Y, MM/D/Y, MM/DD/YY, MM/DD/YYYY
27 ts_olddate_to_newDate <- function(olddate) {
28   tempDate <- strsplit(olddate, split = "/")
29   if (nchar(olddate) == 6) {
30     newDate <- paste("20", tempDate[[1]][[3]], "-0", tempDate[[1]][[1]], "-0", tempDate[[1]][[2]], sep = "")
31   }
32   else if (nchar(olddate) == 7) {
33     if (nchar(tempDate[[1]][[1]]) == 1) {
34       newDate <- paste("20", tempDate[[1]][[3]], "-0", tempDate[[1]][[1]], "-", tempDate[[1]][[2]], sep = "")
35     }
36     else if (nchar(tempDate[[1]][[1]]) == 2) {
37       newDate <- paste("20", tempDate[[1]][[3]], "-", tempDate[[1]][[1]], "-0", tempDate[[1]][[2]], sep = "")
38     }
39   }
40   else if (nchar(olddate) == 8) {
41     newDate <- paste("20", tempDate[[1]][[3]], "-", tempDate[[1]][[1]], "-", tempDate[[1]][[2]], sep = "")
42   }
43   else {
44     return(olddate)
45   }
46   return(newDate)
47 }
48
49 # same as above, but works for covid_19_data dates format MM/DD/YYYY
50 olddate_to_newDate <- function(olddate) {
51   tempDate <- strsplit(olddate, split = "/")
52   newDate <- paste(tempDate[[1]][[3]], "-", tempDate[[1]][[1]], "-", tempDate[[1]][[2]], sep = "")
53   return(newDate)
54 }
55
56+
57+ ##### Cleaning Data #####
58+
59 ## Create csv that is a transpose of time series csv, one column is date, the rest are the countries
60 # Add up all Provinces/Regions so there is only 1 country.
61
62
63 # ts_19_confirmed
64 # creates a transpose of data
65 transpose1 = t(ts_c19_confirmed)
66 # removes columns, Province/State, Lat, Long
67 remCol1 = transpose1[,c(1,3,4),]
68 # fixes row numbers after deleting rows
69 row.names(remCol1) <- NULL
70 # creates a list of new dates from old dates
71 newDateList = lapply(remCol1[1:nrow(remCol1)], ts_olddate_to_newDate)
72
73 # All numbers are stored as characters, so we need to convert to numeric
74 df_clean = remCol1[, -1]
75 colnames(df_clean) <- NULL
76 # now we have a transposed table with country then the rest being the data
77 transpose2 = data.frame(t(df_clean))
78
79 # Fixing South Korea because it's formatted like: Korea, South
80 # This causes problems later
81 # Following statement gives us the row number and we already know the column
82 rowNum = which(apply(transpose2, 1, function(x) any(x %in% "Korea, South")))
83 transpose2[rowNum, 1] = "Korea"
84
85 # the columns we want to modify
86 i = c(2:ncol(transpose2))
87 # converts the columns to numeric
88 transpose2[, i] = apply(transpose2[, i], 2, function(x) as.numeric(as.character(x)))
89
90 # convert into data.table
91 #install.packages("data.table", dependencies=TRUE)
92 library(data.table)
93 df_table = data.table(transpose2)
94 # sum up values by country
95 df_sum = df_table[, lapply(.SD, sum), by=x1]
96
97 # Set up the two data frames and combine
98 finalCountries = data.frame(t(df_sum))
99 finalDate = data.frame(t(data.frame(newDateList)))
100 final = cbind(finalDate, finalCountries)
101 row.names(final) <- NULL
102 colnames(final) <- NULL
103 # change Country/Region to Date
104 final[1,1] = "Date"
105
106 write.table(final, "C:\\Users\\brany\\Documents\\Revature\\Training\\Projects\\P2\\Test\\ts_confirmed_by_country.csv")
```

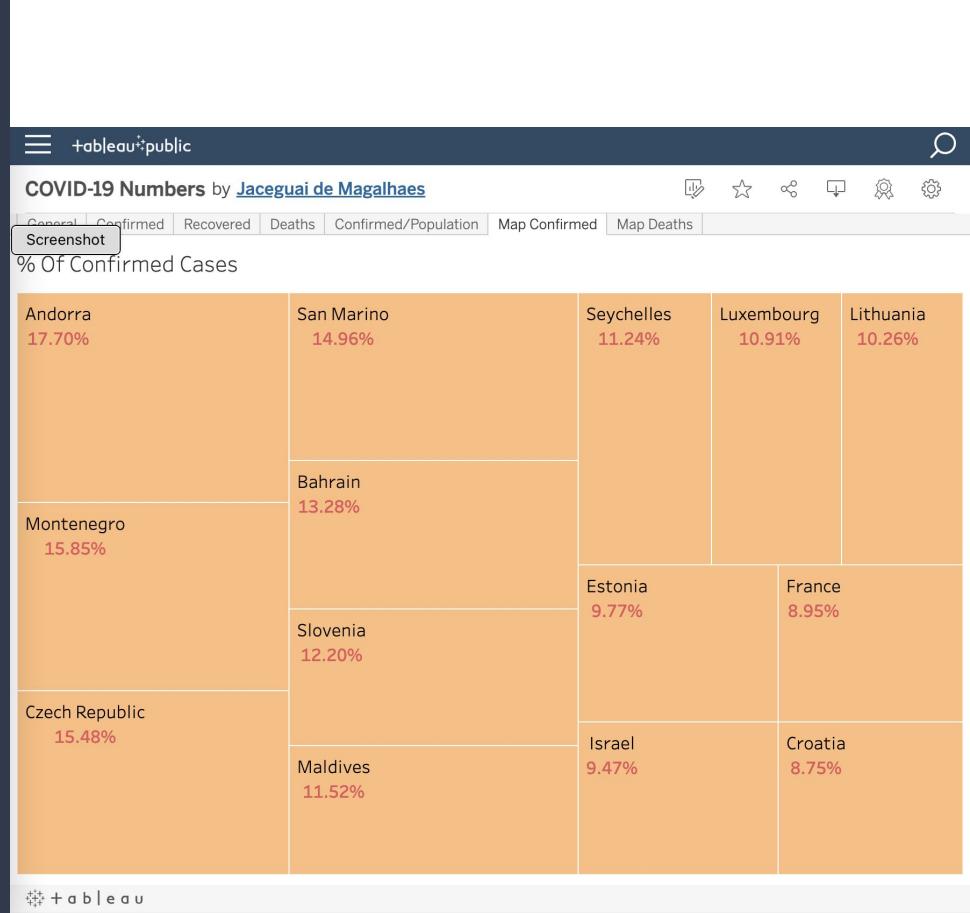
# Data Cleaning/Sanitization

```
225 ## Create csv of covid_19_data, change date to Hive Format, remove Last update (optional), remove decimals for last 3 columns
226 temp = t(covid_19_data)
227 # Create a list of converted dates
228 newList = lapply(temp[1:nrow(temp),2], oldDate_to_newDate)
229
230 # Gets rid of unnecessary decimal by converting to integer
231 newConfirmedList = lapply(temp[1:nrow(temp),6], as.integer)
232 # as.integer doesn't know what to do with "Confirmed" so it gets converted to NA. Adding it back here
233 newConfirmedList[[1]] = temp[[1,6]]
234 # Do the same for Deaths/Recovered
235 newDeathsList = lapply(temp[1:nrow(temp),7], as.integer)
236 newDeathsList[[1]] = temp[[1,7]]
237 newRecoveredList = lapply(temp[1:nrow(temp),8], as.integer)
238 newRecoveredList[[1]] = temp[[1,8]]
239
240 # If we store a list as a dataframe, it will appear as multiple columns instead of 1 column
241 # observationDate, column 2
242 finalDate = data.frame(t(data.frame(newList)))
243 # Confirmed, column 6
244 finalConfirmed = data.frame(t(data.frame(newConfirmedList)))
245 # Deaths, column 7
246 finalDeaths = data.frame(t(data.frame(newDeathsList)))
247 # Recovered, column 8
248 finalRecovered = data.frame(t(data.frame(newRecoveredList)))
249
250 # Store the columns that we didn't alter
251 finalBeginning = data.frame(temp[,1])
252 finalMiddle = data.frame(temp[,c(3,4,5)])
253
254 # Bind everything together
255 final = cbind(finalBeginning, finalDate, finalMiddle, finalConfirmed, finalDeaths, finalRecovered)
256 #row.names(final) <- NULL
257 #colnames(final) <- NULL Don't need these, just to make final row/column headers look nice but we don't save them to csv
258
259 # Save to csv
260 write.table(final, "C:\\\\users\\\\brany\\\\documents\\\\Revature\\\\Training\\\\Projects\\\\P2\\\\Test\\\\covid_19_data_clean.csv", sep = ",", ri
```

Environment		History	Connections	Tutorial
	Import Dataset	899 MB		
R	Global Environment			
Data				
①	covid_19_data	285308 obs. of 8 variables		
①	finalBeginning	285308 obs. of 1 variable		
①	finalConfirmed	285308 obs. of 1 variable		
①	finalDate	285308 obs. of 1 variable		
①	finalDeaths	285308 obs. of 1 variable		
①	finalMiddle	285308 obs. of 3 variables		
①	finalRecovered	285308 obs. of 1 variable		
①	newConfirmedList	Large List (285308 elements, 18.3 MB)		
①	newDateList	Large List (285308 elements, 36.5 MB)		
①	newDeathsList	Large List (285308 elements, 18.3 MB)		
①	newRecoveredList	Large List (285308 elements, 18.3 MB)		
①	temp	Large matrix (2282464 elements, 42.4 MB)		
values				
	tdate1	"1/7/20"		
	tdate2	"1/22/20"		
	tdate3	"11/2/21"		
	tdate4	"11/23/20"		
	tdate5	"01/22/2020"		
	tdate6	"05/03/2020"		
Functions				
	oldDate_to_newDa...	function (oldDate)		
	ts_oldDate_to_ne...	function (oldDate)		

# Q1: What % of the global population makes up the confirmed, death, and recovered cases?

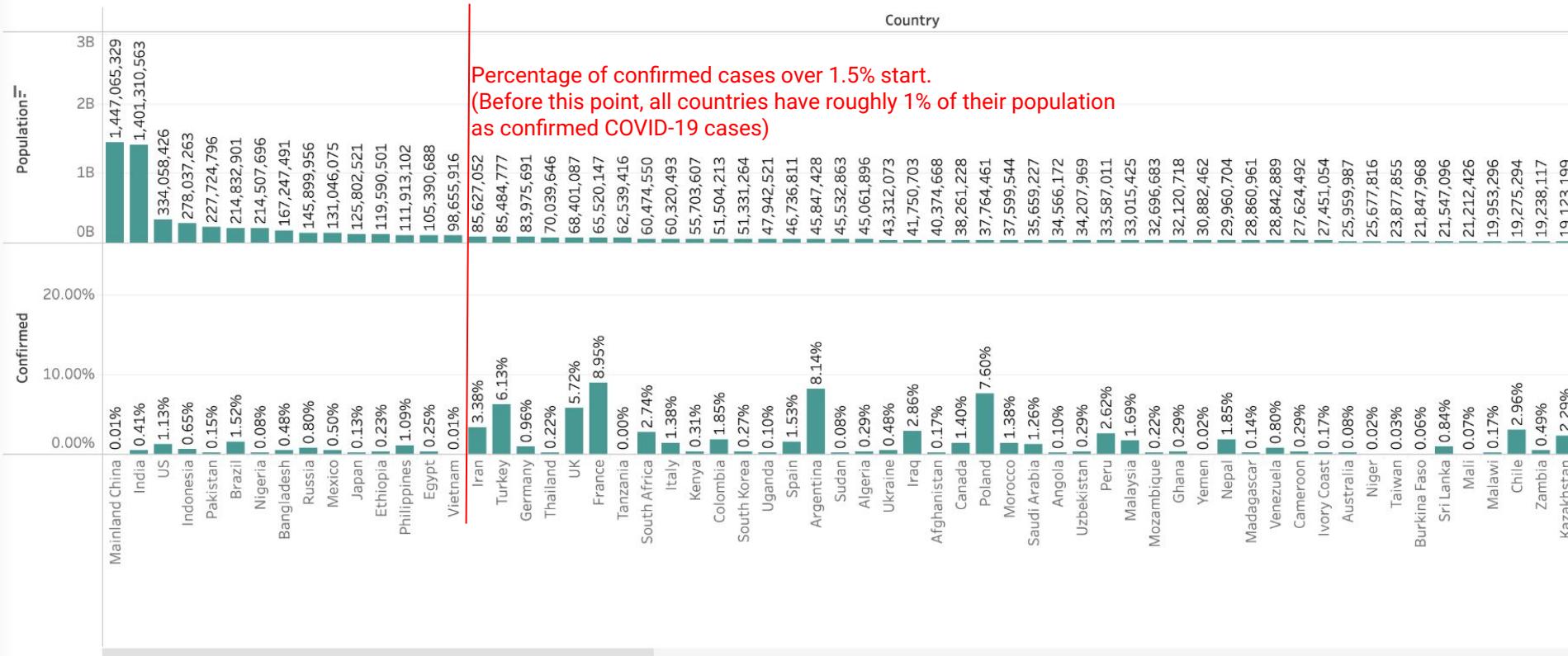
- It looks like the countries with the highest percentage of Confirmed Cases are not among the most populous countries.



## COVID-19 Numbers

Screenshot	Confirmed	Recovered	Deaths	Confirmed/Population	Map Confirmed	Map Deaths
------------	-----------	-----------	--------	----------------------	---------------	------------

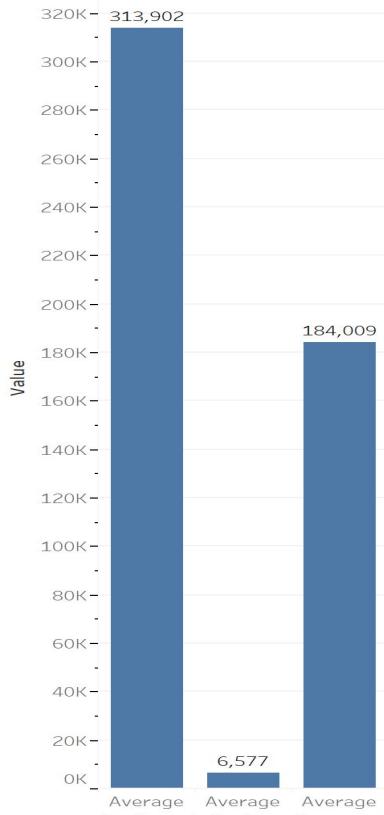
### Confirmed/Population



# Q2: What is the cumulative average of confirmed cases, death cases, and recovered cases?

- It seems that the total cumulative average Confirmed Cases throughout the world does not reach even half a million people. While the cumulative average deaths is not even 5% of the total average confirmed cases. Also the average recovery rate is over 50% of the confirmed cases. It means that a person is far more likely to survive an infection, than he is to die from it.
- Average cases per day (World Wide)
- Date: 05/02/2021

Query 2: Average  
Confirmed, Deaths, and  
Recoveries

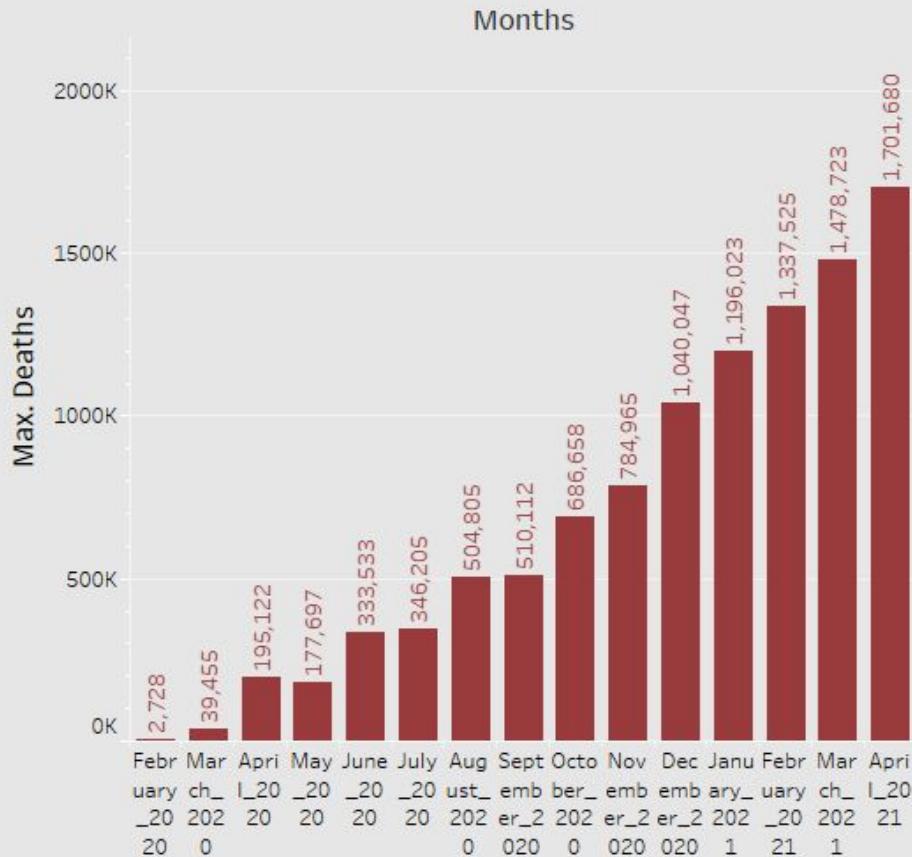


Average Confirmed, Average Deaths and  
Average Recovered.

# Q3: When was the peak of mortality rate during the pandemic?

- It's interesting to see that the peak of the pandemic is actually at the end of the last recorded month. We hypothesize that the peak mortality rate during the pandemic would be closer to the middle of the pandemic, but turns out that the mortality rate increases each month as it goes on

## Query 3: Peak Mortality Rate of the Pandemic

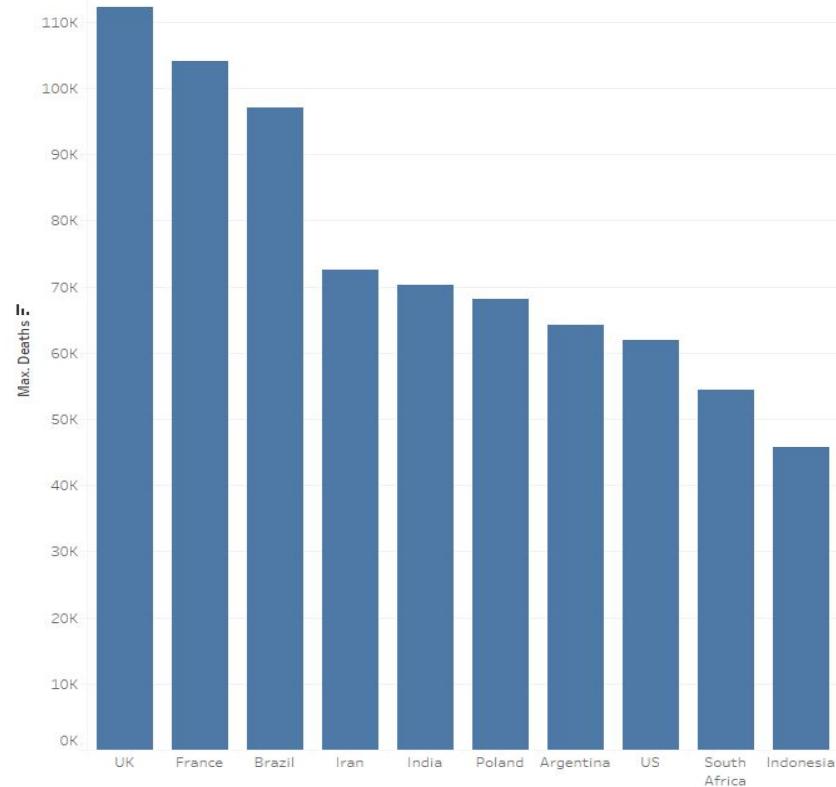


# Q4: What are the top 10 countries with the greatest amount of deaths?

- Even though the source of the pandemic is from Asia. The death rate is increasingly more on other countries around the world.

Sheet 3

Country/Region



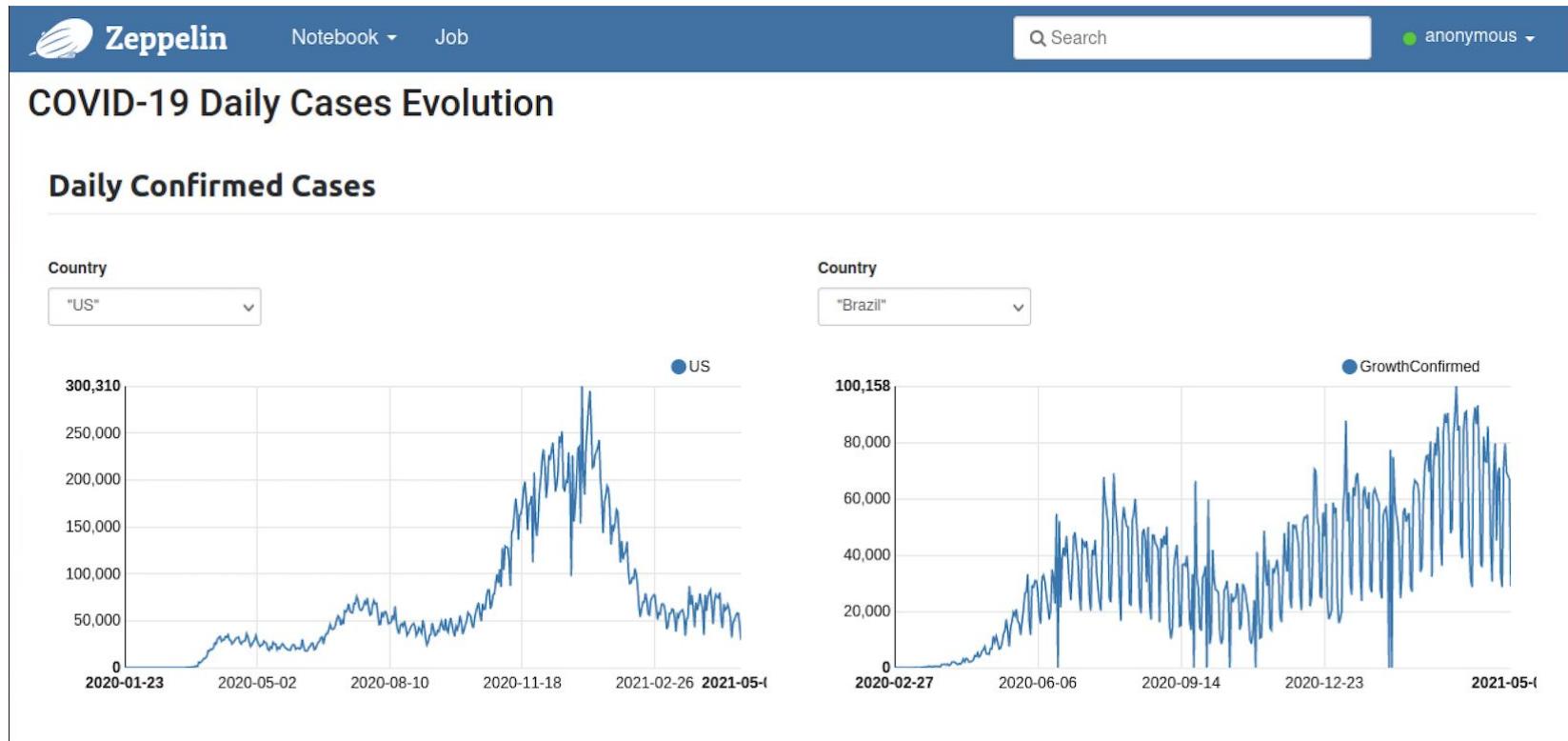
Maximum of Deaths for each Country/Region. The data is filtered on Calculation3, which keeps True.

# Q5: What is the general disease evolution for all countries?

For this query, we built a Zeppelin visualization where the user can select two countries side by side and compare the daily evolution of confirmed, recovered, and dead COVID-19 cases.

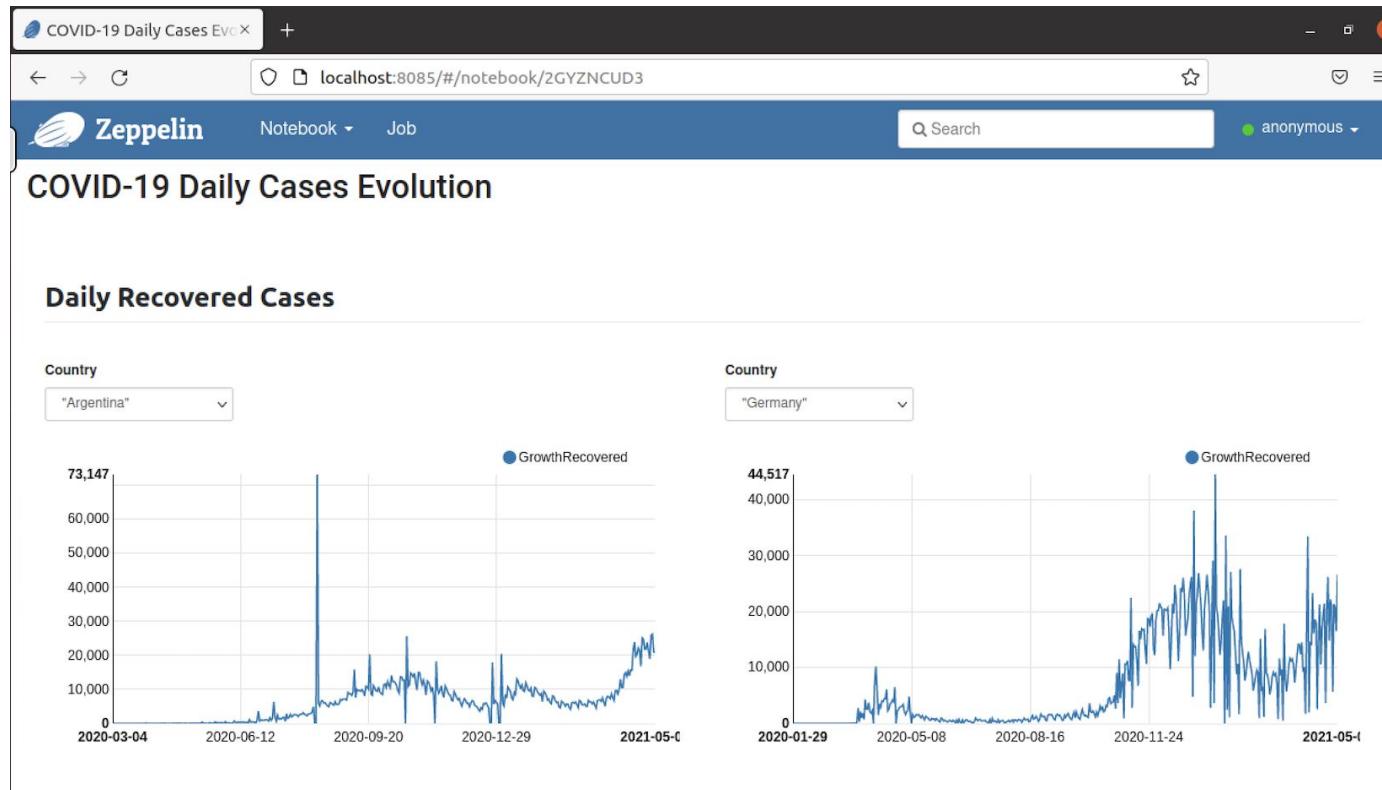
# Q5: What is the general disease evolution for all countries?

- The US had a spike of confirmed cases between Nov/2020 and Feb/2021, while Brazil had a constant trend.



# Q5: What is the general disease evolution for all countries?

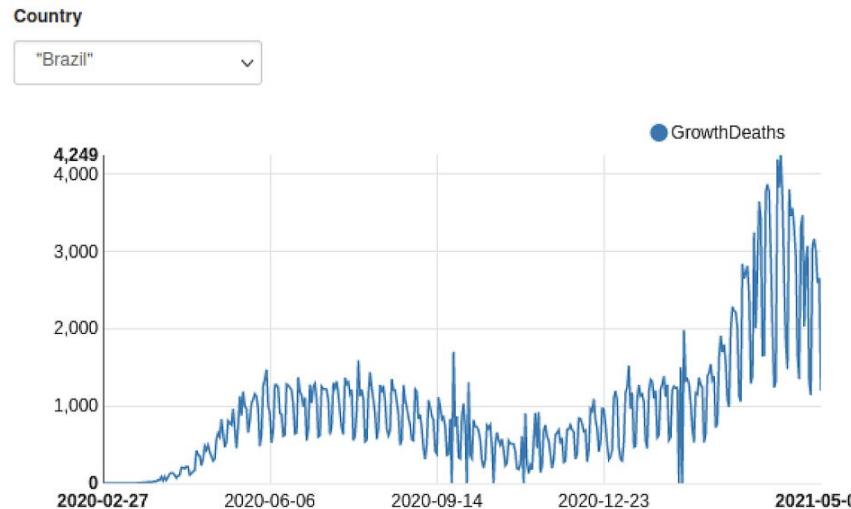
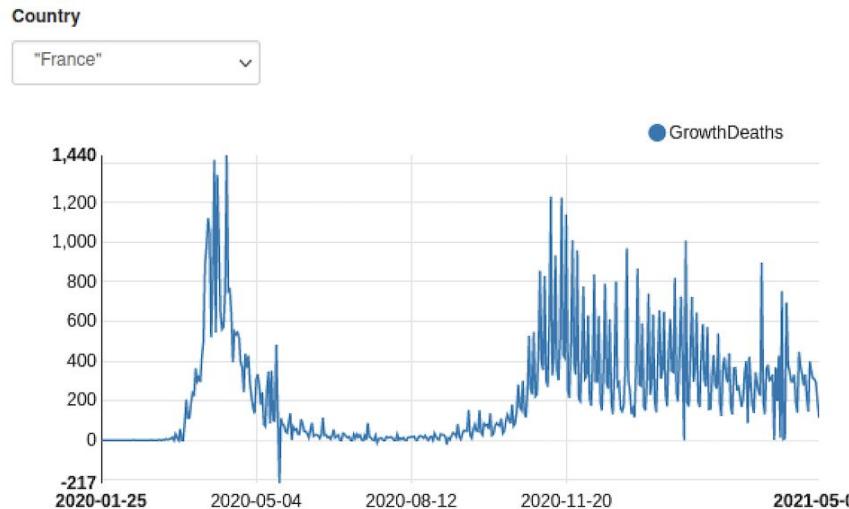
- Argentina and Germany had similar trends regarding recovered cases. Argentina had a huge spike in a day in 2020. That looks like an outlier for Aug-11-2020.



# Q5: What is the general disease evolution for all countries?

- France and Brazil had different trends regarding deaths. France had an initial spike between Jan/2020 and Apr/2020. After that, deaths stabilized and increased in Oct/2020. Brazil had a steady rate, which increased dramatically in Apr/2021.

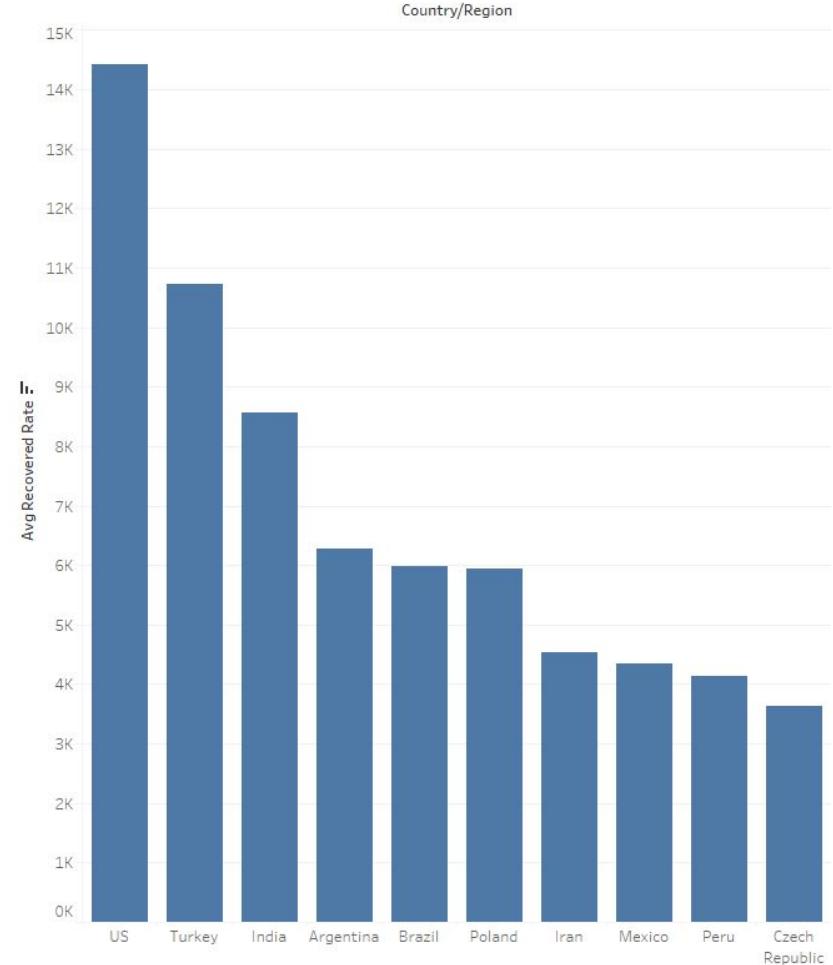
## Daily Covid Deaths



# Q6: What is the average recovered rate by countries?

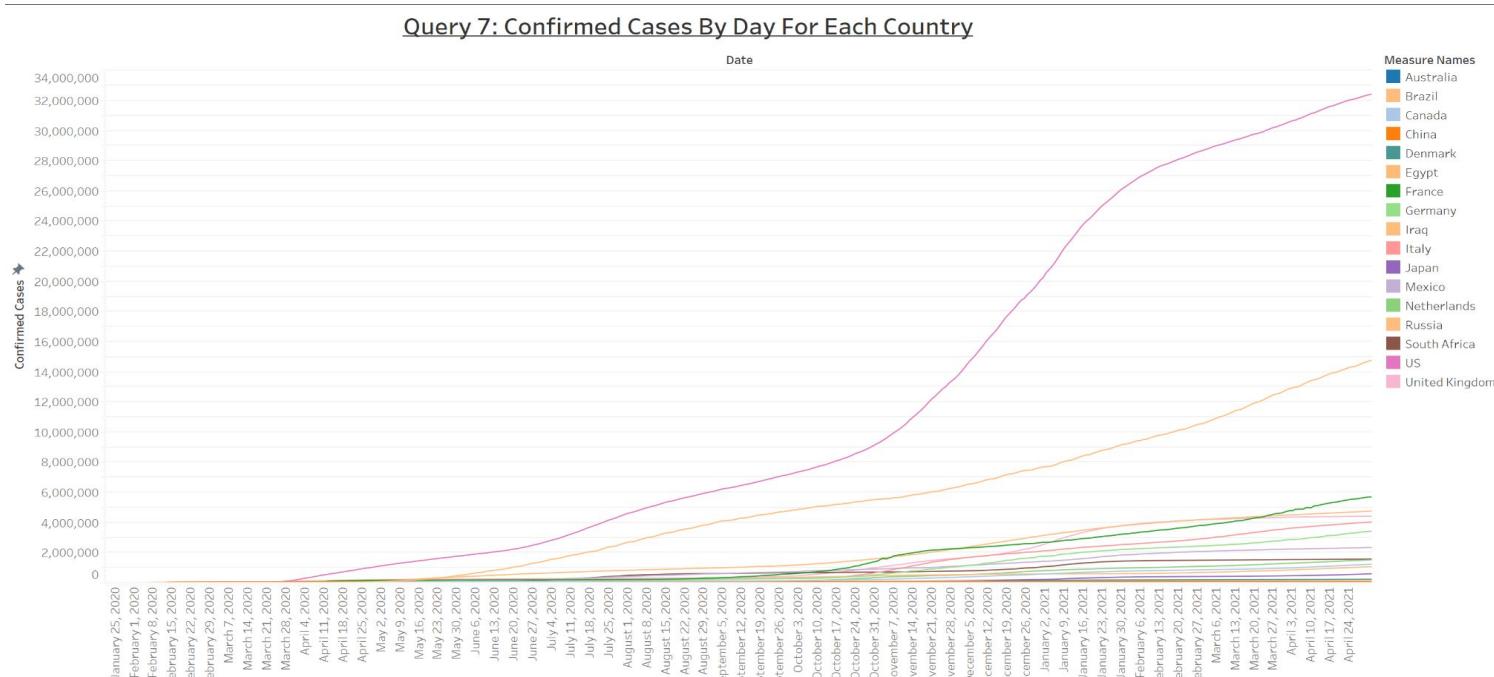
- Top 10 countries avg recovered rate per day
- Ordered by countries
- Pandemic recorded dates  
(01/22/2020-05/02/2021)

Sheet 1



# Q7: Total confirmed cases by countries

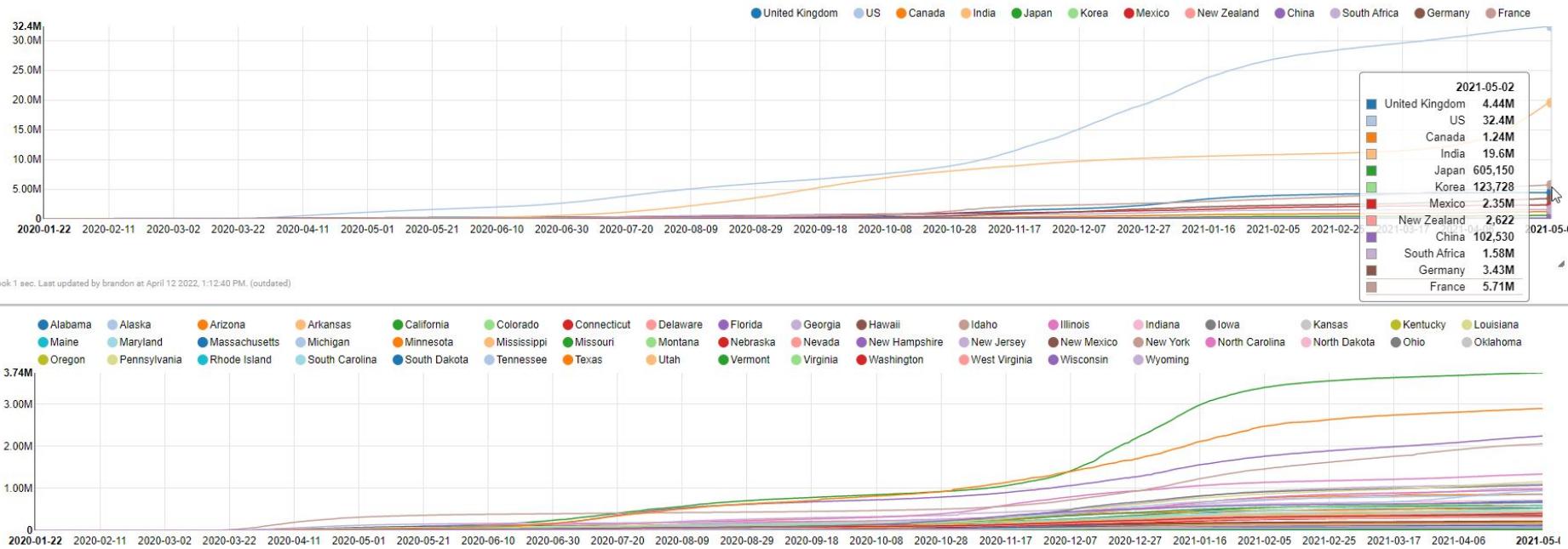
Tableau visualization:



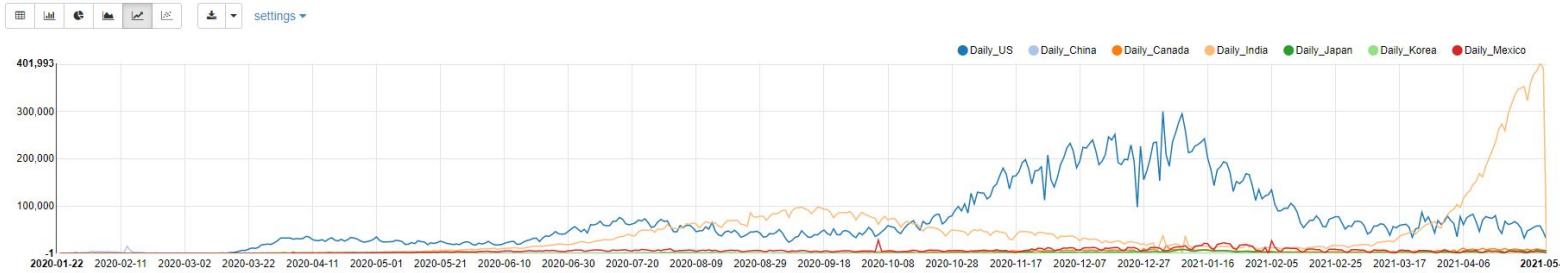
The Above Line Graph shows a visual from the results of running a query 7. Query 7 shows a comparison of the daily amount of confirmed casesfor the countries of Australia, Brazil, Canada, China, Denmark, Egypt, France, Germany, Iraq, Italy, Japan, Mexico, Netherlands, Russia, South Africa, United States, and United Kingdom for the whole duration of the pandemic.

## Q7: Total confirmed cases by country and states

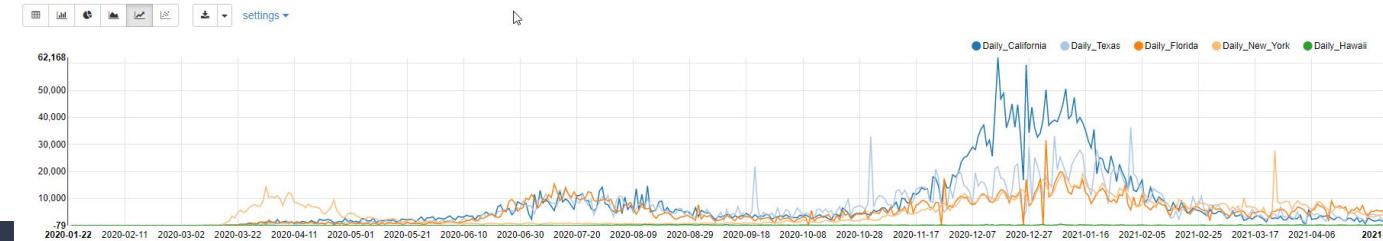
## Zeppelin visualization:



```
%spark.sql
SELECT
  Date,
  US - LAG(US)
    OVER(ORDER BY Date) AS Daily_US,
  China - LAG(China)
    OVER(ORDER BY Date) AS Daily_China,
  Canada - LAG(Canada)
    OVER(ORDER BY Date) AS Daily_Canada,
  India - LAG(India)
    OVER(ORDER BY Date) AS Daily_India,
  Japan - LAG(Japan)
    OVER(ORDER BY Date) AS Daily_Japan,
  Korea - LAG(Korea)
    OVER(ORDER BY Date) AS Daily_Korea,
  Mexico - LAG(Mexico)
    OVER(ORDER BY Date) AS Daily_Mexico
FROM tsConfirmedCountry
ORDER BY Date;
```



```
%spark.sql
SELECT
  Date,
  California - LAG(California)
    OVER(ORDER BY Date) AS Daily_California,
  Texas - LAG(Texas)
    OVER(ORDER BY Date) AS Daily_Texas,
  Florida - LAG(Florida)
    OVER(ORDER BY Date) AS Daily_Florida,
  NewYork - LAG(NewYork)
    OVER(ORDER BY Date) AS Daily_NewYork,
  Hawaii - LAG(Hawaii)
    OVER(ORDER BY Date) AS Daily_Hawaii
FROM tsConfirmedUS
ORDER BY Date;
```



## Q7 Part 2: Daily Confirmed cases by Country/State

# Q8: What is the total numbers of confirmed, death, and recovered cases?

- Clean/sanitized COVID data is cumulative
- Found the sum of the last recorded date
- Showing the total worldwide
- Deaths is about 1% of all cases combined

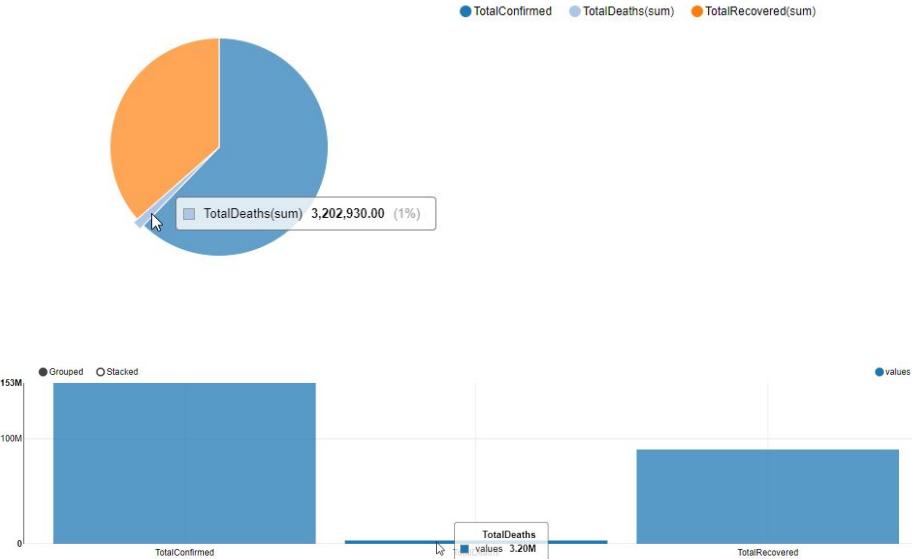
```
%spark
import org.apache.spark.sql.Row
import org.apache.spark.sql.DataFrame, SparkSession
import org.apache.spark.sql.types.{ StructType, StructField, StringType }

val c19Clean = "Data/Project2/clean/covid_19_data_clean.csv"
val q8Path = "Data/Project2/results/q8_TotalCDR/part-00000-71caa67e-bffd-463f-b135-5284caef8e7-c000.csv"
val df_c19Clean: DataFrame = spark.read.options(Map("header"->"true", "inferSchema"->"true", "delimiter"->","")).csv(c19Clean).na.fill(0).na.fill("")
val totalCDR: DataFrame = spark.read.options(Map("header"->"true", "inferSchema"->"true", "delimiter"->","")).csv(q8Path).na.fill(0).na.fill("")

df_c19Clean.registerTempTable("Covid19Data")
```

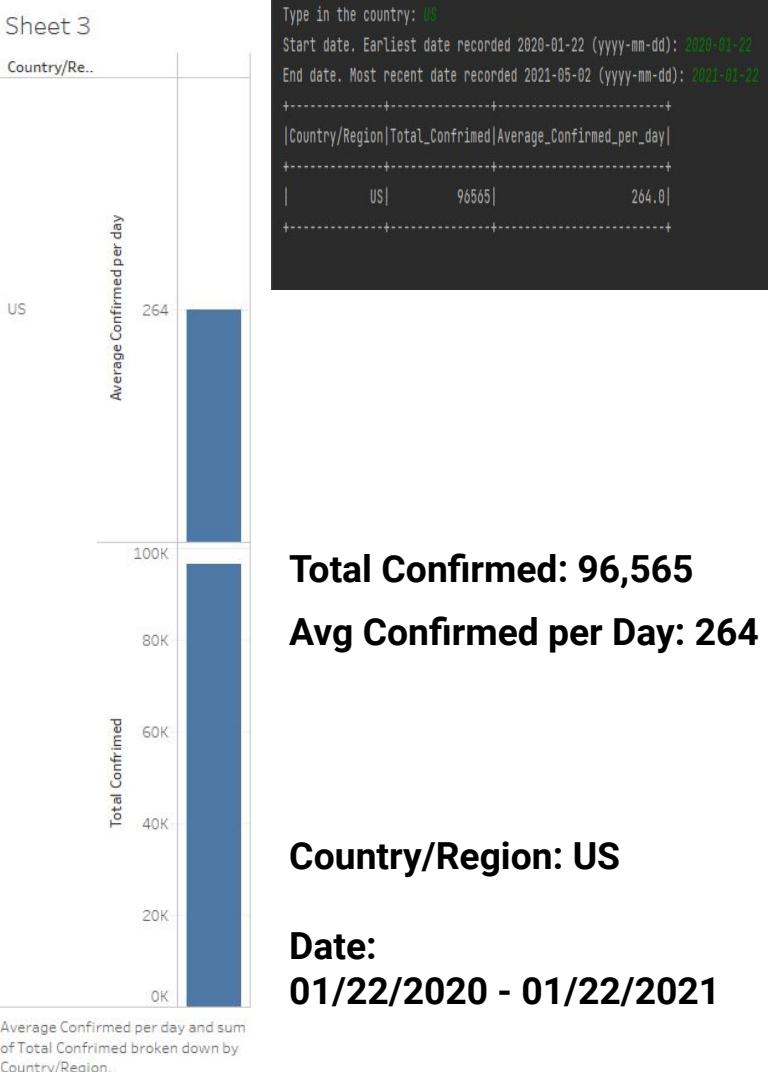
warning: there was one deprecation warning (since 2.0.0); for details, enable `:setting -deprecation` or `:replay -deprecation`  
import org.apache.spark.sql.Row  
import org.apache.spark.sql.DataFrame, SparkSession  
import org.apache.spark.sql.types.{StructType, StructField, StringType}  
c19Clean: String = Data/Project2/clean/covid\_19\_data\_clean.csv  
q8Path: String = Data/Project2/results/q8\_TotalCDR/part-00000-71caa67e-bffd-463f-b135-5284caef8e7-c000.csv  
df\_c19Clean: org.apache.spark.sql.DataFrame = [SNo: int, ObservationDate: string ... 6 more fields]  
totalCDR: org.apache.spark.sql.DataFrame = [TotalConfirmed: int, TotalDeaths: int ... 1 more field]

```
%spark.sql
SELECT sum(Confirmed) AS TotalConfirmed, sum(Deaths) AS TotalDeaths, sum(Recovered) AS TotalRecovered FROM Covid19Data WHERE ObservationDate = "2021-05-02";
```



# Q9: Spread speed for confirmed cases for a determined time and country.

- Implemented search feature
- Grabbed specific data from database
- Selected country and time period
- Show confirmed cases between time period



# Q10: Spread speed for death cases for a determined time and country.

- Similar to Q9
- Implemented search feature
- Selected country and time period
- Show death cases between time period

Confirmed spread speed for determined time and country  
 Type in the country: US  
 Start date. Earliest date recorded 2020-01-22 (yyyy-mm-dd): 2020-01-22  
 End date. Most recent date recorded 2021-05-02 (yyyy-mm-dd): 2021-01-22

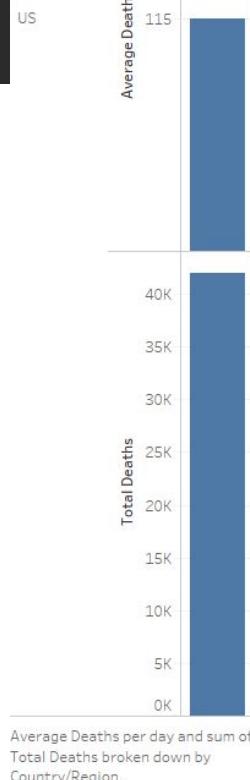
Country/Region	Total_Deaths	Average_Deaths_per_day
US	41974	115.0

**Total Deaths : 41,974**

**Avg Deaths Per Day: 115**

**Country/Region: US**

**Date:**  
**01/22/2020 - 01/22/2021**



# Questions?

## GitHub

<https://github.com/jaceguaidemagalhaes/Big-Data---Project-2>