

...
Soren DeOrlow
IDSN 599, Spring 2021
Deorlow@usc.edu
Final Project Part 1
...

Domain: What is the type/nature of the problem that you are thinking to solve

Through the process of searching for data sources, there has been an abundance of data surrounding covid-19 along with interesting projects attempting to draw new insights from this vast data. I seek to become aware of the work that is happening in this space and through prototyping, gain new insight on data patterns around covid-19.

Dataset: Where did you get your dataset? How big is it – how many rows and attributes? Why/how does it contain the data you need to solve the problem you are thinking of solving?

I have found numerous datasets on Kaggle with fascinating data on covid 19, including lung x-ray imaging and covid vaccine related tweets. The cleanest dataset that I found maps vaccine data globally in relation to total population as well as covid fatalities. It contains 21k rows of data on global patterns. I also discovered that Roche has launched a project called the uncover Covid-19 challenge focused on data exploration and research. This project has launched 18 sub tasks inviting data research contributions to provide new perspective and answers to questions such as, "Predicting illness severity in a particular patient or demographic." I have chosen two similar datasets that I am working with to gain greater insight into the patterns of covid-19 from a public health standpoint and how various populations around the world have responded to vaccines and other measures.

Problem Type: Are you creating a predictor, a classifier, or something else? Why do you think this is the way to go to solve the problem?

I plan to use classification to gain insight into how different parts of the world are approaching vaccination and the velocity of vaccination over time. I also am looking at the relationships of GDP, vaccination and covid morbidity. I'm less concerned with the relationship between covid morbidity, but I will not rule out any insights that might be found. Ultimately, I would like to see changes over time.

Attributes: What are the attributes of the dataset? Which are numeric and which are text? Do you have any missing values for attributes?

As I was exploring various datasets, I was careful to select datasets that were robust and consisted primarily of numeric data, with limited missing attributes. The datasets that I have chosen, list countries and how they are doing in the fight against covid-19. Below is a sample of both dataset including a complete view of the column headers.

Dataset 1

	country	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	New_deaths	population	ratio
0	Afghanistan	AFG	5/11/21	504502	448878	55624	12	40146987	1.118086396
1	Afghanistan	AFG	5/20/21	547901	470341	77560	10	40146987	1.171547444
2	Afghanistan	AFG	5/24/21	573277	476367	96910	10	40146987	1.186557288
3	Afghanistan	AFG	5/26/21	590454	479372	111082	19	40146987	1.194042283
4	Afghanistan	AFG	5/27/21	593313	479574	113739	14	40146987	1.194545434
5	Afghanistan	AFG	5/30/21	600152	480226	119926	20	40146987	1.196169466

