

```
.file/github ==> assign.M.1.assignment.1.covid.data
.institution ==> University of Tennessee
.course ==> COSC.526 Intro. to Data Mining
```

## assign.M1.Assignment.1.covid.data

### 0.Problem summary

For this problem, you will be working with COVID-19 sequence processing data from Kaggle. The dataset contains data about the processing of COVID-19 sequences by different countries over time. It comes as a Comma-Separated Value (CSV) file. It contains the following 6 columns:

- **location:** the country for which the information is provided
- **date:** the date of the data entry
- **variant:** the COVID-19 variant for the data entry
- **num\_sequences:** the number of sequences processed (for the country, variant, and date)
- **num\_sequences\_total:** the total number of sequences available (for the country, variant, and date)
- **perc\_sequences:** the percentage of the available number of sequences that were processed (Note: this value is out of 100)
- **note:** each row (or data entry) in the dataset represents the processing of one variant by one country on one day.

### Objectives:

- 1.Import and manipulate a .csv file
- 2.Assess your Python Programming Skills

=> Other assignments are more challenging. Use this to assess your skills.  
=> Attempt to solve the problems without searching for online assistance.  
=> Prepare questions for class discussion to help source additional tools.

### Codebook and data files

File Name	Purpose\Description
<a href="https://github.com/cosc-526/cosc.526.home.page/blob/main/code_notebook_cosc_526.ipynb">https://github.com/cosc-526/cosc.526.home.page/blob/main/code_notebook_cosc_526.ipynb</a> <save your own copy!>	Course Codebook in Jupyter Notebook name = code.notebook.cosc.526.ipynb
<a href="#">d.M.1.10.assignment.covid.data.variant.s.csv</a>	Course github of source data
<a href="https://www.kaggle.com/yamgwe/omicron-covid19-variant-daily-cases?select=covid-variants.csv">https://www.kaggle.com/yamgwe/omicron-covid19-variant-daily-cases?select=covid-variants.csv</a>	i) Kaggle data homepage ii) grab an api key from this page if using that method to import data

**note.1:** the codebook is formatted differently and below highlights expected outcomes.

**note.2:** the instructions below are an overview and more details are in the notebook.

**note.3:** perform your work in the notebook, and export and submit as a .pdf file.

```
.file/github ==> assign.M.1.assignment.1.covid.data
.institution ==> University of Tennessee
.course ==> COSC.526 Intro. to Data Mining
```

## 0.Task.0 - Import, inspect, and view descriptive statistics

Import data and view descriptive statistics with the pandas library.

- grab data from **Github** URL, .csv. or kaggle api

### Task.0 - Expected outcome:

```
covid19-variants.zip: Skipping, found more recently modified local c
-----
> dataframe fields w pd.head <
-----
  location    date    variant  num_sequences  perc_sequences \
0  Angola  2020-07-06    Alpha             0             0.0
1  Angola  2020-07-06  B.1.1.277             0             0.0
2  Angola  2020-07-06  B.1.1.302             0             0.0
3  Angola  2020-07-06  B.1.1.519             0             0.0
4  Angola  2020-07-06   B.1.160             0             0.0

  num_sequences_total
0                   3
1                   3
2                   3
3                   3
4                   3
-----
==> descriptive statistics <==
-----
      num_sequences  perc_sequences  num_sequences_total
count          100416.0           100416.0           100416.0
mean              72.0             6.0             1510.0
std             1669.0            22.0             8445.0
min               0.0            -0.0              1.0
25%               0.0             0.0             12.0
50%               0.0             0.0             59.0
75%               0.0             0.0             394.0
max            142280.0            100.0            146170.0 1
```

## 1.Task.1 - find uncommon variants

A. The 3 main variants of COVID-19 that we've experienced in the US are:

- \* `Alpha`
- \* `Delta`
- \* `Omicron`

B. Assignment Tasks

1. What other variants are recognized by the World Health Organization (WHO) in this dataset?
2. Sort the variant names alphanumerically and store in a list.
3. Exclude in output `on\_who` and `others` from `variant`.

## 2.Task.2 -

Which variant of COVID-19 has the most sequences processed?

### 3.Task.3 -

Which country processed sequences the best of all variants including the "catch-all" categories?

### 4.Task.4 -

#### Part A

Which country did the best at processing sequences across the Alpha, Delta, and Omicron variants only?

- o hint: output is one country

#### Part B

Determine the ranking of the United States at processing sequences across the Alpha, Delta, and Omicron variants only?

- o hint: output is one country

### 5.Task.5 -

Determine each country's total number of processed sequences for the Omicron variant on December 27, 2021. Sort the output from highest number of processed sequences to the smallest number of processed sequences. Each element in the output should include both the name of the country and the number of processed sequences.

### 6.Task.7 -

Determine the percentage of processed sequences for the Alpha, Delta, and Omicron variants in the United States.

## 7. Additional resources

- <https://github.com/cosc-526/cosc.526.home.page>
- [Jupyter Community Forum](#)