

# Assignment 1:

Arun Agarwal

915774866

## Homework 1

1/19/2022 - 1/23/2022

### Part 1: Data Scraping and Preparation:

#### Task 1: Scrape your Competitor's Data:

```
In [ ]: #Installing Necessary Packages:
pip install bs4
pip install matplotlib
pip install requests
pip install pandas
pip install numpy

In [ ]: #Code added to convert Jupyter Notebook to pdf
python -m pip install notebook==pdf
pyppeteer-install

In [2]: pip install requests

Requirement already satisfied: requests in c:\users\aruna\desktop\lib\site-packages (2.24.0)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\aruna\desktop\lib\site-packages (from requests) (2020.6.20)
Requirement already satisfied: idna<3,>=2.5 in c:\users\aruna\desktop\lib\site-packages (from requests) (2.10)
Requirement already satisfied: chardet<4,>=3.0.2 in c:\users\aruna\desktop\lib\site-packages (from requests) (3.0.4)
Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in c:\users\aruna\desktop\lib\site-packages (from requests) (1.25.11)
Note: you may need to restart the kernel to use updated packages.

In [1]: #Making the necessary imports, as requested in the instructions:
from bs4 import BeautifulSoup as bs
import requests as rq
import pandas as pd
import numpy as np
import matplotlib as mp #For Visualization

In [2]: #Using requests to get page content (as in, HTTP GET)
#I use the requests package get method to get this page content
#We need to add the headers argument to be able to get the data
r = rq.get('https://www.spaceweatherlive.com/en/solar-activity/top-50-solar-flares.html', headers = {
    'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/80.0.4012.91 Safari/537.36'})
#I extract the text from the page using the text function
ourText = r.text

In [3]: #Now, I use the BeautifulSoup package to read and parse the data as html
#There are multiple options to place as the second parametere for bs, such as 'html5lib',
#but I chose to use 'html.parser'
theSoup = bs(ourText, 'html.parser')

In [4]: #I use prettify() function of BeautifulSoup to view the content and find the appropriate table
theSoup.prettify()
```



[illegible]







[illegible]



https://cdaw.gsfc.nasa.gov/CME\_list/UNIVERSAL/2002/04/jsmovie/2002\_04/20020417\_082605\_p932/cdf.html?P0=0&P1=0&P2=0&P3=0&P4=0&P5=0&P6=0&P7=0&P8=0&P9=0&P10=0&P11=0&P12=0&P13=0&P14=0&P15=0&P16=0&P17=0&P18=0&P19=0&P20=0&P21=0&P22=0&P23=0&P24=0&P25=0&P26=0&P27=0&P28=0&P29=0&P30=0&P31=0&P32=0&P33=0&P34=0&P35=0&P36=0&P37=0&P38=0&P39=0&P40=0&P41=0&P42=0&P43=0&P44=0&P45=0&P46=0&P47=0&P48=0&P49=0&P50=0&P51=0&P52=0&P53=0&P54=0&P55=0&P56=0&P57=0&P58=0&P59=0&P60=0&P61=0&P62=0&P63=0&P64=0&P65=0&P66=0&P67=0&P68=0&P69=0&P70=0&P71=0&P72=0&P73=0&P74=0&P75=0&P76=0&P77=0&P78=0&P79=0&P80=0&P81=0&P82=0&P83=0&P84=0&P85=0&P86=0&P87=0&P88=0&P89=0&P90=0&P91=0&P92=0&P93=0&P94=0&P95=0&P96=0&P97=0&P98=0&P99=0&P100=0&P101=0&P102=0&P103=0&P104=0&P105=0&P106=0&P107=0&P108=0&P109=0&P110=0&P111=0&P112=0&P113=0&P114=0&P115=0&P116=0&P117=0&P118=0&P119=0&P120=0&P121=0&P122=0&P123=0&P124=0&P125=0&P126=0&P127=0&P128=0&P129=0&P130=0&P131=0&P132=0&P133=0&P134=0&P135=0&P136=0&P137=0&P138=0&P139=0&P140=0&P141=0&P142=0&P143=0&P144=0&P145=0&P146=0&P147=0&P148=0&P149=0&P150=0&P151=0&P152=0&P153=0&P154=0&P155=0&P156=0&P157=0&P158=0&P159=0&P160=0&P161=0&P162=0&P163=0&P164=0&P165=0&P166=0&P167=0&P168=0&P169=0&P170=0&P171=0&P172=0&P173=0&P174=0&P175=0&P176=0&P177=0&P178=0&P179=0&P180=0&P181=0&P182=0&P183=0&P184=0&P185=0&P186=0&P187=0&P188=0&P189=0&P190=0&P191=0&P192=0&P193=0&P194=0&P195=0&P196=0&P197=0&P198=0&P199=0&P200=0&P201=0&P202=0&P203=0&P204=0&P205=0&P206=0&P207=0&P208=0&P209=0&P210=0&P211=0&P212=0&P213=0&P214=0&P215=0&P216=0&P217=0&P218=0&P219=0&P220=0&P221=0&P222=0&P223=0&P224=0&P225=0&P226=0&P227=0&P228=0&P229=0&P230=0&P231=0&P232=0&P233=0&P234=0&P235=0&P236=0&P237=0&P238=0&P239=0&P240=0&P241=0&P242=0&P243=0&P244=0&P245=0&P246=0&P247=0&P248=0&P249=0&P250=0&P251=0&P252=0&P253=0&P254=0&P255=0&P256=0&P257=0&P258=0&P259=0&P260=0&P261=0&P262=0&P263=0&P264=0&P265=0&P266=0&P267=0&P268=0&P269=0&P270=0&P271=0&P272=0&P273=0&P274=0&P275=0&P276=0&P277=0&P278=0&P279=0&P280=0&P281=0&P282=0&P283=0&P284=0&P285=0&P286=0&P287=0&P288=0&P289=0&P290=0&P291=0&P292=0&P293=0&P294=0&P295=0&P296=0&P297=0&P298=0&P299=0&P300=0&P301=0&P302=0&P303=0&P304=0&P305=0&P306=0&P307=0&P308=0&P309=0&P310=0&P311=0&P312=0&P313=0&P314=0&P315=0&P316=0&P317=0&P318=0&P319=0&P320=0&P321=0&P322=0&P323=0&P324=0&P325=0&P326=0&P327=0&P328=0&P329=0&P330=0&P331=0&P332=0&P333=0&P334=0&P335=0&P336=0&P337=0&P338=0&P339=0&P340=0&P341=0&P342=0&P343=0&P344=0&P345=0&P346=0&P347=0&P348=0&P349=0&P350=0&P351=0&P352=0&P353=0&P354=0&P355=0&P356=0&P357=0&P358=0&P359=0&P360=0&P361=0&P362=0&P363=0&P364=0&P365=0&P366=0&P367=0&P368=0&P369=0&P370=0&P371=0&P372=0&P373=0&P374=0&P375=0&P376=0&P377=0&P378=0&P379=0&P380=0&P381=0&P382=0&P383=0&P384=0&P385=0&P386=0&P387=0&P388=0&P389=0&P390=0&P391=0&P392=0&P393=0&P394=0&P395=0&P396=0&P397=0&P398=0&P399=0&P400=0&P401=0&P402=0&P403=0&P404=0&P405=0&P406=0&P407=0&P408=0&P409=0&P410=0&P411=0&P412=0&P413=0&P414=0&P415=0&P416=0&P417=0&P418=0&P419=0&P420=0&P421=0&P422=0&P423=0&P424=0&P425=0&P426=0&P427=0&P428=0&P429=0&P430=0&P431=0&P432=0&P433=0&P434=0&P435=0&P436=0&P437=0&P438=0&P439=0&P440=0&P441=0&P442=0&P443=0&P444=0&P445=0&P446=0&P447=0&P448=0&P449=0&P450=0&P451=0&P452=0&P453=0&P454=0&P455=0&P456=0&P457=0&P458=0&P459=0&P460=0&P461=0&P462=0&P463=0&P464=0&P465=0&P466=0&P467=0&P468=0&P469=0&P470=0&P471=0&P472=0&P473=0&P474=0&P475=0&P476=0&P477=0&P478=0&P479=0&P480=0&P481=0&P482=0&P483=0&P484=0&P485=0&P486=0&P487=0&P488=0&P489=0&P490=0&P491=0&P492=0&P493=0&P494=0&P495=0&P496=0&P497=0&P498=0&P499=0&P500=0&P501=0&P502=0&P503=0&P504=0&P505=0&P506=0&P507=0&P508=0&P509=0&P510=0&P511=0&P512=0&P513=0&P514=0&P515=0&P516=0&P517=0&P5



(1)







[illegible]

```

task 4: How the NASA table.

In [15]:
#First, we recode the missing entries at NaN. After looking carefully at the NASA data and the data description
#I have found that:
#---- and ---- are missing data entries that can exist in multiple columns (for Flare and CME data),
#???? is a missing data entry that can appear in the frequency columns and indicates that the starting
#and ending frequencies are not determined,
# -/- is a missing data entry that can appear in the date column for CME,
# and ---- is a missing data entry that can appear in the time column for CME.
#Replacing these missing entries:
tableNASAdf = tableNASAdf.replace('----', np.nan)
tableNASAdf = tableNASAdf.replace('----', np.nan)
tableNASAdf = tableNASAdf.replace('----', np.nan)
tableNASAdf = tableNASAdf.replace('????', np.nan)
tableNASAdf = tableNASAdf.replace('-/-', np.nan)
tableNASAdf = tableNASAdf.replace('----', np.nan)
tableNASAdf = tableNASAdf.replace('----', np.nan)
tableNASAdf = tableNASAdf.replace('360h', value = '360') #there was a typo in the table in which I

tableNASAdf

Out[15]:


|     | start_date | start_time | end_date | end_time | start_frequency | end_frequency | flare_location | flare_region | flare_classification | cme_date | cme_time |
|-----|------------|------------|----------|----------|-----------------|---------------|----------------|--------------|----------------------|----------|----------|
| 0   | 1997/04/01 | 1400       | 04/01    | 1415     | 8000            | 4000          | S25E16         | 8026         | M1.3                 | M1.3     | 04/01    |
| 1   | 1997/04/07 | 1430       | 04/07    | 1730     | 11000           | 1000          | S28E19         | 8027         | C6.8                 | C6.8     | 04/07    |
| 2   | 1997/05/12 | 0515       | 05/14    | 1600     | 12000           | 80            | N21W08         | 8038         | C1.3                 | C1.3     | 05/12    |
| 3   | 1997/05/21 | 2020       | 05/21    | 2200     | 5000            | 500           | N05W12         | 8040         | M1.3                 | M1.3     | 05/21    |
| 4   | 1997/09/23 | 2153       | 09/23    | 2216     | 6000            | 2000          | S29E25         | 8088         | C1.4                 | C1.4     | 09/23    |
| ... | ...        | ...        | ...      | ...      | ...             | ...           | ...            | ...          | ...                  | ...      | ...      |
| 517 | 2017/09/17 | 1145       | 09/17    | 1235     | 16000           | 900           | S08E170        | NaN          | NaN                  | NaN      | 09/17    |
| 518 | 2017/10/18 | 0548       | 10/18    | 1240     | 1600            | 400           | S06E123        | NaN          | NaN                  | NaN      | 10/18    |
| 519 | 2019/05/03 | 2352       | 05/04    | 0016     | 13000           | 2300          | N12E82         | 12740        | C1.0                 | C1.0     | 05/03    |
| 520 | 2020/11/29 | 1307       | 11/29    | 1523     | 14000           | 850           | S23E89         | NaN          | M4.4                 | M4.4     | 11/29    |
| 521 | 2020/12/07 | 1618       | 12/08    | 0200     | 14000           | 160           | S25W08         | 12790        | C7.4                 | C7.4     | 12/07    |


522 rows x 14 columns

In [16]:
#The CFA column (cme_angle) contains angles in degrees for most rows, except for halo flares,
#which are coded as Halo.
#I create a new column that indicates if a row corresponds to a halo flare or not,
#and then replace halo entries in the cme_angle column with NaN:
#Creating and filling new 'is_halo' column as indicated:
tableNASAdf['is_halo'] = tableNASAdf['cme_angle'].map(lambda x: x == 'Halo')
#Replacing "Halo" with "NaN"
tableNASAdf = tableNASAdf.replace('Halo', np.nan)
tableNASAdf

Out[16]:


|     | start_date | start_time | end_date | end_time | start_frequency | end_frequency | flare_location | flare_region | flare_classification | cme_date | cme_time |
|-----|------------|------------|----------|----------|-----------------|---------------|----------------|--------------|----------------------|----------|----------|
| 0   | 1997/04/01 | 1400       | 04/01    | 1415     | 8000            | 4000          | S25E16         | 8026         | M1.3                 | M1.3     | 04/01    |
| 1   | 1997/04/07 | 1430       | 04/07    | 1730     | 11000           | 1000          | S28E19         | 8027         | C6.8                 | C6.8     | 04/07    |
| 2   | 1997/05/12 | 0515       | 05/14    | 1600     | 12000           | 80            | N21W08         | 8038         | C1.3                 | C1.3     | 05/12    |
| 3   | 1997/05/21 | 2020       | 05/21    | 2200     | 5000            | 500           | N05W12         | 8040         | M1.3                 | M1.3     | 05/21    |
| 4   | 1997/09/23 | 2153       | 09/23    | 2216     | 6000            | 2000          | S29E25         | 8088         | C1.4                 | C1.4     | 09/23    |
| ... | ...        | ...        | ...      | ...      | ...             | ...           | ...            | ...          | ...                  | ...      | ...      |
| 517 | 2017/09/17 | 1145       | 09/17    | 1235     | 16000           | 900           | S08E170        | NaN          | NaN                  | NaN      | 09/17    |
| 518 | 2017/10/18 | 0548       | 10/18    | 1240     | 1600            | 400           | S06E123        | NaN          | NaN                  | NaN      | 10/18    |
| 519 | 2019/05/03 | 2352       | 05/04    | 0016     | 13000           | 2300          | N12E82         | 12740        | C1.0                 | C1.0     | 05/03    |
| 520 | 2020/11/29 | 1307       | 11/29    | 1523     | 14000           | 850           | S23E89         | NaN          | M4.4                 | M4.4     | 11/29    |
| 521 | 2020/12/07 | 1618       | 12/08    | 0200     | 14000           | 160           | S25W08         | 12790        | C7.4                 | C7.4     | 12/07    |


522 rows x 15 columns

In [17]:
#A helper function to use in map function below to determine if a given string is an integer/number:
def isInteger(string):
    try:
        #Check if the string is a float because NaN values are floats; integers will work here as well
        float(string)
        return True
    except ValueError:
        return False

In [18]:
#The width column indicates if the given value is a lower bound.
#I create a new column that indicates if width is given as a lower bound, and remove any non-numeric part
for the width columns:
#Creating and filling new 'width_lower_bound' column as indicated:
tableNASAdf['width_lower_bound'] = tableNASAdf['cme_width'].map(lambda x: str(x)(0) == '>')
#Removing any non-numeric part of the width column:
tableNASAdf['cme_width'] = tableNASAdf['cme_width'].map(lambda x: x if isInteger(x) else x(1))
#Where, I display the first 50 columns to clearly demonstrate that we have completed the task as requested.
tableNASAdf.head(50)

Out[18]:


|   | start_date | start_time | end_date | end_time | start_frequency | end_frequency | flare_location | flare_region | flare_classification | cme_date | cme_time |
|---|------------|------------|----------|----------|-----------------|---------------|----------------|--------------|----------------------|----------|----------|
| 0 | 1997/04/01 | 1400       | 04/01    | 1415     | 8000            | 4000          | S25E16         | 8026         | M1.3                 | M1.3     | 04/01    |
| 1 | 1997/04/07 | 1430       | 04/07    | 1730     | 11000           | 1000          | S28E19         | 8027         | C6.8                 | C6.8     | 04/07    |
| 2 | 1997/05/12 | 0515       | 05/14    | 1600     | 12000           | 80            | N21W08         | 8038         | C1.3                 | C1.3     | 05/12    |
| 3 | 1997/05/21 | 2020       | 05/21    | 2200     | 5000            | 500           | N05W12         | 8040         | M1.3                 | M1.3     | 05/21    |
| 4 | 1997/09/23 | 2153       | 09/23    | 2216     | 6000            | 2             |                |              |                      |          |          |


```



	start_datetime	end_datetime	start_frequency	end_frequency	flare_location	flare_region	flare_classification	cme_datetime	cme_angle
0	1997-04-01 14:00:00	1997-04-01 14:15:00	8000	4000	S2E16	8026	M1.3	1997-04-01 15:18:00	74
1	1997-04-07 14:30:00	1997-04-07 17:30:00	1000	1000	S2E19	8078	C6.8	1997-04-07 14:27:00	NaN
2	1997-05-12 05:15:00	1997-05-14 16:00:00	12000	80	N21W08	8038	C1.3	1997-05-12 05:30:00	NaN
3	1997-05-21 20:00:00	1997-05-21 22:00:00	5000	500	N05W12	8040	M1.3	1997-05-21 21:00:00	263
4	1997-09-23 21:53:00	1997-09-23 22:16:00	6000	2000	S2E25	8088	C1.4	1997-09-23 22:02:00	133
...	...	...	...	...	...	...	...	...	...
517	2017-09-18 11:45:00	2017-09-18 12:35:00	16000	900	S08E170	NaN	NaN	2017-09-18 12:00:00	NaN
518	2017-10-18 05:48:00	2017-10-18 12:40:00	16000	400	S06E123	NaN	NaN	2017-10-18 08:00:00	85
519	2019-05-03 21:52:00	2019-05-04 00:16:00	13005	2300	N12E82	12740	C1.0	2019-05-03 23:24:00	90
520	2019-11-29 13:07:00	2019-11-29 15:23:00	14000	850	S23W08	NaN	NaN	2019-11-29 13:25:00	NaN
521	2020-12-07 16:18:00	2020-12-08 02:00:00	14000	160	S25W08	12790	C7.4	2020-12-07 16:24:00	NaN

522 rows x 13 columns

## Part 2: Analysis:

### Task 5: Replication:

```
In [23]: #We need to get JUST the top 50 solar flares from the NASA based on their classification
#(with X28 being the highest). Thus, we will create a copy of our tableNASAdef for temporary use
#Finally, to remove the NaN values for their flare_classification, we will use the following code:
tableNASAdefcopy = tableNASAdef.copy()
tableNASAdefcopy = tableNASAdefcopy.dropna(subset = ['flare_classification'])
tableNASAdefcopy
```

	start_datetime	end_datetime	start_frequency	end_frequency	flare_location	flare_region	flare_classification	cme_datetime	cme_angle
0	1997-04-01 14:00:00	1997-04-01 14:15:00	8000	4000	S2E16	8026	M1.3	1997-04-01 15:18:00	74
1	1997-04-07 14:30:00	1997-04-07 17:30:00	1000	1000	S2E19	8078	C6.8	1997-04-07 14:27:00	NaN
2	1997-05-12 05:15:00	1997-05-14 16:00:00	12000	80	N21W08	8038	C1.3	1997-05-12 05:30:00	NaN
3	1997-05-21 20:00:00	1997-05-21 22:00:00	5000	500	N05W12	8040	M1.3	1997-05-21 21:00:00	263
4	1997-09-23 21:53:00	1997-09-23 22:16:00	6000	2000	S2E25	8088	C1.4	1997-09-23 22:02:00	133
...	...	...	...	...	...	...	...	...	...
517	2017-09-18 11:45:00	2017-09-18 12:35:00	16000	900	S08E170	NaN	NaN	2017-09-18 12:00:00	NaN
518	2017-10-18 05:48:00	2017-10-18 12:40:00	16000	400	S06E123	NaN	NaN	2017-10-18 08:00:00	85
519	2019-05-03 21:52:00	2019-05-04 00:16:00	13005	2300	N12E82	12740	C1.0	2019-05-03 23:24:00	90
520	2019-11-29 13:07:00	2019-11-29 15:23:00	14000	850	S23W08	NaN	NaN	2019-11-29 13:25:00	NaN
521	2020-12-07 16:18:00	2020-12-08 02:00:00	14000	160	S25W08	12790	C7.4	2020-12-07 16:24:00	NaN

407 rows x 13 columns

```
In [24]: #Now we will grab flares/corow that have a flare_classification with X (and store it in a new dataframe):
#topFlaredef = tableNASAdefcopy.loc[tableNASAdefcopy['flare_classification'].str.contains('X')]
#For ease of use, I will get rid of the X character at the beginning and add it back after sorting:
topFlaredef['flare_classification'] = topFlaredef['flare_classification'].str.strip('X').astype(float)
#Now, sorting the rows according to the flare_classification:
topFlaredef = topFlaredef.sort_values('flare_classification', ascending = False)
#Placing just the top 50:
topFlaredef = topFlaredef.head(50)
#Placing back the 'X' character for the flare_classification:
topFlaredef['flare_classification'] = topFlaredef['flare_classification'].astype(str)
topFlaredef['flare_classification'] = 'X' + topFlaredef['flare_classification']
topFlaredef
```

	start_datetime	end_datetime	start_frequency	end_frequency	flare_location	flare_region	flare_classification	cme_datetime	cme_angle
240	2019-11-04 20:00:00	2019-11-05 00:00:00	10000	200	S19W83	10486	X2.0	2019-11-04 19:54:00	NaN
117	2001-04-02 22:05:00	2001-04-03 02:30:00	14000	250	N19W72	9993	X2.0	2001-04-02 22:06:00	261
233	2003-10-28 11:10:00	2003-10-30 00:00:00	14000	40	S16E08	10486	X17.0	2003-10-28 11:30:00	NaN
126	2001-04-15 14:05:00	2001-04-16 13:00:00	14000	40	S20W85	9415	X14.0	2001-04-15 14:06:00	245
234	2003-10-29 20:55:00	2003-10-30 00:30:00	11000	50	S15W02	10486	X10.0	2003-10-29 20:54:00	NaN
8	1997-11-06 12:20:00	1997-11-07 08:30:00	14000	100	S18W63	8100	X9.4	1997-11-06 12:10:00	NaN
514	2017-09-06 12:05:00	2017-09-07 08:00:00	14000	70	S08W33	12673	X9.3	2017-09-06 12:24:00	NaN
328	2006-12-05 10:50:00	2006-12-05 20:00:00	14000	250	S07E68	10930	X8.0	NaT	NaN
237	2003-11-02 17:30:00	2003-11-03 01:00:00	12000	250	S19W56	10486	X8.3	2003-11-02 17:30:00	NaN
515	2017-09-10 16:02:00	2017-09-11 06:50:00	14000	150	S09W52	9415	X8.3	2017-09-10 16:00:00	NaN
288	2005-01-20 01:08:00	2005-01-20 16:30:00	14000	25	N14W61	10720	X7.1	2005-01-20 06:54:00	NaN
359	2017-08-09 08:20:00	2017-08-09 08:35:00	16000	4000	N17W69	11263	X6.9	2017-08-09 08:12:00	NaN
331	2006-12-06 19:00:00	2006-12-09 00:00:00	16000	30	S05E64	10930	X6.5	NaT	NaN
317	2005-09-09 19:45:00	2005-09-09 22:00:00	10000	50	S12E67	10808	X6.2	2005-09-09 19:48:00	NaN
82	2000-07-14 10:30:00	2000-07-15 14:30:00	14000	80	N22W07	9077	X5.7	2000-07-14 10:54:00	NaN
121	2001-04-06 19:35:00	2001-04-07 01:50:00	14000	230	S21E31	9415	X5.6	2001-04-06 19:30:00	NaN
375	2012-03-07 01:00:00	2012-03-08 19:00:00	16000	30	N17E27	11429	X5.4	2012-03-07 02:24:00	NaN
135	2001-08-25 16:50:00	2001-08-25 23:00:00	8000	170	S17E34	9591	X5.3	2001-08-25 16:50:00	NaN
443	2014-02-25 05:06:00	2014-02-25 11:28:00	14000	100	S12E82	11990	X4.9	2014-02-25 07:25:00	NaN
193	2002-07-23 00:50:00	2002-07-23 04:00:00	11000	400	S13E72	10039	X4.8	2002-07-23 04:26:00	NaN
104	2000-11-26 17:00:00	2000-11-26 17:15:00	14000	7000	N18W38	9236	X4.0	2000-11-26 17:15:00	NaN
239	2003-11-03 10:00:00	2003-11-03 12:30:00	6000	400	N08W77	10488	X3.9	2003-11-03 10:06:00	293
286	2005-01-17 10:00:00	2005-01-17 10:35:00	6100	1500	N15W25	10720	X3.8	2005-01-17 09:54:00	NaN
222	2003-05-28 01:00:00	2003-05-29 03:30:00	1000	200	S07W20	10365	X3.6	2003-05-28 01:50:00	NaN
332	2006-12-13 02:45:00	2006-12-13 10:40:00	12000	150	S06W23	10930	X3.4	2006-12-13 02:54:00	NaN
160	2001-12-28 20:35:00	2001-12-29 03:00:00	14000	350	S26E90	9756	X3.4	2001-12-28 20:30:00	NaN
192	2002-07-20 21:30:00	2002-07-20 22:20:00	10000	2000	S13E90	10039	X3.3	2002-07-20 22:06:00	NaN
404	2013-05-14 01:16:00	2013-05-14 08:20:00	16000	240	N08E77	11748	X3.2	2013-05-14 01:25:00	NaN
201	2002-08-24 01:45:00	2002-08-24 03:25:00	5000	400	S02W81	10069	X3.1	2002-08-24 01:27:00	NaN
403	2013-05-13 16:55:00	2013-05-13 19:10:00	16000	300	N11E85	11748	X2.8	2013-05-13 16:50:00	NaN
487	2015-05-05 22:24:00	2015-05-05 23:14:00	14000	500	N15E79	12339	X2.7	2015-05-05 22:42:00	NaN
19	1998-05-06 08:25:00	1998-05-06 08:35:00	14000	5000	S10W63	8218	X2.7	1998-05-06 08:29:00	309
238	2003-11-03 01:05:00	2003-11-03 01:25:00	3000	1500	N11W65	10486	X2.7	2003-11-03 01:59:00	304
284	2005-01-15 23:00:00	2005-01-17 00:00:00	3000	40	N15W05	10720	X2.6	2005-01-15 23:06:00	NaN
142	2001-09-24 10:45:00	2001-09-25 10:00:00	7000	30	S16E23	9632	X2.6	2001-09-24 10:30:00	NaN
9	1997-11-27 13:30:00	1997-11-27 14:00:00	14000	7000	N17E63	8113	X2.6	1997-11-27 13:10:00	98
276	2004-11-10 02:25:00	2004-11-10 03:40:00	14000	1000	N09W49	10696	X2.5	2004-11-10 02:26:00	NaN
123	2001-04-10 05:24:00	2001-04-10 00:00:00	14000	100	S23W09	9415	X2.3	2001-04-10 05:30:00	NaN
99	2000-11-24 15:25:00	2000-11-24 22:00:00	14000	200	N22W07	9236	X2.3	2000-11-24 15:30:00	NaN
73	2000-06-06 15:20:00	2000-06-08 09:00:00	14000	40	N20E18	9026	X2.3	2000-06-06 15:54:00	NaN
345	2011-02-15 02:10:00	2011-02-15 07:00:00	16000	400	S20W12	11158	X2.2	2011-02-15 02:24:00	NaN
318	2005-09-10 21:45:00	2005-09-11 01:00:00	14000	200	S13E47	10808	X2.1	2005-09-09 21:50:00	NaN
361	2011-09-06 23:30:00	2011-09-06 15:40:00	16000	150	N14W18	11283	X2.1	2011-09-06 23:05:00	NaN
420	2013-10-25 15:08:00	2013-10-25 22:30:00	16000	200	S06E69	11882	X2.1	2013-10-25 15:12:00	NaN
7	1997-11-04 06:00:00	1997-11-05 04:30:00	14000	100	S14W33	8100	X2.1	1997-11-04 06:11:00	NaN
98	2000-11-24 05:10:00	2000-11-24 15:00:00	14000	100	N20W53	9236	X2.0	2000-11-24 05:30:00	NaN
125	2001-04-12 10:20:00	2001-04-12 10:40:00	14000	7000	S19W43	9415	X2.0	2001-04-12 10:37:00	NaN
274	2004-11-07 16:25:00	2004-11-08 12:00:00	14000	60	N09W17	10696	X2.0	2004-11-07 16:20:00	NaN
285	2000-07-19 09:25:00	2000-07-19 16:00:00	14000	30	N15W25	10720	X2.0	2000-07-19 09:30:00	NaN
102	2001-11-25 19:00:00	2001-11-25 19:35:00	6000	2000	N20W23	9236	X1.9	2001-11-25 19:31:00	NaN

```
In [25]: #To get a table that is directly comparable to the SpaceWeatherlive table (myTable),
#we will just look at the 'flare_classification', 'start_datetime', 'cme_datetime', 'end_datetime', and 'flare_region'
topFlaredef = topFlaredef.reset_index()
#Finally, to remove the NaN values for their flare_classification, we will use the following code:
#Renaming some of the columns to match the names seen in the SWL data:
topFlaredef['flare_classification'] = topFlaredef['flare_classification'].str.strip('X').astype(float)
#If an not used, we will not need the NaN column in the SpaceWeatherlive table because they are already sorted
#in decreasing order, so we will drop that column in a copy of the dataframe.
tableSWL = myTable.drop(columns = 'rank')
#Comparing the two tables side by side (in a concatenated table):
dCompare = pd.concat([tableSWL, topFlaredef], axis = 1, keys = ['SWL Data', 'NASA Data'])
dCompare
```

404	2013-08-24 01:16:00	2002-08-24 08:20:00	16000	240	N08E77	11748	X3.2	2013-08-24 01:25:00	NaN	
201	2013-08-24 01:45:00	2002-08-24 03:25:00	5000	400	S02W81	10069	X3.1	2013-08-24 01:27:00	NaN	
403	2013-05-13 16:15:00	2013-05-13 19:10:00	16000	300	N1E85	11748	X2.8	2013-05-13 16:07:00	NaN	
487	2015-05-22 22:24:00	2015-05-23 23:14:00	1500	500	N15E79	12339	X2.7	2015-05-23 22:24:00	NaN	
19	1998-05-06 08:25:00	1998-05-06 08:35:00	14000	5000	S11W65	8210	X2.7	1998-05-06 08:29:00	309	
238	2013-11-03 01:15:00	2013-11-03 01:25:00	3000	1500	N10W83	10488	X2.7	2013-11-03 01:15:00	304	
284	2005-01-15 23:00:00	2005-01-17 00:00:00	3000	40	N15W05	10720	X2.6	2005-01-15 23:06:00	NaN	
142	2001-09-24 10:45:00	2001-09-25 20:00:00	7000	30	S16E23	9632	X2.6	2001-09-24 10:30:00	NaN	
9	1997-11-27 13:30:00	1997-11-27 14:00:00	14000	7000	N17E63	8113	X2.6	1997-11-27 13:56:00	98	
276	2004-11-10 02:25:00	2004-11-10 03:40:00	14000	1000	N09W49	10696	X2.5	2004-11-10 02:30:00	NaN	
123	2001-04-10 05:30:00	2001-04-11 00:00:00	14000	100	S23W09	9415	X2.3	2001-04-10 05:30:00	NaN	
99	2000-11-24 15:25:00	2000-11-24 22:00:00	14000	200	N22W07	9236	X2.3	2000-11-24 15:30:00	NaN	
73	2000-06-16 09:00:00	2000-06-18 09:00:00	14000	40	N20E18	9026	X2.3	2000-06-16 09:00:00	NaN	
345	2011-02-15 02:00:00	2011-02-15 07:00:00	16000	400	S20W12	11158	X2.2	2011-02-15 02:24:00	NaN	
318	2005-09-11 21:45:00	2005-09-11 01:00:00	14000	200	S13E47	10808	X2.1	2005-09-10 21:52:00	NaN	
361	2011-09-06 23:30:00	2011-09-07 15:40:00	16000	150	N14W18	11283	X2.1	2011-09-06 23:05:00	NaN	
420	2013-10-25 15:08:00	2013-10-25 22:32:00	16000	200	S06E69	11882	X2.1	2013-10-25 15:12:00	NaN	
7	1997-11-04 06:00:00	1997-11-05 04:30:00	14000	100	S14W33	8100	X2.1	1997-11-04 06:10:00	NaN	
98	2000-11-24 05:10:00	2000-11-24 15:00:00	14000	100	N20W05	9236	X2.0	2000-11-24 05:30:00	NaN	
125	2001-04-12 10:20:00	2001-04-12 10:40:00	14000	7000	S19W43	9415	X2.0	2001-04-12 10:31:00	NaN	
274	2004-11-07 16:20:00	2004-11-08 20:00:00	14000	60	N09W17	10696	X2.0	2004-11-07 16:54:00	NaN	
285	2005-01-17 09:25:00	2005-01-17 16:00:00	14000	30	N15W25	10720	X2.0	2005-01-17 09:30:00	NaN	
102	2000-11-25 19:00:00	2000-11-25 19:35:00	6000	2000	N20W23	9236	X1.9	2000-11-25 19:31:00	NaN	
Out [25]:										
<pre>#To get a table that is directly comparable to the SpaceWeatherlive table (myTable), #we will just look at the flare_classification, start_datetime, cme_datetime, end_datetime, and flare_region topLabeleddf = topLabeled.reset_index() topLabeleddf['x_classification'] = topLabeleddf['start_datetime', 'cme_datetime', 'end_datetime', 'flare_r #Renaming some of the columns to match the names seen in the SWL data topLabeleddf = topLabeleddf.rename(columns={'flare_classification': 'x_classification', 'cme_datetime': 'max_da #We do not really need the rank column in the rank column in the rank column because they are already sorted #in decreasing order, so we will drop that column in a copy of the dataframe. finalSWL = myTable.drop(columns = ['rank'])</pre>										
<pre>#Comparing the two tables side by side (in a concatenated table): dfCompare = pd.concat([tableSWL, topLabeleddf], axis=1, keys=['SWL Data', 'NASA Data']) dfCompare</pre>										
Out [26]:										
	x_classification	start_datetime	max_datetime	end_datetime	region	x_classification	start_datetime	max_datetime	end_datetime	region
0	X28+	2013-11-04 19:29:00	2013-11-04 19:53:00	2013-11-04 20:06:00	486	X28.0	2013-11-04 20:06:00	2013-11-04 20:10:00	2013-11-05 00:00:00	10486
1	X20+	2001-04-02 21:32:00	2001-04-02 21:51:00	2001-04-02 22:00:00	9393	X20.0	2001-04-02 22:05:00	2001-04-02 22:06:00	2001-04-03 00:30:00	9393
2	X172+	2003-10-28 09:51:00	2003-10-28 11:10:00	2003-10-28 11:24:00	486	X17.0	2003-10-28 11:10:00	2003-10-28 11:30:00	2003-10-30 00:00:00	10486
3	X17+	2005-09-07 17:17:00	2005-09-07 17:40:00	2005-09-07 18:03:00	808	X14.0	2001-04-15 14:05:00	2001-04-15 14:06:00	2001-04-16 13:00:00	9415
4	X14.4	2001-04-15 13:19:00	2001-04-15 13:50:00	2001-04-15 13:55:00	9415	X10.0	2003-10-29 20:55:00	2003-10-29 20:54:00	2003-10-30 00:00:00	10486
5	X10	2003-10-29 23:30:00	2003-10-29 23:40:00	2003-10-29 23:50:00	486	X9.4	1997-11-06 12:20:00	1997-11-06 12:10:00	1997-11-07 08:30:00	8100
6	X9.4	1997-11-06 11:49:00	1997-11-06 11:55:00	1997-11-06 12:01:00	8100	X9.3	2017-09-06 12:05:00	2017-09-06 12:24:00	2017-09-07 08:00:00	12673
7	X9.3	2017-09-06 11:55:00	2017-09-06 12:00:00	2017-09-06 12:10:00	2673	X9.0	2006-12-05 19:30:00	NaN	2006-12-05 20:00:00	10930
8	X9	2006-12-05 10:18:00	2006-12-05 10:55:00	2006-12-05 10:46:00	930	X8.3	2013-11-02 17:30:00	2013-11-02 17:30:00	2013-11-03 01:00:00	10486
9	X8.3	2013-11-02 17:03:00	2013-11-02 17:25:00	2013-11-02 17:30:00	486	X8.3	2017-09-10 16:02:00	2017-09-10 16:00:00	2017-09-11 06:50:00	NaN
10	X8.2	2017-09-10 15:35:00	2017-09-10 16:06:00	2017-09-10 16:31:00	2673	X7.1	2005-01-20 07:15:00	2005-01-20 06:54:00	2005-01-20 16:30:00	10720
11	X7.1	2005-01-20 06:36:00	2005-01-20 07:01:00	2005-01-20 07:26:00	720	X6.9	2011-08-09 08:20:00	2011-08-09 08:12:00	2011-08-09 08:35:00	11263
12	X6.9	2011-08-09 07:48:00	2011-08-09 08:05:00	2011-08-09 08:08:00	1263	X6.5	2006-12-06 19:00:00	NaN	2006-12-09 00:00:00	10930
13	X6.5	2006-12-06 18:29:00	2006-12-06 18:47:00	2006-12-06 19:00:00	930	X6.2	2005-09-09 19:45:00	2005-09-09 19:48:00	2005-09-09 22:00:00	10808
14	X6.2	2005-09-09 19:13:00	2005-09-09 19:20:00	2005-09-09 19:24:00	808	X5.7	2000-07-14 10:30:00	2000-07-14 10:54:00	2000-07-15 01:50:00	9077
15	X6.2	2001-12-13 14:20:00	2001-12-13 14:30:00	2001-12-13 14:35:00	9733	X5.6	2001-04-06 19:30:00	2001-04-06 19:30:00	2001-04-06 01:50:00	9415
16	X5.7	2000-07-14 10:03:00	2000-07-14 10:24:00	2000-07-14 10:44:00	9077	X5.4	2012-03-07 01:00:00	2012-03-07 00:24:00	2012-03-08 19:00:00	11429
17	X5.6	2001-04-06 19:10:00	2001-04-06 19:21:00	2001-04-06 19:31:00	9415	X5.3	2001-08-25 16:50:00	2001-08-25 16:50:00	2001-08-25 23:00:00	9591
18	X5.4	2012-03-07 00:02:00	2012-03-07 00:24:00	2012-03-07 00:40:00	1429	X4.9	2014-02-25 00:56:00	2014-02-25 01:25:00	2014-02-25 11:28:00	11990
19	X5.4	2005-08-09 20:52:00	2005-08-09 21:06:00	2005-08-09 21:17:00	486	X4.8	2002-07-23 00:50:00	2002-07-23 00:42:00	2002-07-23 04:00:00	9236
20	X5.4	2003-10-23 08:19:00	2003-10-23 08:35:00	2003-10-23 08:49:00	486	X4.0	2000-11-26 17:00:00	2000-11-26 17:06:00	2000-11-26 17:15:00	10339
21	X5.3	2001-08-25 16:23:00	2001-08-25 16:42:00	2001-08-25 16:50:00	9591	X3.9	2003-11-03 10:30:00	2003-11-03 10:56:00	2003-11-03 11:00:00	10488
22	X4.9	2014-02-25 00:38:00	2014-02-25 00:42:00	2014-02-25 01:03:00	1990	X3.8	2005-01-17 05:40:00	2005-01-17 05:45:00	2005-01-17 05:50:00	10720
23	X4.9	1998-08-18 22:10:00	1998-08-18 22:19:00	1998-08-18 22:28:00	8307	X3.6	2003-05-28 01:00:00	2003-05-28 00:50:00	2003-05-29 00:30:00	10365
24	X4.8	2002-07-23 00:18:00	2002-07-23 00:35:00	2002-07-23 00:47:00	39	X3.4	2006-12-13 02:45:00	2006-12-13 02:54:00	2006-12-13 10:40:00	10930
25	X4	2000-11-26 16:34:00	2000-11-26 16:48:00	2000-11-26 16:56:00	9236	X3.4	2001-12-28 20:35:00	2001-12-28 20:30:00	2001-12-29 03:00:00	9756
26	X3.9	2003-11-03 09:43:00	2003-11-03 09:55:00	2003-11-03 10:19:00	488	X3.3	2002-07-20 21:30:00	2002-07-20 22:06:00	2002-07-20 22:20:00	10039
27	X3.9	1998-08-19 21:35:00	1998-08-19 21:45:00	1998-08-19 21:50:00	8307	X3.2	2013-05-14 01:16:00	2013-05-14 01:25:00	2013-05-14 08:20:00	11748
28	X3.8	2005-01-17 06:59:00	2005-01-17 06:42:00	2005-01-17 06:49:00	720	X3.1	2002-08-24 01:27:00	2002-08-24 01:27:00	2002-08-24 02:58:00	10696
29	X3.7	1998-11-22 06:30:00	1998-11-22 06:42:00	1998-11-22 07:03:00	8384	X2.8	2013-05-13 16:15:00	2013-05-13 16:07:00	2013-05-13 16:15:00	11748
30	X3.6	2005-09-09 09:42:00	2005-09-09 09:59:00	2005-09-09 10:08:00	808	X2.7	2015-05-25 22:24:00	2015-05-25 22:24:00	2015-05-25 23:14:00	12339
31	X3.6	2004-07-16 13:49:00	2004-07-16 13:55:00	2004-07-16 14:01:00	649	X2.7	1998-05-06 08:25:00	1998-05-06 08:29:00	1998-05-06 08:35:00	8210
32	X3.6	2003-05-28 00:17:00	2003-05-28 00:27:00	2003-05-28 00:39:00	365	X2.7	2003-11-03 01:15:00	2003-11-03 01:59:00	2003-11-03 01:25:00	10488
33	X3.4	2006-12-13 06:20:00	2006-12-13 06:42:00	2006-12-13 07:03:00	930	X2.6	2005-01-15 16:15:00	2005-01-15 16:07:00	2005-01-17 05:50:00	10365







```
flare_location flare_region flare_classification cme_datetime \
274 N09W17 10996 2.0 2004-11-07 16:54:00

cme_angle cme_width cme_speed is_halo width_lower_bound
274 NaN 360 360 1759 True False
Current Flare Classification from SWH: 2.7
Closest Matches Found:
102 1.9
112 2.0
285 2.0
Name: flare_classification, dtype: float64
Current Start Datetime from SWH: 1998-05-06 07:58:00
Closest Start Datetime from NASA: 2000-11-25 19:30:00
Current Start Datetime from SWH: 1998-05-06 09:20:00
Closest Start Datetime from NASA: 2000-11-25 19:35:00
The Best Match:
start_datetime end_datetime start_frequency end_frequency \
102 2000-11-25 19:00:00 2000-11-25 19:35:00 6000 2000

flare_location flare_region flare_classification cme_datetime \
102 N20W23 9236 1.9 2000-11-25 19:31:00

cme_angle cme_width cme_speed is_halo width_lower_bound
102 NaN 360 671 True False
Current Flare Classification from SWH: 2.6
Closest Matches Found:
49 1.8
100 1.8
125 2.0
189 1.8
285 2.0
Name: flare_classification, dtype: float64
Current Start Datetime from SWH: 2005-01-17 09:25:00
Closest Start Datetime from NASA: 2005-01-17 09:25:00
Current Start Datetime from SWH: 2005-01-17 23:31:00
Closest Start Datetime from NASA: 2005-01-17 16:00:00
The Best Match:
start_datetime end_datetime start_frequency end_frequency \
285 2005-01-17 09:25:00 2005-01-17 16:00:00 14000 30

flare_location flare_region flare_classification cme_datetime \
285 N15W25 10720 2.0 2004-11-07 09:30:00

cme_angle cme_width cme_speed is_halo width_lower_bound
285 NaN 360 2094 True False
Current Flare Classification from SWH: 2.6
Closest Matches Found:
49 1.8
100 1.8
125 2.0
189 1.8
Name: flare_classification, dtype: float64
Current Start Datetime from SWH: 2001-09-24 09:32:00
Closest Start Datetime from NASA: 2001-04-12 10:20:00
Current Start Datetime from SWH: 2001-09-24 11:09:00
Closest Start Datetime from NASA: 2001-04-12 10:40:00
The Best Match:
start_datetime end_datetime start_frequency end_frequency \
125 2001-04-12 10:20:00 2001-04-12 10:40:00 14000 7000

flare_location flare_region flare_classification cme_datetime \
125 S19W43 9415 2.0 2001-04-12 10:31:00

cme_angle cme_width cme_speed is_halo width_lower_bound
125 NaN 360 1184 True False
```

```
In [29]: #Final dataset/output for Task 6 (with SWRank as the last column):
out_NASA
```

	start_datetime	end_datetime	start_frequency	end_frequency	flare_location	flare_region	flare_classification	cme_datetime	cme_angle	cme_width
0	1997-04-01 14:00:00	1997-04-01 14:15:00	8000	4000	S25E16	8026	M1.3	1997-04-01 15:18:00	74	NaN
1	1997-04-07 14:30:00	1997-04-07 17:30:00	11000	1000	S28E19	8027	C6.8	1997-04-07 14:27:00	NaN	NaN
2	1997-05-12 05:15:00	1997-05-14 16:00:00	12000	80	N21W08	8038	C1.3	1997-05-12 05:30:00	NaN	NaN
3	1997-05-21 20:20:00	1997-05-21 22:00:00	5000	500	N05W12	8040	M1.3	1997-05-21 21:00:00	263	NaN
4	1997-09-23 21:53:00	1997-09-23 22:16:00	6000	2000	S29E25	8088	C1.4	1997-09-23 22:02:00	133	NaN
...	...	...	...	...	...	...	...	...	...	...
517	2017-09-17 11:45:00	2017-09-17 12:35:00	16000	900	S08E170	NaN	NaN	2017-09-17 12:00:00	NaN	NaN
518	2017-10-18 05:48:00	2017-10-18 12:40:00	16000	400	S06E123	NaN	NaN	2017-10-18 08:00:00	85	NaN
519	2019-05-03 23:52:00	2019-05-04 00:16:00	13000	2300	N12E82	12740	C1.0	2019-05-03 23:24:00	90	NaN
520	2020-11-29 13:07:00	2020-11-29 15:23:00	14000	850	S23E89	NaN	M4.4	2020-11-29 13:25:00	NaN	NaN
521	2020-12-07 16:18:00	2020-12-08 02:00:00	14000	160	S25W08	12790	C7.4	2020-12-07 16:24:00	NaN	NaN

522 rows × 14 columns

```
In [30]: #To more clearly see the matches, I am outputting a shortened out_NASA dataset with just the matching row
import math

out_NASA_matches = out_NASA.loc[out_NASA['SWRank'] >= 0]
#We notice that the 502000th row values are floats in out_NASA, which needs to be the case because NaN values
#must be floats. Since this dataset only have non-null values, we can show the ranks as ints.
out_NASA_matches = out_NASA_matches.astype({'SWRank': 'int32'})
print("The size of this dataset is 50, as it should be: ", out_NASA_matches.shape)
out_NASA_matches
```

The size of this dataset is 50, as it should be: (50, 14)

	start_datetime	end_datetime	start_frequency	end_frequency	flare_location	flare_region	flare_classification	cme_datetime	cme_angle	cme_width
7	1997-11-04 14:00:00	1997-11-05 14:15:00	14000	100	S14W33	8100	X2.1	1997-11-04 06:18:00	NaN	NaN
8	1997-11-06 08:25:00	1997-11-07 02:30:00	14000	100	S18W63	8100	X9.4	1997-11-06 03:24:00	NaN	NaN
9	1997-11-27 13:30:00	1997-11-27 14:05:00	14000	7000	N17E63	8113	X2.6	1997-11-27 13:56:00	98	NaN
19	1998-05-06 08:25:00	1998-05-06 08:35:00	14000	5000	S11W65	8210	X2.7	1998-05-06 08:29:00	309	NaN
73	2000-06-06 15:20:00	2000-06-08 09:00:00	14000	40	N20E18	9026	X2.3	2000-06-06 15:54:00	NaN	NaN
82	2000-07-14 10:30:00	2000-07-15 14:30:00	14000	80	N22W07	9077	X5.7	2000-07-14 10:54:00	NaN	NaN
98	2000-11-24 05:00:00	2000-11-24 15:00:00	14000	100	N20W05	9236	X2.0	2000-11-24 05:30:00	NaN	NaN
99	2000-11-24 15:25:00	2000-11-24 22:00:00	14000	200	N22W07	9236	X2.3	2000-11-24 15:30:00	NaN	NaN
102	2000-11-25 19:00:00	2000-11-25 19:35:00	6000	2000	N20W23	9236	X1.9	2000-11-25 19:31:00	NaN	NaN
104	2000-11-26 17:00:00	2000-11-26 17:15:00	14000	7000	N18W38	9236	X4.0	2000-11-26 17:06:00	NaN	NaN
117	2001-04-02 22:05:00	2001-04-03 02:30:00	14000	250	N19W72	9293	X3.0	2001-04-02 22:06:00	261	NaN
121	2001-04-06 19:35:00	2001-04-07 01:50:00	14000	230	S21E31	9415	X5.6	2001-04-06 19:30:00	NaN	NaN
123	2001-04-10 05:24:00	2001-04-11 00:00:00	14000	100	S23W09	9415	X2.3	2001-04-10 05:30:00	NaN	NaN
125	2001-04-12 10:20:00	2001-04-12 10:40:00	14000	7000	S19W43	9415	X2.0	2001-04-12 10:31:00	NaN	NaN
126	2001-04-15 14:05:00	2001-04-16 13:50:00	14000	40	S20W85	9415	X1.4	2001-04-15 14:06:00	245	NaN
135	2001-08-25 16:50:00	2001-08-25 23:00:00	8000	170	S17E34	9591	X5.3	2001-08-25 16:50:00	NaN	NaN
142	2001-09-24 10:45:00	2001-09-25 20:00:00	7000	30	S16E23	9632	X2.6	2001-09-24 10:30:00	NaN	NaN
160	2001-12-28 20:35:00	2001-12-29 03:00:00	14000	350	S26E90	9756	X3.4	2001-12-28 20:30:00	NaN	NaN
192	2002-07-23 05:00:00	2002-07-23 04:00:00	10000	2000	S13E90	10039	X3.3	2002-07-23 04:20:00	NaN	NaN
193	2002-08-24 01:45:00	2002-08-24 03:25:00	5000	400	S02W81	10069	X3.1	2002-08-24 01:27:00	NaN	NaN
222	2003-05-28 01:00:00	2003-05-29 00:30:00	1000	200	S07W20	10365	X1.7	2003-05-28 00:50:00	NaN	NaN
233	2003-10-28 11:00:00	2003-10-30 00:00:00	14000	40	S16E08	10486	X3.6	2003-10-28 11:30:00	NaN	NaN
234	2003-10-29 20:55:00	2003-10-30 00:00:00	11000	500	S15W02	10486	X1.0	2003-10-29 20:54:00	NaN	NaN
237	2003-11-02 01:50:00	2003-11-03 01:25:00	12000	250	S14W56	10486	X8.3	2003-11-02 01:50:00	NaN	NaN
238	2003-11-03 01:50:00	2003-11-03 01:25:00	3000	1500	N10W83	10488	X2.7	2003-11-03 01:50:00	304	NaN
239	2003-11-03 10:00:00	2003-11-03 12:30:00	6000	400	N08W77	10488	X3.9	2003-11-03 10:06:00	293	NaN
240	2003-11-04 20:00:00	2003-11-05 00:00:00	10000	200	S19W83	10486	X2.8	2003-11-04 19:54:00	NaN	NaN
274	2004-11-07 16:25:00	2004-11-08 20:00:00	14000	60	N09W17	10696	X2.0	2004-11-07 16:54:00	NaN	NaN
276	2004-11-10 02:25:00	2004-11-10 05:40:00	14000	1000	N09W49	10696	X2.5	2004-11-10 02:26:00	NaN	NaN
284	2005-01-15 23:00:00	2005-01-17 16:00:00	3000	40	N15W05	10720	X2.6	2005-01-15 23:06:00	NaN	NaN
285	2005-01-17 09:25:00	2005-01-17 16:00:00	14000	30	N15W25	10720	X2.0	2005-01-17 09:25:00	NaN	NaN
286	2005-01-17 10:00:00	2005-01-17 10:35:00	6100	1500	N15W25	10720	X3.8	2005-01-17 09:54:00	NaN	NaN
288	2005-01-20 07:15:00	2005-01-20 16:30:00	14000	25	N14W61	10720	X7.1	2005-01-20 06:54:00	NaN	NaN
317	2005-09-09 19:45:00	2005-09-09 22:00:00	10000	50	S12E67	10808	X6.2	2005-09-09 19:48:00	NaN	NaN
318	2005-09-10 21:45:00	2005-09-10 01:00:00	14000	200	S13E47	10808	X2.1	2005-09-10 21:52:00	NaN	NaN
328	2006-12-05 16:00:00	2006-12-05 20:00:00	14000	250	S07E68	10930	X9.0	NaT	NaN	NaN
331	2006-12-06 01:50:00	2006-12-09 00:00:00	16000	30	S05E64	10930	X6.5	NaT	NaN	NaN
332	2006-12-13 02:45:00	2006-12-13 10:40:00	12000	150	S06W23	10930	X3.4	2006-12-13 02:54:00	NaN	NaN
345	2011-02-15 02:10:00	2011-02-15 07:00:00	16000	400	S20W12	11158	X2.2	2011-02-15 02:24:00	NaN	NaN
359	2011-08-09 08:20:00	2011-08-09 08:35:00	16000	4000	N17W69	11263	X6.9	2011-08-09 08:20:00	NaN	NaN
361	2011-09-06 22:30:00	2011-09-07 15:40:00	16000	150	N14W18	11283	X2.1	2011-09-06 23:05:00	NaN	NaN
375	2012-03-07 01:00:00	2012-03-08 19:00:00	16000	30	N17E27	11429	X5.4	2012-03-07 01:20:00	NaN	NaN
403	2013-05-13 16:50:00	2013-05-13 19:10:00	16000	300	N11E85	11748	X2.8	2013-05-13 16:07:00	NaN	NaN
404	2013-05-14 01:16:00	2013-05-14 08:20:00	16000	240	N08E77	11748	X3.2	2013-05-14 01:25:00	NaN	NaN
420	2013-10-25 15:08:00	2013-10-25 22:32:00	16000	200	S06E69	11882	X4.1	2013-10-25 15:12:00	NaN	NaN
443	2014-02-25 00:56:00	2014-02-25 11:28:00	14000	100	S12E82	11990	X2.9	2014-02-25 00:54:00	NaN	NaN
487	2015-05-05 22:24:00	2015-05-05 23:14:00	14000	500	N15E79	12339	X2.7	2015-05-05 22:24:00	NaN	NaN
514	2017-09-06 12:25:00	2017-09-07 08:00:00	16000	70	S08W33	12673	X9.3	2017-09-06 12:00:00	NaN	NaN
515	2017-09-10 16:00:00	2017-09-11 06:50:00	16000	150	S09W92	NaN	X8.3	2017-09-10 16:00:00	NaN	NaN

Task 7: Attributes Visualization:

```
In [31]: <matplotlib inline
```

```
In [32]: #Sorting the top 50 flares from NASA dataset (top50)
topFlareIdx = tableNASAIdxCopy.loc[tableNASAIdxCopy['flare_classification'].str.contains('X')]
topFlareIdx['flare_classification'] = topFlareIdx['flare_classification'].str.strip('X').astype(float)
topFlareIdx = topFlareIdx.sort_values('flare_classification', ascending = False)
topFlareIdx = topFlareIdx.head(50)
topFlareIdx['flare_classification'] = topFlareIdx['flare_classification'].astype(str)
topFlareIdx['flare_classification'] = "X" + topFlareIdx['flare_classification']

!python-input-32-c7f73ff2bfb8d3: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame.
Try using .loc[row_index,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
topFlareIdx['flare_classification'] = topFlareIdx['flare_classification'].str.strip('X').astype(float)
```

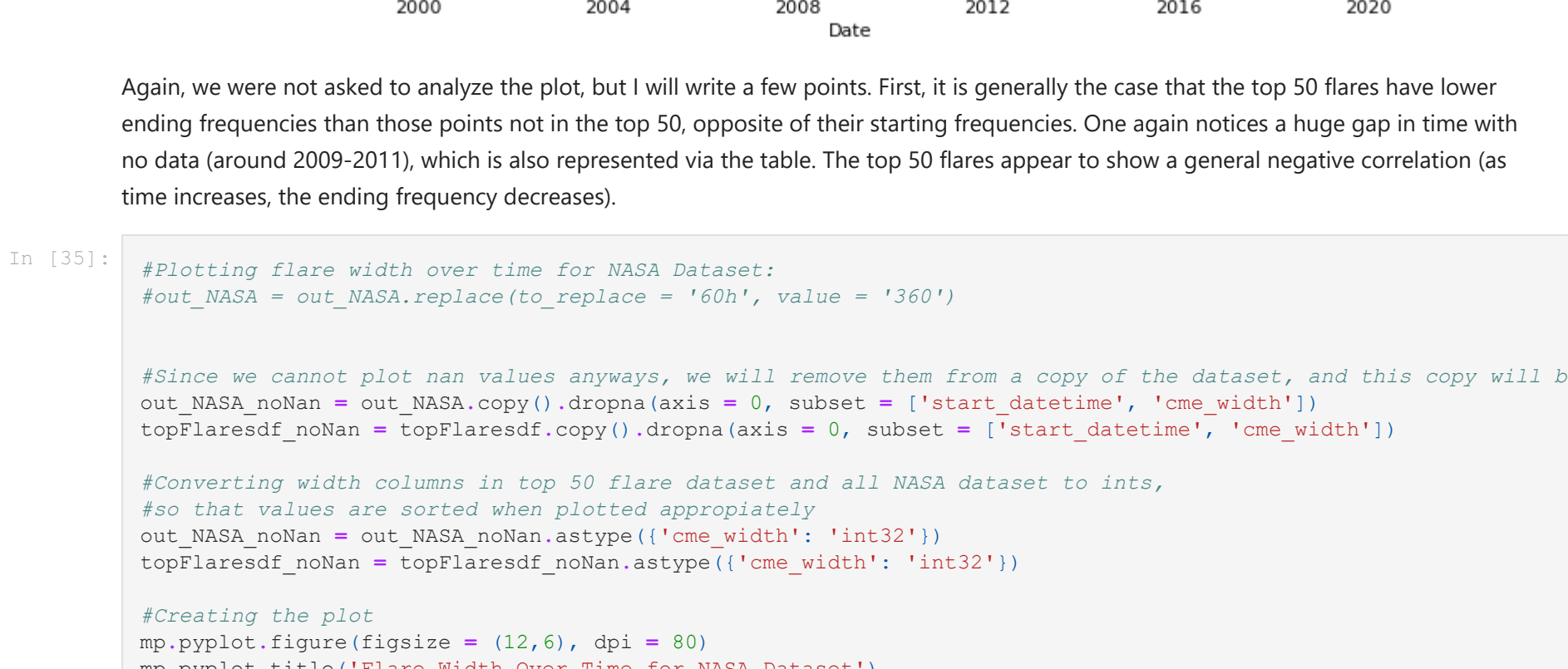
```
In [33]: #Plotting start frequency over time for NASA Dataset:

#Since we cannot plot nan values anyway, we will remove them from a copy of the dataset, and this copy will be
out_NASA_noNaN = out_NASA.copy().dropna(axis = 0, subset = ['start_datetime', 'start_frequency'])

#Converting frequency columns in top 50 flare dataset and all NASA dataset to int, so that values are sorted
#so that values are sorted when plotted appropriately
out_NASA_noNaN = out_NASA_noNaN.astype({'start_datetime': 'int32'})
topFlareIdx = topFlareIdx.astype({'start_datetime': 'int32'})

#Creating the plot
mp.pyplot.figure(figsize = (12,6), dpi = 80)
mp.pyplot.title('Starting Frequency Over Time for NASA Dataset')
mp.pyplot.xlabel('Date')
mp.pyplot.ylabel('Starting Frequency')
mp.pyplot.plot(out_NASA_noNaN['start_datetime'], out_NASA_noNaN['start_frequency'], 'b', label = 'All NASA data')
mp.pyplot.plot(topFlareIdx['start_datetime'], topFlareIdx['start_frequency'], 'g', label = 'Top 50')
mp.pyplot.legend(loc = 'upper right')
```

Out [33]: <matplotlib.legend.Legend at 0x2362c421640>



We were not asked to analyze the plot, but I will write a few points. First, it is generally the case that the top 50 flares have higher starting frequencies than those points not in the top 50. They also seem to be more clustered towards time before 2008. One also notices a huge gap in time with no data (around 2009-2011), which is also represented by the table. The top 50 flares appear to show a general negative correlation as time increases, the ending frequency decreases.

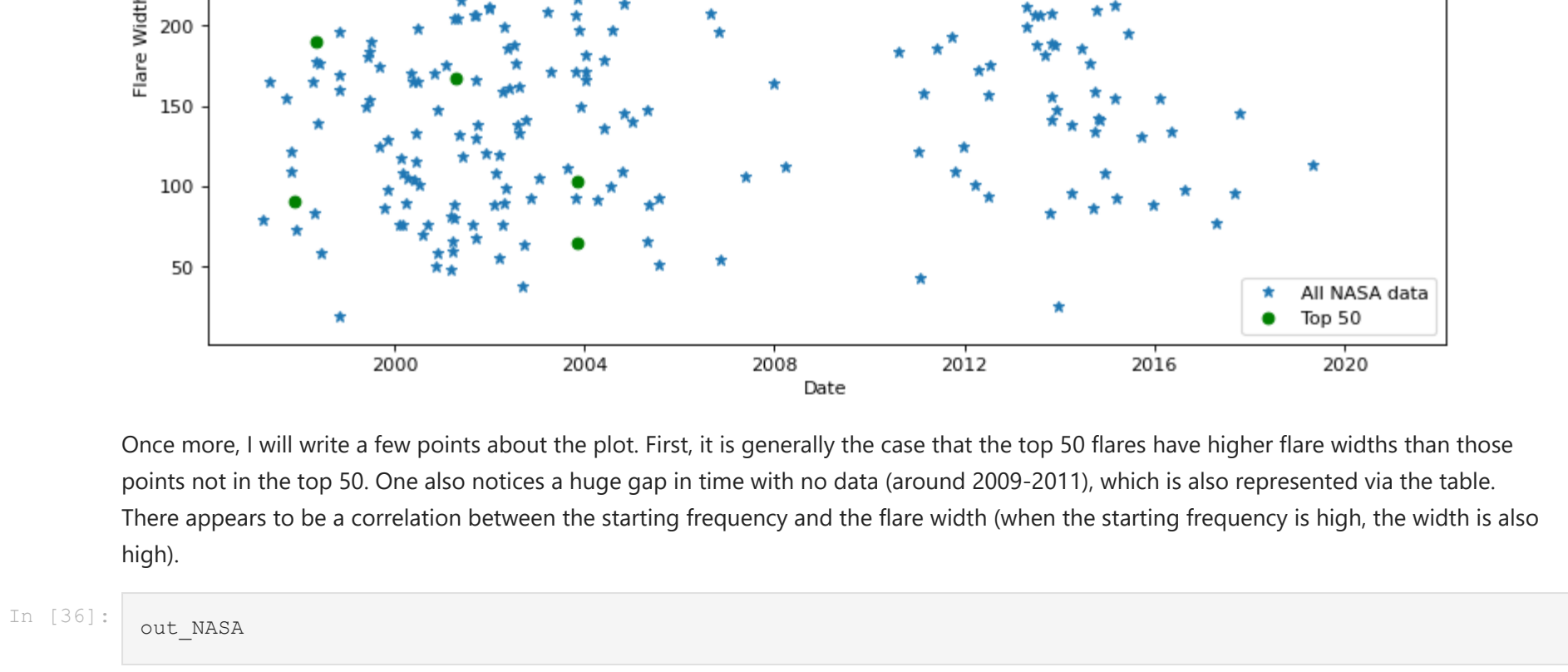
```
In [34]: #Plotting end frequency over time for NASA Dataset:

#Since we cannot plot nan values anyway, we will remove them from a copy of the dataset, and this copy will be
out_NASA_noNaN = out_NASA.copy().dropna(axis = 0, subset = ['start_datetime', 'end_frequency'])

#Converting frequency columns in top 50 flare dataset and all NASA dataset to int,
#so that values are sorted when plotted appropriately
out_NASA_noNaN = out_NASA_noNaN.astype({'end_frequency': 'int32'})
topFlareIdx = topFlareIdx.astype({'end_frequency': 'int32'})

#Creating the plot
mp.pyplot.figure(figsize = (12,6), dpi = 80)
mp.pyplot.title('Ending Frequency Over Time for NASA Dataset')
mp.pyplot.xlabel('Date')
mp.pyplot.ylabel('Ending Frequency')
mp.pyplot.plot(out_NASA_noNaN['start_datetime'], out_NASA_noNaN['end_frequency'], 'b', label = 'All NASA data')
mp.pyplot.plot(topFlareIdx['start_datetime'], topFlareIdx['end_frequency'], 'g', label = 'Top 50')
mp.pyplot.legend(loc = 'upper right')
```

Out [34]: <matplotlib.legend.Legend at 0x2362c49f6d0>



Again, we were not asked to analyze the plot, but I will write a few points. First, it is generally the case that the top 50 flares have lower ending frequencies than those points not in the top 50, opposite of their starting frequencies. One again notices a huge gap in time with no data (around 2009-2011), which is also represented by the table. The top 50 flares appear to show a general negative correlation as time increases, the ending frequency decreases.

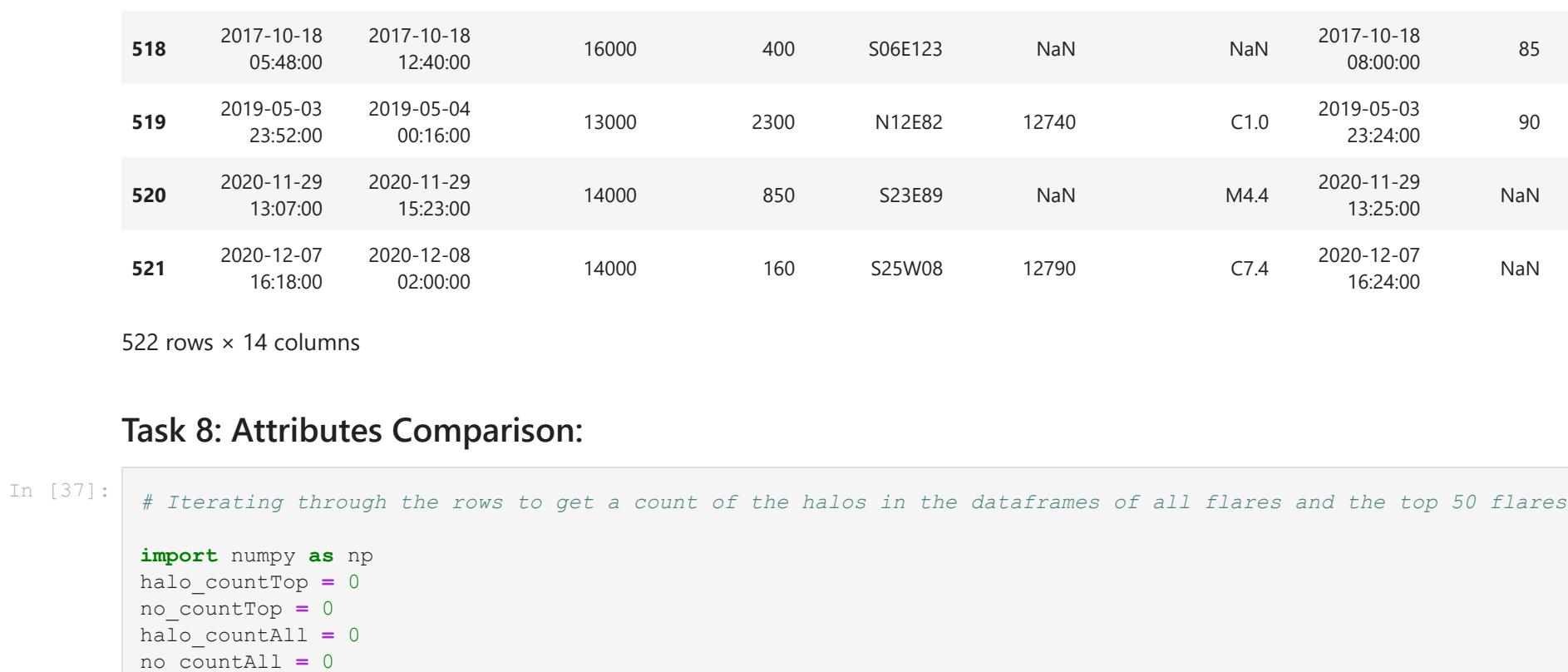
```
In [35]: #Plotting flare width over time for NASA Dataset:
four_NASA = out_NASA.replace(to_replace = '60h', value = '360')

#Since we cannot plot nan values anyway, we will remove them from a copy of the dataset, and this copy will be
out_NASA_noNaN = out_NASA.copy().dropna(axis = 0, subset = ['start_datetime', 'cme_width'])
topFlareIdx_noNaN = topFlareIdx.copy().dropna(axis = 0, subset = ['start_datetime', 'cme_width'])

#Converting width columns in top 50 flare dataset and all NASA dataset to int,
#so that values are sorted when plotted appropriately
out_NASA_noNaN = out_NASA_noNaN.astype({'cme_width': 'int32'})
topFlareIdx_noNaN = topFlareIdx_noNaN.astype({'cme_width': 'int32'})

#Creating the plot
mp.pyplot.figure(figsize = (12,6), dpi = 80)
mp.pyplot.title('Flare Width Over Time for NASA Dataset')
mp.pyplot.xlabel('Date')
mp.pyplot.ylabel('Flare Width')
mp.pyplot.plot(out_NASA_noNaN['start_datetime'], out_NASA_noNaN['cme_width'], 'b', label = 'All NASA data')
mp.pyplot.plot(topFlareIdx_noNaN['start_datetime'], topFlareIdx_noNaN['cme_width'], 'g', label = 'Top 50')
mp.pyplot.legend(loc = 'lower right')
```

Out [35]: <matplotlib.legend.Legend at 0x2362c49f6d0>



Once again, I will write a few points about the plot. First, it is generally the case that the top 50 flares have higher flare widths than those points not in the top 50. One also notices a huge gap in time with no data (around 2009-2011), which is also represented by the table. There appears to be a correlation between the starting frequency and the flare width (when the starting frequency is high, the width is also high).

```
In [36]: out_NASA</
```