
Welcome to Data Science Online Bootcamp

Week 1

**Day 2-3 : Pandas Recap & Working
with CSV Files**



Democratizing Data Science Learning

Learning Objectives

1. Pandas Recap

2. Working with CSV files



Pandas Session Recap

What we learnt so far?

- Pandas Objects/Data Structures (refer pre-session slides for clear explanation)
 - Series
 - Dataframe
- Data indexing and selection
 - **Use case:** Helps fetch a data record when we are dealing with large volumes of data. For ex: If you have a data with 10 million records, you can easily fetch information of a particular serial no/index no with pandas. This operation may not be feasible with MS Excel while dealing with such large volumes



Pandas Session Recap

- Data Wrangling & Handling Missing Values - read the notebook.
 - **Use case:** Often Data Scientists get unclean data with a lot missing values and we need to have a solution to deal with it. Pandas facilitate in handling this issue
- Pandas String Operations
 - **Use case:** Helps in Handling Missing Values



Next steps: Pandas Session Recap

How would you **eat an elephant?**

... one bite at a time!

Disclaimer: The writer is a vegetarian by no means the visual is meant to hurt animals, it is just given as an analogy.

Next steps: Pandas Session Recap

- Similar to the analogy given in previous slide. Let's attack one concept at a time in Pandas.
- **Day 2:** We can learn about Pandas Objects and Index. Alongside, we can try solving some exercises given in the notebook. You can target to solve Exercises 1-3
- **Day 3:** We can learn about handling missing data and solving exercises from 4-6.
- This can be stretched till Day 4 too. Overall this could be a little hectic but with our strong desire to learn data science, we are confident we can overcome the hurdles that come midway! Happy learning :)

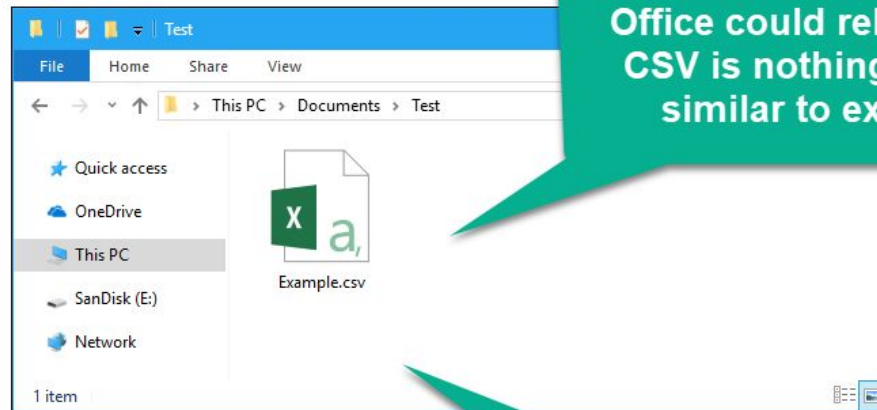


Working with CSVs



Democratizing Data Science Learning

What is a CSV file?



Windows users using MS Office could relate to this. CSV is nothing but a file similar to excel files

Only catch is, in CSV you store data with comma separated values

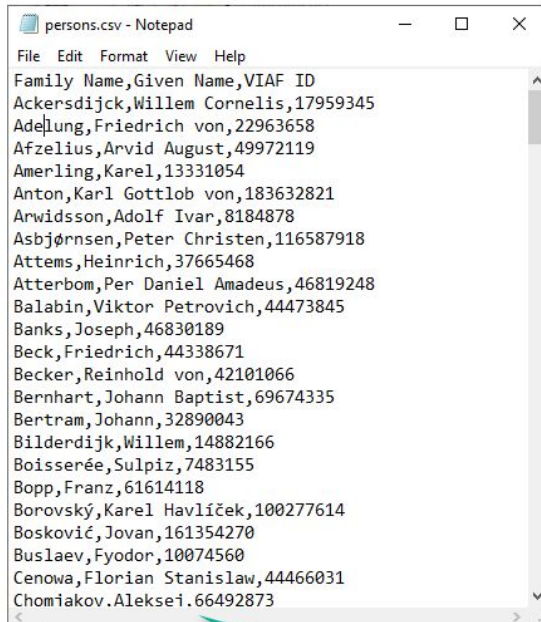


What is a CSV ?

- CSV files are normally created by programs that handle large amounts of data. They are a convenient way to export data from spreadsheets and databases as well as import or use it in other programs.
- CSV (Comma Separated Values) is a simple file format used to store tabular data, such as a spreadsheet or database.
- A CSV file stores tabular data (numbers and text) in plain text.
- Each line of the file is a data record.
- Each record consists of one or more fields, separated by commas.
- The use of the comma as a field separator is the source of the name for this file format.

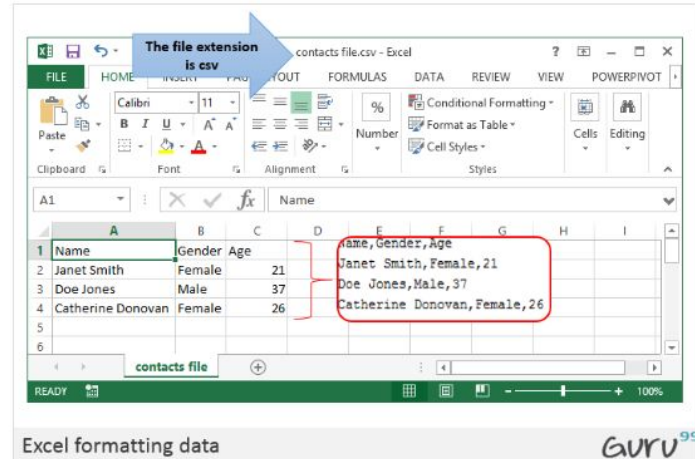


How does it look like?



```
File Edit Format View Help
Family Name,Given Name,VIAF ID
Ackersdijk,Willem Cornelis,17959345
AdeJung,Friedrich von,22963658
Afzelius,Arvid August,49972119
Amerling,Karel,13331054
Anton,Karl Gottlob von,183632821
Arwidsson,Adolf Ivar,8184878
Asbjørnsen,Peter Christen,116587918
Attems,Heinrich,37665468
Atterbom,Per Daniel Amadeus,46819248
Balabin,Viktor Petrovich,44473845
Banks,Joseph,46830189
Beck,Friedrich,44338671
Becker,Reinhold von,42101066
Bernhart,Johann Baptist,69674335
Bertram,Johann,32890043
Bilderdijk,Willem,14882166
Boisserée,Sulpiz,7483155
Bopp,Franz,61614118
Borovský,Karel Havlíček,100277614
Bosković,Jovan,161354270
Buslaev,Fyodor,10074560
Cenowa,Florian Stanislaw,44466031
Chomiakov,Aleksei,66492873
```

Looks like this
when opened on
Notepad



The file extension is csv

	A	B	C	D	E	F	G	H	I
1	Name	Gender	Age		Name, Gender, Age				
2	Janet Smith	Female	21		Janet Smith, Female, 21				
3	Doe Jones	Male	37		Doe Jones, Male, 37				
4	Catherine Donovan	Female	26		Catherine Donovan, Female, 26				
5									
6									

Excel formatting data

Guru⁹⁹

Looks like this
when opened on
MS Excel



Working with CSV files in Python

- For working CSV files in python, there is an inbuilt module named csv.
- However, a common method for working with CSV files is using Pandas. It makes importing and analyzing data much easier.
- One crucial feature of Pandas is its ability to write and read Excel, CSV, and many other types of files.



Pandas read_csv

- Functions like the Pandas read_csv() method enable you to work with files effectively.
- The read_csv() function reads the CSV file into a DataFrame object.
- A CSV file is similar to a two-dimensional table and the DataFrame object represents two dimensional tabular view.
- The most basic way to read a csv file in Pandas:

```
# Import pandas  
import pandas as pd
```

```
# reading csv file  
pd.read_csv("filename.csv")
```

- Read the next slide to understand how to provide filename



Pandas read_csv

```
In [4]: # Import pandas
import pandas as pd
```

```
In [5]: pwd #Let's get to know which directory/folder we are working with
```

```
Out[5]: 'C:\\Users\\chanukya\\Documents'
```

```
In [9]: # reading csv file
pd.read_csv("test1.csv") # we have a file with name test1.csv in documents folder so it is easy to directly call the name
```

```
Out[9]:
```

	username	password	firstname	lastname	email	cohort1
0		NaN	Chanukya Patnaik	Learner		Beginner

```
In [10]: #Let's say the file is in a different directory Eg: on my desktop
```

```
In [11]: pd.read_csv("C:/Users/chanukya/Desktop/test1.csv") #we can directly put the folder location as mentioned here within parentheses
```

```
Out[11]:
```

	username	password	firstname	lastname	email	cohort1
0		NaN	Chanukya Patnaik	Learner		Beginner



Pandas read_csv

- There are many other things one can do through this function only to change the returned object completely.
- For instance, one can read a csv file not only locally, but from a URL through read_csv or one can choose what columns needed to export so that we don't have to edit the array later.
- These modifications can be done by the various arguments it takes.
- We don't need to memorise all the arguments though, let's have a look at few important ones in the next two slides.



Pandas to_csv

- The easiest way to write DataFrames to CSV files is using the Pandas to_csv function
- Syntax:

```
# DataFrame to CSV file  
# df is the name of the DataFrame here  
df.to_csv('file_name.csv')
```

Where df is the name of your DataFrame

- If you want to export without the index, simply add index=False

```
# Specify index as False to import without index  
df.to_csv('file_name.csv', index=False)
```



Pandas to_csv with example

Read the
comments
carefully

```
In [15]: a = pd.read_csv("C:/Users/chanukya/Desktop/test1.csv") #I am just importing a dataframe and storing it as "a"
```

```
In [17]: a.to_csv("C:/Users/chanukya/Desktop/test2.csv") #converting the dataframe into csv and storing it with a new filename test2.csv  
# here "a" is the name of the dataframe which was imported, however, you can create your own dataframe with pandas  
# and export it to your desired local folder using the above line of code
```



Comprehensive Tutorial

- Must read-
A comprehensive tutorial on pandas for beginners.

<https://www.learndatasci.com/tutorials/python-pandas-tutorial-complete-introduction-for-beginners/>



Let's Practice!

Now let's practice what we've learnt on a real dataset

The dataset can be downloaded from

<https://data.un.org/Data.aspx?q=greenhouse+gas+co2+emissions&d=GHG&f=seriesID%3aGHG>

The screenshot shows the UN Data website interface. At the top, a purple bar contains the word "Statistics". Below it, a light blue box displays the title "Greenhouse Gas (GHGs) Emissions without Land Use, Land-Use Change and Forestry (LULUCF), in kilotons" and the source "Source: Greenhouse Gas Inventory Data | United Nations Framework Convention on Climate Change". On the left, a "Select filters:" panel shows a list of countries with checkboxes, including Australia, Austria, Belarus, Belgium, and Bulgaria. In the center, a "Download" button is visible. A modal dialog box titled "Select download format:" is open, showing two sections: "Structured" with an "XML" option, and "Value separated" with three options: "Comma" (highlighted with a green icon), "Semicolon" (with a blue icon), and "Pipe" (with a blue icon). A teal callout bubble points to the "Comma" option with the text "Click on Comma to download a comma separated file". The background shows a table with a header "Country or Area" and several rows of data, all of which are "Australia".

Click on Comma to download a comma separated file



Let's Practice!

- Extract the zip file.
- For Jupyter Notebook users:
Place the CSV file that you just extracted in the same folder you're running your Jupyter Notebook from or just copy the folder path and paste it as given in the examples in previous slides.
- For Google Colab users:
 - To upload from your local drive, start with the following code:

```
from google.colab import files
uploaded = files.upload()
```
 - It will prompt you to select a file. Click on "Choose Files" then select and upload the file. Wait for the file to be 100% uploaded. You should see the name of the file once Colab has uploaded it.
 - [Read this article here for more help](#)



```
from google.colab import files
uploaded = files.upload()
```



Choose Files UNdata_Exp...02870.csv

● **UNdata_Export_20200528_120502870.csv**(text/csv) - 45915 bytes, last modified: 28/05/2020 - 100% done
Saving UNdata_Export_20200528_120502870.csv to UNdata_Export_20200528_120502870.csv

Let's Practice!

- Load the CSV file into a variable - `greenhouse_data`
- Use the `description` function to understand how the data looks like
- Print its first 10 rows using `head()`
- Print its last 10 rows using `tail()`
- Check if there are any null values
- Store the DataFrame to a CSV file named 'file2.csv'

Optional Video Tutorial

Ignore matplotlib library for the moment and you may follow the rest



PYTHON + PANDAS TUTORIAL

pt.2

```
read_csv()  
to_csv()
```

That's it for the day. Thank you!

Feel free to post any queries in the #help channel on Slack



Democratizing Data Science Learning