

Data Pre-Processing



What is our GOAL for this MODULE?

The goal of this module is to learn about observing and preprocessing data.

What did we ACHIEVE in the class TODAY?

- We took two different datasets, merged them together and processed the data.

Which CONCEPTS/CODING BLOCKS did we cover today?

- Logic Building
- Critical Thinking
- Python

How did we DO the activities?

1. Take a look at the previous scraped data from last class.

<https://raw.githubusercontent.com/whitehatjr/web-scrapping-2/master/final.csv>

2. Observe that the headers have 11 items while the data has 13 items.
3. Find the missing headers and add them.
4. Go to a new website and download the csv from there.

<https://exoplanetarchive.ipac.caltech.edu/cgi-bin/TblView/nph-tblView?app=ExoTbls&config=planets>

5. Observe the two datasets and try to find a pattern in order to merge them.
6. Conclude that both the datasets are arranged in alphabetical order.
7. Observe that the second dataset has all the lowercase planet names at the bottom of the list separately.

Home About Us Data Tools Support Login													
Select Columns Download Table Plot Table View Documentation User Preferences ⚠ This service is being retired. Please use the Planetary Systems table.													
Confirmed Planets (retiring)													
Host Name	Planet Letter	Planet Name	Discovery Method	Controversial Flag	Number of Planets in System	Orbital Period [days]	Orbit Semi-Major Axis [au]	Eccentricity	Inclination [deg]	Planet Mass or M*sin(i) [Jupiter mass]	Planet Mass or M*sin(i) Provenance	PI	
XO-4	b	XO-4 b	Transit	0	1	4.125080±0.000004		0.0	88.80±0.60	1.42±0.19	Mass	1.25	
XO-5	b	XO-5 b	Transit	0	1	4.1877558±0.0000006	0.0515±0.0005	0	86.8±0.2	1.19±0.03	Mass	1.14	
XO-6	b	XO-6 b	Transit	0	1	3.7650007±0.0000081	0.0815±0.0077	0	86.0±0.2	<4	Mass	2.07	
XO-7	b	XO-7 b	Transit	0	1	2.8641424±0.0000043	0.04421±0.00062	0.038±0.033	83.45±0.29	0.709±0.034	Mass	1.37	
YZ Cet	b	YZ Cet b	Radial Velocity	0	3	1.96876±0.00021	0.01557±0.00052	0.0±0.1		0.0024±0.0004	Msin(i)		
YZ Cet	c	YZ Cet c	Radial Velocity	0	3	3.06008±0.00022	0.02090±0.00070	0.040±0.110		0.00308±0.00044	Msin(i)		
YZ Cet	d	YZ Cet d	Radial Velocity	0	3	4.65627±0.00042	0.02764±0.00093	0.129±0.096		0.00359±0.00053	Msin(i)		
α Ari	b	α Ari b	Radial Velocity	0	1	380.8±0.3	1.2	0.25±0.03		1.8±0.2	Msin(i)		
α Tau	b	α Tau b	Radial Velocity	0	1	628.96±0.90	1.46±0.27	0.10±0.05		6.47±0.53	Msin(i)		
β Cnc	b	β Cnc b	Radial Velocity	0	1	605.2±4.0	1.7±0.1	0.08±0.02		7.8±0.8	Msin(i)		
β Pic	b	β Pic b	Imaging	0	2	7665.0 ^{+7300.0} _{-730.0}	9.10 ^{+5.30} _{-0.50}	0.080 ^{+0.320} _{-0.030}	88.9±0.7				
β Pic	c	β Pic c	Radial Velocity	0	2	1200	2.7	0.24	90	9	Mass		
β UMi	b	β UMi b	Radial Velocity	0	1	522.3±2.7	1.4±0.1	0.19±0.02		6.1±1.0	Msin(i)		
ε CrB	b	ε CrB b	Radial Velocity	0	1	417.9±0.5	1.3	0.11±0.03		6.7±0.3	Msin(i)		
ε Eri	b	ε Eri b	Radial Velocity	0	1	2690±30		0.070 ^{+0.060} _{-0.050}	89.0±42.0	0.780 ^{+0.380} _{-0.120}	Mass		
ε Ind A	b	ε Ind A b	Radial Velocity	0	1	16510 ⁺²¹⁰⁰ ₋₁₇₄₀	11.55 ^{+0.98} _{-0.86}	0.26 ^{+0.07} _{-0.03}	64.25 ^{+13.80} _{-8.09}	3.25 ^{+0.39} _{-0.65}	Mass		
ε Tau	b	ε Tau b	Radial Velocity	0	1	594.9±5.3	1.93±0.03	0.151±0.023		7.6±0.2	Msin(i)		
γ 1 Leo	b	γ 1 Leo b	Radial Velocity	0	1	428.5±1.25	1.19±0.02	0.144±0.046		8.78±1.0	Msin(i)		
γ Cep	b	γ Cep b	Radial Velocity	0	1	903.3±1.5	2.05±0.06	0.049±0.034		1.85±0.16	Msin(i)		
γ Lib	b	γ Lib b	Radial Velocity	0	2	415.2 ^{+1.8} _{-1.9}	1.24±0.10	0.21±0.10		1.02±0.14	Msin(i)		
γ Lib	c	γ Lib c	Radial Velocity	0	2	964.6±3.1	2.17±0.10	0.057 ^{+0.034} _{-0.032}		4.56 ^{+0.45} _{-0.43}	Msin(i)		
ι Dra	b	ι Dra b	Radial Velocity	0	1	511.098±0.089	1.275±0.074	0.7124±0.0039		8.82±0.72	Msin(i)		
κ And	b	κ And b	Imaging	0	1		55±2			13.616 ^{+23.042} _{-1.047}	Mass		

Showing records 4241 to 4263 of 4284 (4284 total) DOI 10.26133/NEA1

[Clear Checked](#) [Check All](#) [Reset Filters](#)

8. Create a script that can arrange the second dataset alphabetically irrespective of

- them being lowercase or uppercase.
9. Create a new directory and create a virtual environment.
 10. Move the downloaded CSV to this new folder and rename it to dataset_1.csv.
 11. Think about the logic.
 - **Fun Fact: All the lowercase letters have a higher ASCII value (lowercase a has an ASCII value of 97) and all the uppercase letters have a lower ASCII value (uppercase A has an ASCII value of 65).**
 12. Using this fun fact, we conclude that we will convert all planet names (3rd column of the CSV) into either uppercase or lowercase and then sort the data.
 13. Write the code on the logic we thought

```
import csv

data = []

with open("dataset_2.csv", "r") as f:
    csvreader = csv.reader(f)
    for row in csvreader:
        data.append(row)

headers = data[0]
planet_data = data[1:]

#Converting all planet names to lowercase
for data_point in planet_data:
    data_point[2] = data_point[2].lower()

#Sorting planet names in alphabetical order
planet_data.sort(key=lambda planet_data: planet_data[2])

with open("dataset_2_sorted.csv", "a+") as f:
    csvwriter = csv.writer(f)
    csvwriter.writerow(headers)
    csvwriter.writerows(planet_data)
```

14. Move the CSV of the data we scraped and rename it to dataset_1.csv.
15. Write a script to merge the two datasets and create a final.csv.

```
import csv
```

```
dataset_1 = []
dataset_2 = []

with open("dataset_1.csv", "r") as f:
    csvreader = csv.reader(f)
    for row in csvreader:
        dataset_1.append(row)

with open("dataset_2_sorted.csv", "r") as f:
    csvreader = csv.reader(f)
    for row in csvreader:
        dataset_2.append(row)

headers_1 = dataset_1[0]
planet_data_1 = dataset_1[1:]

headers_2 = dataset_2[0]
planet_data_2 = dataset_2[1:]

headers = headers_1 + headers_2
planet_data = []
for index, data_row in enumerate(planet_data_1):
    planet_data.append(planet_data_1[index] + planet_data_2[index])

with open("final.csv", "a+") as f:
    csvwriter = csv.writer(f)
    csvwriter.writerow(headers)
    csvwriter.writerows(planet_data)
```

What's NEXT?

In the next class, we will learn about data-cleaning. It is the final step before we start performing statistical analysis on our data.

EXTEND YOUR KNOWLEDGE:

You can read the following blog on data processing to understand more:

<https://bigdataanalyticsnews.com/data-preparation-why-is-it-important/>