

Amol Torar

Follow

Dec 5, 2021

5 min read

Listen

# 10 Python operations for 70% of your Data Analysis

Image designed by Unsplash

Day to day work of a Data Scientist involves a variety of tasks such as data pre-processing, data analysis, machine learning model creation, model deployment.

If you are starting on your journey of becoming a Data Scientist, then the first step is to acquire the skill of data manipulation because it is required in almost every Data Science project. Data Manipulation starts with reading your data and converting it into a form where you can answer your questions from data. Python programming language has Pandas library written for data manipulation and data analysis work.

In this blog, I will give you an overview to top 10 python (pandas) operations that every aspiring Data Scientist should know:

1. **Reading dataset**: Data is the ingredient of every analysis. Knowing how to read data from different file formats such as : csv, excel, text etc is one of the 1st steps that you should be adept at as a Data Scientist. Below is an example of how you can read a csv file, containing Covid-19 data, using pandas.

```
import pandas as pd
# reading the countries_data file along with the location within read_csv function.
countries_df = pd.read_csv('C:/Users/annol/Desktop/Courses/Python for Data Science/Code/countries_data')
# showing the first 5 rows of the dataframe
countries_df.head()
```

Following is the output of countries\_df.head() using which we can see the first 5 rows of a dataframe :

|   | Country     | CountryCode | Bug         | NewConfirmed | TotalConfirmed | NewDeaths | TotalDeaths | NewRecovered | TotalRecovered | Date                 | Premium |
|---|-------------|-------------|-------------|--------------|----------------|-----------|-------------|--------------|----------------|----------------------|---------|
| 0 | Afghanistan | AF          | afghanistan | 75           | 38716          | 0         | 1420        | 404          | 31638          | 2020-09-14T18:31:55Z | 0       |
| 1 | Albania     | AL          | albania     | 168          | 11353          | 4         | 334         | 75           | 6569           | 2020-09-14T18:31:55Z | 0       |
| 2 | Algeria     | DZ          | algeria     | 247          | 48264          | 7         | 1612        | 162          | 34037          | 2020-09-14T18:31:55Z | 0       |
| 3 | Andorra     | AD          | andorra     | 0            | 1344           | 0         | 53          | 0            | 943            | 2020-09-14T18:31:55Z | 0       |
| 4 | Angola      | AO          | angola      | 53           | 3388           | 2         | 134         | 12           | 1301           | 2020-09-14T18:31:55Z | 0       |

2. **Summary Statistics** : Once you have read the dataset, the next step is to understand the data by looking at the data summaries such as count, mean, standard deviation(std), 25th percentile etc. of numeric columns such as NewConfirmed, TotalConfirmed etc and frequency, top occurring value etc of categorical columns such as Country, CountryCode etc.

Using the describe function we can get the summary of continuous variables of the dataset as shown below :

```
1 #get summary of continuous variables
2 countries_df.describe()
```

|       | NewConfirmed | TotalConfirmed | NewDeaths   | TotalDeaths   | NewRecovered | TotalRecovered |
|-------|--------------|----------------|-------------|---------------|--------------|----------------|
| count | 1.880000e+02 | 1.880000e+02   | 188.000000  | 188.000000    | 188.000000   | 1.880000e+02   |
| mean  | 3.527606e+04 | 1.882256e+05   | 19.372340   | 4913.930851   | 988.175532   | 1.043737e+05   |
| std   | 4.664484e+05 | 8.155154e+05   | 95.039647   | 19412.818049  | 6026.992378  | 4.348660e+05   |
| min   | 0.000000e+00 | 0.000000e+00   | 0.000000    | 0.000000      | 0.000000     | 0.000000e+00   |
| 25%   | 1.000000e+00 | 2.245750e+03   | 0.000000    | 37.750000     | 0.000000     | 1.205000e+03   |
| 50%   | 5.300000e+01 | 1.021750e+04   | 0.000000    | 208.000000    | 15.500000    | 6.556000e+03   |
| 75%   | 4.780000e+02 | 6.937975e+04   | 5.250000    | 1120.000000   | 263.750000   | 4.391675e+04   |
| max   | 6.396100e+06 | 6.519673e+06   | 1136.000000 | 194071.000000 | 77512.000000 | 3.780107e+06   |

Within the describe function, we can set the argument “include = ‘all’” to get summaries of both continuous and categorical variables.

```
1 #get summary of continuous and categorical variables
2 countries_df.describe(include = 'all')
```

|        | Country    | CountryCode | Bug           | NewConfirmed | TotalConfirmed | NewDeaths   | TotalDeaths   | NewRecovered | TotalRecovered | Date                 | Premium |
|--------|------------|-------------|---------------|--------------|----------------|-------------|---------------|--------------|----------------|----------------------|---------|
| count  | 188        | 187         | 188           | 1.880000e+02 | 1.880000e+02   | 188.000000  | 188.000000    | 188.000000   | 1.880000e+02   | 188                  | 188     |
| unique | 188        | 187         | 188           | NaN          | NaN            | NaN         | NaN           | NaN          | NaN            | NaN                  | 1       |
| top    | Montenegro | ES          | serbia-bosnia | NaN          | NaN            | NaN         | NaN           | NaN          | NaN            | 2020-09-14T18:31:55Z | 0       |
| freq   | 1          | 1           | 1             | NaN          | NaN            | NaN         | NaN           | NaN          | NaN            | 188                  | 188     |
| mean   | NaN        | NaN         | NaN           | 3.527606e+04 | 1.882256e+05   | 19.372340   | 4913.930851   | 988.175532   | 1.043737e+05   | NaN                  | NaN     |
| std    | NaN        | NaN         | NaN           | 4.664484e+05 | 8.155154e+05   | 95.039647   | 19412.818049  | 6026.992378  | 4.348660e+05   | NaN                  | NaN     |
| min    | NaN        | NaN         | NaN           | 0.000000e+00 | 0.000000e+00   | 0.000000    | 0.000000      | 0.000000     | 0.000000e+00   | NaN                  | NaN     |
| 25%    | NaN        | NaN         | NaN           | 1.000000e+00 | 2.245750e+03   | 0.000000    | 37.750000     | 0.000000     | 1.205000e+03   | NaN                  | NaN     |
| 50%    | NaN        | NaN         | NaN           | 5.300000e+01 | 1.021750e+04   | 0.000000    | 208.000000    | 15.500000    | 6.556000e+03   | NaN                  | NaN     |
| 75%    | NaN        | NaN         | NaN           | 4.780000e+02 | 6.937975e+04   | 5.250000    | 1120.000000   | 263.750000   | 4.391675e+04   | NaN                  | NaN     |
| max    | NaN        | NaN         | NaN           | 6.396100e+06 | 6.519673e+06   | 1136.000000 | 194071.000000 | 77512.000000 | 3.780107e+06   | NaN                  | NaN     |

*Want to get an in-depth understanding of python for Data Analysis ? You can follow the official python documentation or you can enroll into my course on Udemy and become certified in python for Data Analysis. Use the following link to avail 70% discount : <https://bit.ly/3148Qq6>*

3. **Data selection and filtering** : Not all the rows and columns of a dataset are required for the analysis. You would need to select the columns of interest and filter some rows based on the question that you are trying to answer.

For example, we can select Country and NewConfirmed columns using the following code :

```
1 # selecting Country and NewConfirmed columns
2 countries_df[['Country','NewConfirmed']]
```

| Country        | NewConfirmed |
|----------------|--------------|
| 0 Afghanistan  | 75           |
| 1 Albania      | 168          |
| 2 Algeria      | 247          |
| 3 Andorra      | 0            |
| 4 Angola       | 53           |
| ...            | ...          |
| Viet Nam       | 3            |
| Western Sahara | 0            |
| Yemen          | 2            |
| Zambia         | 73           |
| Zimbabwe       | 18           |

We can also filter the data for United States of America as country. Using loc, we can filter a column based on some value as shown below :

```
1 # filtering USA using country column
2 countries_df.loc[countries_df['Country'] == 'United States of America']
```

|     | Country                  | CountryCode | Bug           | NewConfirmed | TotalConfirmed | NewDeaths | TotalDeaths | NewRecovered | TotalRecovered | Date                 | Premium |
|-----|--------------------------|-------------|---------------|--------------|----------------|-----------|-------------|--------------|----------------|----------------------|---------|
| 179 | United States of America | US          | united states | 2450         | 62795e+3       | 372       | 199071      | 19708        | 241180         | 2020-09-14T18:31:55Z | 0       |

4. **Aggregation** : Finding numeric summaries such as count, sum, mean etc at different variable groupings is data aggregation. It is one of the most frequently performed task of a Data Scientist.

We can find the total of NewConfirmed cases across the countries using aggregation. Aggregation is performed using groupby and agg functions. Within groupby function, we provide the level at which we want to perform the aggregation (Country column) and within the aggregation function we provide the column name (NewConfirmed) and mathematical operation (sum) we want to perform on the column.

```
1 # total NewConfirmed cases across countries
2 countries_df.groupby(['Country']).agg({'NewConfirmed':'sum'})
```

| NewConfirmed   |     |
|----------------|-----|
| Country        |     |
| Afghanistan    | 75  |
| Albania        | 168 |
| Algeria        | 247 |
| Andorra        | 0   |
| Angola         | 53  |
| ...            | ... |
| Viet Nam       | 3   |
| Western Sahara | 0   |
| Yemen          | 2   |
| Zambia         | 73  |
| Zimbabwe       | 18  |

5. **Join** : Combining 2 datasets to create one single dataset is done using Join operation. Many times, different information is present in different datasets, for example, one dataset could contain the count of Covid-19 cases across different countries and another dataset could contain the latitude and longitude information of the different countries. Now, if we need to combine these 2 informations then we can perform a join operation as shown below :

```
1 #reading countries lat and lon data
2 countries_lat_lon = pd.read_excel('C:/Users/annol/Desktop/Courses/Python for Data Science/Code/countries_lat_lon')
3
4 # joining the 2 dataframe : countries_df and countries_lat_lon
5 # syntax : pd.merge(left_df, right_df, on = 'on_column', how = 'type_of_join')
6 joined_df = pd.merge(countries_df, countries_lat_lon, on = 'CountryCode', how = 'inner')
7 joined_df
```

| CountryCode | Latitude   | Longitude | Country     | Bug         | NewConfirmed | TotalConfirmed | NewDeaths | TotalDeaths | NewR |
|-------------|------------|-----------|-------------|-------------|--------------|----------------|-----------|-------------|------|
| 0 AF        | 33.939110  | 67.709953 | Afghanistan | afghanistan | 75           | 38716          | 0         | 1420        |      |
| 1 AL        | 41.153332  | 20.168331 | Albania     | albania     | 168          | 11353          | 4         | 334         |      |
| 2 AO        | -11.202692 | 17.873887 | Angola      | angola      | 53           | 3388           | 2         | 134         |      |

6. **Built-in functions** : Knowing the mathematical built-in functions such as min(), max(), mean(), sum() etc is very helpful for performing different analysis. We can apply these functions directly on a dataframe simply by calling them. These functions can be used standalone on a column or within the aggregation function as shown below :

```
1 # finding sum of NewConfirmed cases of all the countries
2 countries_df['NewConfirmed'].sum()
3 # Output : 6,631,899
4
5 # finding the sum of NewConfirmed cases across different countries
6 countries_df.groupby(['Country']).agg({'NewConfirmed':'sum'})
7
8 # Output
9 # Country
10 #Afghanistan 75
11 #Albania 168
12 #Algeria 247
13 #Andorra 0
14 #Angola 53
```

7. **User defined functions** : Functions that we write on our own are user-defined functions. We can execute the codes within these functions, whenever needed, by calling that function. For example, we can create a function to add 2 numbers as shown below :

```
1 # User defined function is created using 'def' keyword, followed by function definition - 'addition'
2 # and 2 arguments num1 and num2
3 def addition(num1, num2):
4     return num1+num2
5
6 # calling the function using function name and providing the arguments
7 print(addition(1,2))
8 #output : 3
```

8. **Pivot** : Pivoting is converting the unique values within the rows of a column into multiple new columns. This is advance data manipulation technique. Using pivot\_table() function on the Covid-19 dataset, we can convert the country names into individual new columns :

```
1 # using pivot_table to convert values within the Country column into individual columns and
2 # filling the values corresponding to these columns with numeric variable - NewConfirmed
3 pivot_df = pd.pivot_table(countries_df, columns = 'Country', values = 'NewConfirmed')
4 pivot_df
```

|              | Country | Afghanistan | Albania | Algeria | Andorra | Angola | Antigua and Barbuda | Argentina | Armenia | Australia | Austria | ... | United Kingdom |
|--------------|---------|-------------|---------|---------|---------|--------|---------------------|-----------|---------|-----------|---------|-----|----------------|
| NewConfirmed |         | 75          | 168     | 247     | 0       | 53     | 0                   | 9056      | 187     | 41        | 463     | ... | 3338           |

9. **Iterating over dataframe** : Many times it is required to iterate through the index and rows of a data frame. We can iterate through the dataframe using iterrows function :

```
1 # iterating over the index and row of a dataframe using iterrows() function
2 for index, row in countries_df.iterrows():
3     print('Index is ' + str(index))
4     print('Country is ' + str(row['Country']))
5
6 # Output :
7 # Index is 0
8 # Country is Afghanistan
9 # Index is 1
10 # Country is Albania
11 # .....
```

10. **String operations** : Many times we deal with string columns in our dataset, in such cases it is important to know some basic string operations such as how to convert a string into upper case, lower case and how to find the length of a string in a column.

```
1 # country column to upper case
2 countries_df['Country_upper'] = countries_df['Country'].str.upper()
3
4 # country column to lower case
5 countries_df['CountryCode_lower']=countries_df['CountryCode'].str.lower()
6
7 # finding length of characters in the country column
8 countries_df['len'] = countries_df['Country'].str.len()
9
10 countries_df.head()
```

| NewDeaths | TotalDeaths | NewRecovered | TotalRecovered | Date                 | Premium | Country_upper | CountryCode_lower | len |
|-----------|-------------|--------------|----------------|----------------------|---------|---------------|-------------------|-----|
| 0         | 1420        | 404          | 31638          | 2020-09-14T18:31:55Z | 0       | AFGHANISTAN   | af                | 11  |
| 4         | 334         | 75           | 6569           | 2020-09-14T18:31:55Z | 0       | ALBANIA       | al                | 7   |
| 7         | 1612        | 162          | 34037          | 2020-09-14T18:31:55Z | 0       | ALGERIA       | dz                | 7   |
| 0         | 53          | 0            | 943            | 2020-09-14T18:31:55Z | 0       | ANDORRA       | ad                | 7   |
| 2         | 134         | 12           | 1301           | 2020-09-14T18:31:55Z | 0       | ANGOLA        | ao                | 6   |

Knowing how to perform these 10 operations will cater to almost 70% of your data manipulation needs.

818

2

Sign up for CrunchX

By CodeX

A weekly newsletter on what's going on around the tech and programming space [Take a look.](#)

Get this newsletter

Emails will be sent to jeabinay@gmail.com.

Not you?

More from CodeX

Everything connected with Tech & Code. Follow to join our 900K+ monthly readers

Haekal · Dec 5, 2021

## The Fall of the Greatest Modern Inventor

The Failure of Steve Jobs The story about Steve Jobs always fascinating even the story about his failure and this is the story when Steve Jobs facing his lowest point in life. After a decade working on the company that he...

Mac Book · 3 min read

Share your ideas with millions of readers.

Write on Medium

Jeffrey Cloe · Dec 5, 2021

## The 2021 Macbook Pro Keeps Surprising Me!

An Upgrade that Continues to Impress — I must admit I had been eying the newer Macbook Pro for the last couple of years. I remember the complaints regarding the butterfly keyboard and Apple eventually listeni...

Apple · 3 min read

Zllin · Dec 5, 2021

## Telegram HTTP API With Python — Sending Messages Programmatically

Perhaps you're training a machine learning model and wish to notify yourself through Telegram when it's done training. Perhaps you have so...

Python · 4 min read

Norberto Santiago · Dec 5, 2021

## Seek and Destroy — The JavaScript Algorithm

During our algorithm practice, some exercises are boring, others will be fun, and others you can't help but blog about them. That's the case with this fun algorithm I came across a few days ago. Seek and Destroy, the...

JavaScript · 3 min read

Andrew Zuo · Dec 5, 2021

## Why Smart People Are Often Wrong On The Internet

By smart here I don't mean you think they're smart. I mean smart as in the ability to build up a logical argument. Even if a logical argument is flawed. The idea for this post really came out of conspiracy theorists on...

Cryptocurrency · 3 min read

Read more from CodeX