



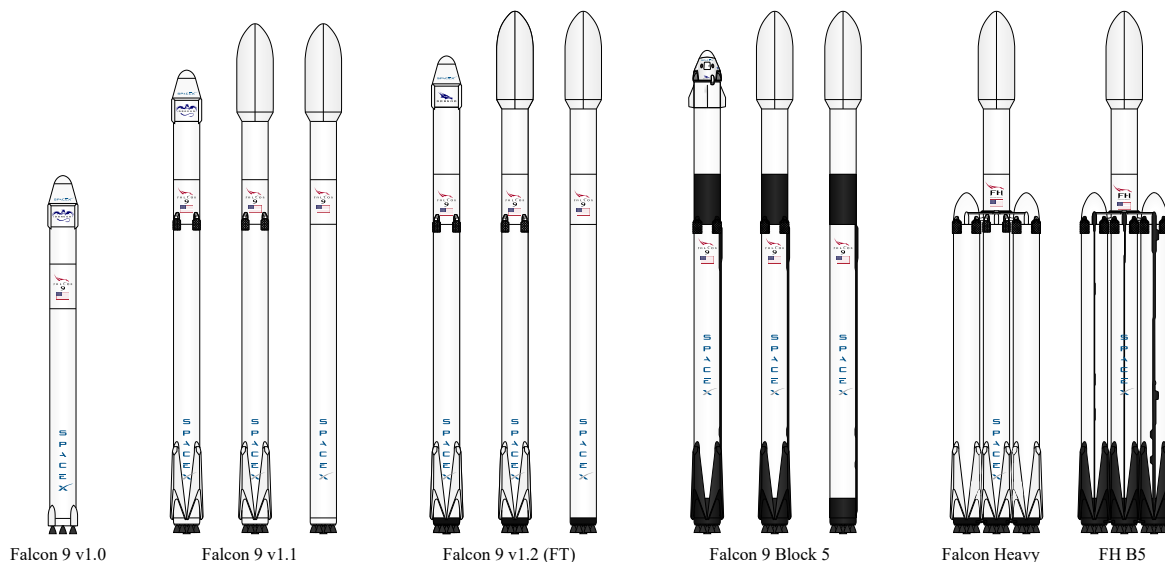
## ✓ Space X Falcon 9 First Stage Landing Prediction

### ✓ Web scraping Falcon 9 and Falcon Heavy Launches Records from Wikipedia

Estimated time needed: **40** minutes

In this lab, you will be performing web scraping to collect Falcon 9 historical launch records from a Wikipedia page titled [List of Falcon 9 and Falcon Heavy launches](https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)

[https://en.wikipedia.org/wiki/List\\_of\\_Falcon\\_9\\_and\\_Falcon\\_Heavy\\_launches](https://en.wikipedia.org/wiki/List_of_Falcon_9_and_Falcon_Heavy_launches)



Falcon 9 first stage will land successfully



Several examples of an unsuccessful landing are shown here:



More specifically, the launch records are stored in a HTML table shown below:

2020 [edit]

In late 2019, [Gwynne Shotwell](#) stated that SpaceX hoped for as many as 24 launches for Starlink satellites in 2020,<sup>[490]</sup> in addition to 14 or 15 non-Starlink launches. At 26 launches, 13 of which for Starlink satellites, Falcon 9 had its most prolific year, and Falcon rockets were second most prolific rocket family of 2020, only behind China's [Long March](#) rocket family.<sup>[491]</sup>

[hide] Flight No.	Date and time (UTC)	Version, Booster <sup>[b]</sup>	Launch site	Payload <sup>[c]</sup>	Payload mass	Orbit	Customer	Launch outcome	Booster landing
78	7 January 2020, 02:19:21 <sup>[492]</sup>	F9 B5 <span>△</span> B1049.4	CCAFS, SLC-40	Starlink 2 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[8]</sup>	LEO	SpaceX	Success	Success (drone ship)
Third large batch and second operational flight of Starlink constellation. One of the 60 satellites included a test coating to make the satellite less reflective, and thus less likely to interfere with ground-based astronomical observations. <sup>[493]</sup>									
79	19 January 2020, 15:30 <sup>[494]</sup>	F9 B5 <span>△</span> B1046.4	KSC, LC-39A	Crew Dragon in-flight abort test <sup>[495]</sup> (Dragon C205.1)	12,050 kg (26,570 lb)	Sub-orbital <sup>[496]</sup>	NASA (CTS) <sup>[497]</sup>	Success	No attempt
An atmospheric test of the <i>Dragon 2</i> abort system after <i>Max Q</i> . The capsule fired its <i>SuperDraco</i> engines, reached an apogee of 40 km (25 mi), deployed parachutes after reentry, and <b>splashed down</b> in the ocean 31 km (19 mi) downrange from the launch site. The test was previously slated to be accomplished with the <i>Crew Dragon Demo-1</i> capsule; <sup>[498]</sup> but that test article exploded during a ground test of SuperDraco engines on 20 April 2019. <sup>[419]</sup> The abort test used the capsule originally intended for the first crewed flight. <sup>[499]</sup> As expected, the booster was destroyed by aerodynamic forces after the capsule aborted. <sup>[500]</sup> First flight of a Falcon 9 with only one functional stage — the second stage had a <i>mass simulator</i> in place of its engine.									
80	29 January 2020, 14:07 <sup>[501]</sup>	F9 B5 <span>△</span> B1051.3	CCAFS, SLC-40	Starlink 3 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[8]</sup>	LEO	SpaceX	Success	Success (drone ship)
Third operational and fourth large batch of Starlink satellites, deployed in a circular 290 km (180 mi) orbit. One of the fairing halves was caught, while the other was fished out of the ocean. <sup>[502]</sup>									
81	17 February 2020, 15:05 <sup>[503]</sup>	F9 B5 <span>△</span> B1056.4	CCAFS, SLC-40	Starlink 4 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[8]</sup>	LEO	SpaceX	Success	Failure (drone ship)
Fourth operational and fifth large batch of Starlink satellites. Used a new flight profile which deployed into a 212 km × 386 km (132 mi × 240 mi) elliptical orbit instead of launching into a circular orbit and firing the second stage engine twice. The first stage booster failed to land on the drone ship <sup>[504]</sup> due to incorrect wind data. <sup>[505]</sup> This was the first time a flight proven booster failed to land.									
82	7 March 2020, 04:50 <sup>[506]</sup>	F9 B5 <span>△</span> B1059.2	CCAFS, SLC-40	SpaceX CRS-20 (Dragon C112.3 <span>△</span> )	1,977 kg (4,359 lb) <sup>[507]</sup>	LEO (ISS)	NASA (CRS)	Success	Success (ground pad)
Last launch of phase 1 of the CRS contract. Carries <i>Bartolomeo</i> , an <i>ESA</i> platform for hosting external payloads onto ISS. <sup>[508]</sup> Originally scheduled to launch on 2 March 2020, the launch date was pushed back due to a second stage engine failure. SpaceX decided to swan out the second stage instead of replacing the faulty part. <sup>[509]</sup> It was SpaceX's 60th successful landing of a first stage booster, the third flight of the <i>Dragon</i> C112 and the last launch of the <i>Dragon</i> spacecraft.									

83	18 March 2020, 12:16 <sup>[510]</sup>	F9 B5 Δ B1048.5	KSC, LC-39A	Starlink 5 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[5]</sup>	LEO	SpaceX	Success	Failure (drone ship)
Fifth operational launch of Starlink satellites. It was the first time a first stage booster flew for a fifth time and the second time the fairings were reused (Starlink flight in May 2019). <sup>[511]</sup> Towards the end of the first stage burn, the booster suffered premature shut down of an engine, the first of a <i>Merlin 1D</i> variant and first since the CRS-1 mission in October 2012. However, the payload still reached the targeted orbit. <sup>[512]</sup> This was the second Starlink launch booster landing failure in a row, later revealed to be caused by residual cleaning fluid trapped inside a sensor. <sup>[513]</sup>									
84	22 April 2020, 19:30 <sup>[514]</sup>	F9 B5 Δ B1051.4	KSC, LC-39A	Starlink 6 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[5]</sup>	LEO	SpaceX	Success	Success (drone ship)

## ▼ Objectives

Web scrap Falcon 9 launch records with BeautifulSoup :

- Extract a Falcon 9 launch records HTML table from Wikipedia
- Parse the table and convert it into a Pandas data frame

First let's import required packages for this lab

```
1 !pip3 install beautifulsoup4
2 !pip3 install requests
```

```
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.10/dist-packa
Requirement already satisfied: soupsieve>1.2 in /usr/local/lib/python3.10/dist-packag
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (2
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-package
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-p
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-p
```

Import Libraries

```
1 import sys
2
3 import requests
4 from bs4 import BeautifulSoup
5 import re
6 import unicodedata
7 import pandas as pd
```

and we will provide some helper functions for you to process web scraped HTML table

```
1 def date_time(table_cells):
2     """
3     This function returns the data and time from the HTML table cell
4     Input: the element of a table data cell extracts extra row
5     """
6     return [data_time.strip() for data_time in list(table_cells.strings)][0:2]
7
8 def booster_version(table_cells):
```

```

8 def booster_version(table_cells):
9     """
10    This function returns the booster version from the HTML table cell
11    Input: the element of a table data cell extracts extra row
12    """
13    out=''.join([booster_version for i,booster_version in enumerate( table_cells.strin
14    return out
15
16 def landing_status(table_cells):
17     """
18    This function returns the landing status from the HTML table cell
19    Input: the element of a table data cell extracts extra row
20    """
21    out=[i for i in table_cells.strings][0]
22    return out
23
24
25 def get_mass(table_cells):
26    mass=unicodedata.normalize("NFKD", table_cells.text).strip()
27    if mass:
28        mass.find("kg")
29        new_mass=mass[0:mass.find("kg")+2]
30    else:
31        new_mass=0
32    return new_mass
33
34
35 def extract_column_from_header(row):
36     """
37    This function returns the landing status from the HTML table cell
38    Input: the element of a table data cell extracts extra row
39    """
40    if (row.br):
41        row.br.extract()
42    if row.a:
43        row.a.extract()
44    if row.sup:
45        row.sup.extract()
46
47    column_name = ' '.join(row.contents)
48
49    # Filter the digit and empty names
50    if not(column_name.strip().isdigit()):
51        column_name = column_name.strip()
52        return column_name
53

```

To keep the lab tasks consistent, you will be asked to scrape the data from a snapshot of the List of Falcon 9 and Falcon Heavy launches Wikipage updated on 9th June 2021

```
1 static_url = "https://en.wikipedia.org/w/index.php?title=List_of_Falcon_9_and_Falcon_H
```

Next, request the HTML page from the above URL and get a `response` object

## ✓ TASK 1: Request the Falcon9 Launch Wiki page from its URL

Perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

```
1 # use requests.get() method with the provided static_url
2 # assign the response to a object
3 response = requests.get(static_url)
```

Create a `BeautifulSoup` object from the HTML response

```
1 # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
2 soup = BeautifulSoup(response.content, 'html.parser')
```

Print the page title to verify if the `BeautifulSoup` object was created properly

```
1 # Use soup.title attribute
2 soup.title
<title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

## ✓ TASK 2: Extract all column/variable names from the HTML table header

Next, we want to collect all relevant column names from the HTML table header

Let's try to find all tables on the wiki page first. If you need to refresh your memory about `BeautifulSoup`, please check the external reference link towards the end of this lab

```
1 # Use the find_all function in the BeautifulSoup object, with element type `table`
2 # Assign the result to a list called `html_tables`
3 html_tables = soup.find_all('table')
```

Starting from the third table is our target table contains the actual launch records.

```
1 # Let's print the third table and check its content
2 first_launch_table = html_tables[2]
```

```
3 print(first_launch_table)
```

```
<table class="wikitable plainrowheaders collapsible" style="width: 100%;">
<tbody><tr>
<th scope="col">Flight No.
</th>
<th scope="col">Date and<br/>time (<a href="/wiki/Coordinated_Universal_Time" title="
</th>
<th scope="col"><a href="/wiki/List_of_Falcon_9_first-stage_boosters" title="List of
</th>
<th scope="col">Launch site
</th>
<th scope="col">Payload<sup class="reference" id="cite_ref-Dragon_12-0"><a href="#cit
</th>
<th scope="col">Payload mass
</th>
<th scope="col">Orbit
</th>
<th scope="col">Customer
</th>
<th scope="col">Launch<br/>outcome
</th>
<th scope="col"><a href="/wiki/Falcon_9_first-stage_landing_tests" title="Falcon 9 fi
</th></tr>
<tr>
<th rowspan="2" scope="row" style="text-align:center;">1
</th>
<td>4 June 2010,<br/>18:45
</td>
<td><a href="/wiki/Falcon_9_v1.0" title="Falcon 9 v1.0">F9 v1.0</a><sup class="refere
</td>
<td><a href="/wiki/Cape_Canaveral_Space_Force_Station" title="Cape Canaveral Space Fo
</td>
<td><a href="/wiki/Dragon_Spacecraft_Qualification_Unit" title="Dragon Spacecraft Qua
</td>
<td>
</td>
<td><a href="/wiki/Low_Earth_orbit" title="Low Earth orbit">LEO</a>
</td>
<td><a href="/wiki/SpaceX" title="SpaceX">SpaceX</a>
</td>
<td class="table-success" style="background: #9EFF9E; vertical-align: middle; text-al
</td>
<td class="table-failure" style="background: #FFC7C7; vertical-align: middle; text-al
</td></tr>
<tr>
<td colspan="9">First flight of Falcon 9 v1.0.<sup class="reference" id="cite_ref-sfn
</td></tr>
<tr>
<th rowspan="2" scope="row" style="text-align:center;">2
</th>
<td>8 December 2010,<br/>15:43<sup class="reference" id="cite_ref-spaceflightnow_Clar
</td>
<td><a href="/wiki/Falcon_9_v1.0" title="Falcon 9 v1.0">F9 v1.0</a><sup class="refere
</td>
<td><a href="/wiki/Cape Canaveral Space Force Station" title="Cape Canaveral Space Fo
```

```

</td>
<td><a href="/wiki/SpaceX_Dragon" title="SpaceX Dragon">Dragon</a> <a class="mw-redir
</td>
<td>

```

You should be able to see the column names embedded in the table header elements `<th>` as follows:

```

<tr>
<th scope="col">Flight No.
</th>
<th scope="col">Date and<br/>time (<a href="/wiki/Coordinated_Universal_Time" title="Coordinated U
</th>
<th scope="col"><a href="/wiki/List_of_Falcon_9_first-stage_boosters" title="List of Falcon 9 first
</th>
<th scope="col">Launch site
</th>
<th scope="col">Payload<sup class="reference" id="cite_ref-Dragon_12-0"><a href="#cite_note-Dragon
</th>
<th scope="col">Payload mass
</th>
<th scope="col">Orbit
</th>
<th scope="col">Customer
</th>
<th scope="col">Launch<br/>outcome
</th>
<th scope="col"><a href="/wiki/Falcon_9_first-stage_landing_tests" title="Falcon 9 first-stage lan
</th></tr>

```

Next, we just need to iterate through the `<th>` elements and apply the provided `extract_column_from_header()` to extract column name one by one

```

1 column_names = []
2
3 # Apply find_all() function with `th` element on first_launch_table
4 # Iterate each th element and apply the provided extract_column_from_header() to get a
5 # Append the Non-empty column name (`if name is not None and len(name) > 0`) into a li
6 table = first_launch_table.find_all('th')
7 for row in table:
8     name = extract_column_from_header(row)
9     if name is not None and len(name) > 0:

```

```

7         if name is not None and len(name) > 0:
10             column_names.append(name)

```

Check the extracted column names

```

1 print(column_names)
    ['Flight No.', 'Date and time ( )', 'Launch site', 'Payload', 'Payload mass', 'Orbit'

```

## ✓ TASK 3: Create a data frame by parsing the launch HTML tables

Create an empty dictionary with keys from the extracted column names in the previous task. Later, this dictionary will be converted into a Pandas dataframe

```

1 launch_dict= dict.fromkeys(column_names)
2
3 # Remove an irrelevant column
4 del launch_dict['Date and time ( )']
5
6 # Let's initial the launch_dict with each value to be an empty list
7 launch_dict['Flight No.'] = []
8 launch_dict['Launch site'] = []
9 launch_dict['Payload'] = []
10 launch_dict['Payload mass'] = []
11 launch_dict['Orbit'] = []
12 launch_dict['Customer'] = []
13 launch_dict['Launch outcome'] = []
14 # Added some new columns
15 launch_dict['Version Booster']=[]
16 launch_dict['Booster landing']=[]
17 launch_dict['Date']=[]
18 launch_dict['Time']=[]

```

Next, we just need to fill up the `launch_dict` with launch records extracted from table rows.

Usually, HTML tables in Wiki pages are likely to contain unexpected annotations and other types of noises, such as reference links `B0004.1[8]`, missing values `N/A [e]`, inconsistent formatting, etc.

To simplify the parsing process, we have provided an incomplete code snippet below to help you to fill up the `launch_dict`. Please complete the following code snippet with TODOs or you can choose to write your own logic to parse all launch tables:



```
1 extracted_row = 0
2 #Extract each table
3 for table_number,table in enumerate(soup.find_all('table',"wikitable plainrowheaders c
4     # get table row
5     for rows in table.find_all("tr"):
6         #check to see if first table heading is as number corresponding to launch a nu
7         if rows.th:
8             if rows.th.string:
9                 flight_number=rows.th.string.strip()
10                flag=flight_number.isdigit()
11            else:
12                flag=False
13        #get table element
14        row=rows.find_all('td')
15        #if it is number save cells in a dictionary
16        if flag:
17            extracted_row += 1
18            1# Flight Number value
19            # Append the flight_number into launch_dict with key `Flight No.`
20            launch_dict["Flight No."].append(flight_number)
21            #print(flight_number)
22            datatimelist=date_time(row[0])
23
24            2# Date value
25            # Append the date into launch_dict with key `Date`
26            date = datatimelist[0].strip(',')
27            launch_dict["Date"].append(date)
28            #print(date)
29
30            3# Time value
31            # Append the time into launch_dict with key `Time`
32            time = datatimelist[1]
33            launch_dict["Time"].append(time)
34            #print(time)
35
36            4# Booster version
37            # Append the bv into launch_dict with key `Version Booster`
38            bv=booster_version(row[1])
39            if not(bv):
40                bv=row[1].a.string
41            #print(bv)
42            launch_dict["Version Booster"].append(bv)
43
44            6# Launch Site
45            # Append the bv into launch_dict with key `Launch Site`
46            launch_site = row[2].a.string
47            launch_dict["Launch site"].append(launch_site)
48            #print(launch_site)
49
50            7# Payload
```

```

51     # Append the payload into launch_dict with key `Payload`
52     payload = row[3].a.string
53     launch_dict["Payload"].append(payload)
54     #print(payload)
55
56     8# Payload Mass
57     # Append the payload_mass into launch_dict with key `Payload mass`
58     payload_mass = get_mass(row[4])
59     launch_dict["Payload mass"].append(payload_mass)
60     #print(payload)
61
62     9# Orbit
63     # Append the orbit into launch_dict with key `Orbit`
64     orbit = row[5].a.string
65     launch_dict["Orbit"].append(orbit)
66     #print(orbit)
67
68     10# Customer
69     # Append the customer into launch_dict with key `Customer`
70     customer = row[6]
71     launch_dict["Customer"].append(customer)
72     #print(customer)
73
74     11# Launch outcome
75     # Append the launch_outcome into launch_dict with key `Launch outcome`
76     launch_outcome = list(row[7].strings)[0]
77     launch_dict["Launch outcome"].append(launch_outcome)
78     #print(launch_outcome)
79
80     12# Booster landing
81     # Append the launch_outcome into launch_dict with key `Booster landing`
82     booster_landing = landing_status(row[8])
83     launch_dict["Booster landing"].append(booster_landing)
84     #print(booster_landing)

```

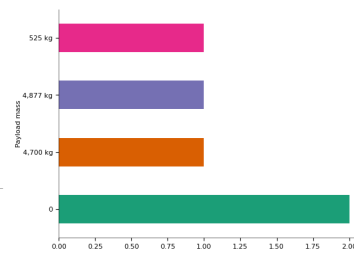
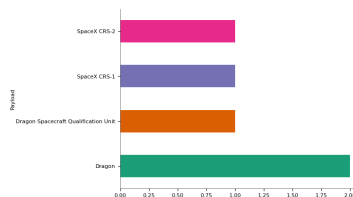
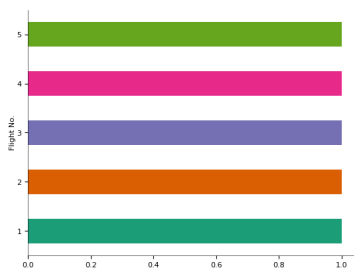
```
1 df=pd.DataFrame(launch_dict)
```

```
1 df.head()
```

	Flight No.	Launch site	Payload	Payload mass	Orbit	Customer	Launch outcome	Version Booster	Boost landi
0	1	CCAFS	Dragon Spacecraft Qualification Unit	0	LEO	[[SpaceX], \n]	Success\n	F9 v1.0B0003.1	Fail
1	2	CCAFS	Dragon	0	LEO	[[.mw- parser- output n\nainlist	Success	F9 v1.0B0003.1	Fail

						[[NASA], (, [COTS], )\n]	Success	v1.0B0005.1	F9 attemp
2	3	CCAFS	Dragon	525 kg	LEO				
3	4	CCAFS	SpaceX CRS-1	4,700 kg	LEO	[[NASA], (, [CRS], )\n]	Success\n	v1.0B0006.1	F9 atten
4	5	CCAFS	SpaceX CRS-2	4,877 kg	LEO	[[NASA], (, [CRS], )\n]	Success\n	v1.0B0007.1	F9 attemp

### Categorical distributions



```

-----
TypeError                                Traceback (most recent call last)
/usr/local/lib/python3.10/dist-packages/pandas/core/algorithms.py in
safe_sort(values, codes, use_na_sentinel, assume_unique, verify)
    1542         try:
-> 1543             sorter = values.argsort()
    1544             ordered = values.take(sorter)

```

**TypeError:** '<' not supported between instances of 'Tag' and 'Tag'

During handling of the above exception, another exception occurred:

```

-----
TypeError                                Traceback (most recent call last)

```

16 frames

**TypeError:** '<' not supported between instances of 'Tag' and 'Tag'

During handling of the above exception, another exception occurred:

```

-----
TypeError                                Traceback (most recent call last)

```

<string> in <module>

```

/usr/local/lib/python3.10/dist-packages/numpy/core/fromnumeric.py in _wrapit(obj,
method, *args, **kws)

```

```

    43     except AttributeError:
    44         wrap = None
--> 45     result = getattr(asarray(obj), method)(*args, **kws)
    46     if wrap:
    47         if not isinstance(result, mu.ndarray):

```

**TvpeError:** '<' not supported between instances of 'Tag' and 'Tag'

```
from matplotlib import pyplot as plt
import seaborn as sns
_df_3.groupby('Customer').size().plot(kind='barh',
color=sns.palettes.mpl_palette('Dark2'))
plt.gca().spines[['top', 'right',]].set_visible(False)
```

### Time series



```
-----
TypeError                                Traceback (most recent call last)
/usr/local/lib/python3.10/dist-packages/pandas/core/algorithms.py in
safe_sort(values, codes, use_na_sentinel, assume_unique, verify)
    1542         try:
-> 1543             sorter = values.argsort()
    1544             ordered = values.take(sorter)
```

TypeError: '<' not supported between instances of 'Tag' and 'Tag'

During handling of the above exception, another exception occurred:

```
-----
TypeError                                Traceback (most recent call last)
_____ ^ 16 frames _____
TypeError: '<' not supported between instances of 'Tag' and 'Tag'
```

During handling of the above exception, another exception occurred:

```
TypeError                                Traceback (most recent call last)
<string> in <module>

/usr/local/lib/python3.10/dist-packages/numpy/core/fromnumeric.py in _wrapit(obj,
method, *args, **kwargs)
    43     except AttributeError:
    44         wrap = None
--> 45     result = getattr(asarray(obj), method)(*args, **kwargs)
    46     if wrap:
    47         if not isinstance(result, mu.ndarray):
```

TypeError: '<' not supported between instances of 'Tag' and 'Tag'

```
from matplotlib import pyplot as plt
import seaborn as sns
def _plot_series(series, series_name, series_index=0):
    palette = list(sns.palettes.mpl_palette('Dark2'))
    counted = (series['Date']
                .value_counts()
                .reset_index(name='counts')
                .rename({'index': 'Date'}, axis=1)
                .sort_values('Date', ascending=True))
```

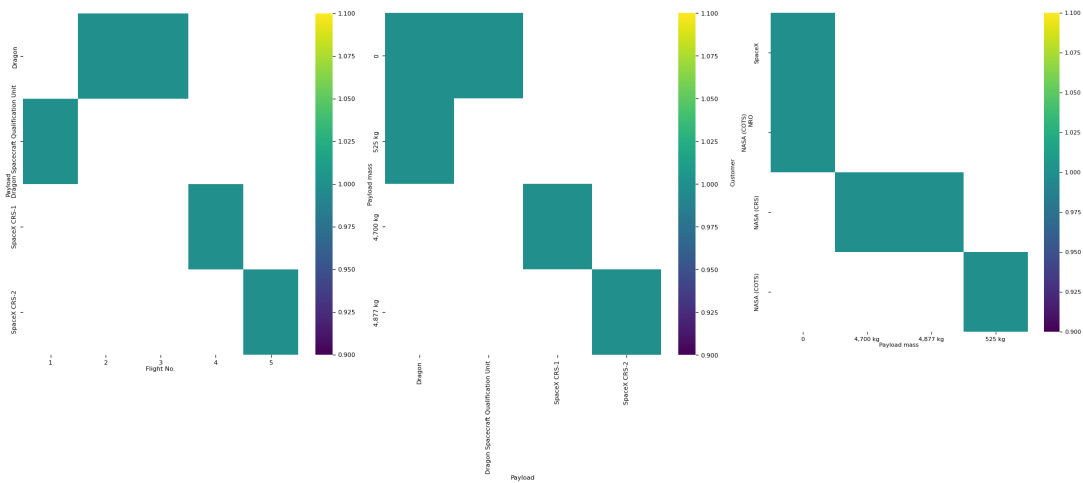
```

xs = counted['Date']
ys = counted['counts']
plt.plot(xs, ys, label=series_name, color=palette[series_index % len(palette)])

fig, ax = plt.subplots(figsize=(10, 5.2), layout='constrained')
df_sorted = _df_7.sort_values('Date', ascending=True)
for i, (series_name, series) in enumerate(df_sorted.groupby('Customer')):
    _plot_series(series, series_name, i)
    fig.legend(title='Customer', bbox_to_anchor=(1, 1), loc='upper left')
sns.despine(fig=fig, ax=ax)
plt.xlabel('Date')
_ = plt.ylabel('count()')

```

## 2-d categorical distributions



Next  
steps:

Generate  
code with df



recommended

Explain  
error

Explain  
error

Explain  
error

After you have fill in the parsed launch record values into `launch_dict`, you can create a dataframe from it.

```
1 df= pd.DataFrame({ key:pd.Series(value) for key, value in launch_dict.items() })
```

We can now export it to a **CSV** for the next section, but to make the answers consistent and in case you have difficulties finishing this lab.

Following labs will be using a provided dataset to make each lab independent.

```
df.to_csv('spacex_web_scraped.csv', index=False)
```

## ▼ Authors

View full

[Yan Luo](#)

[Nayef Abou Tayoun](#)

## ▼ Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2021-06-09	1.0	Yan Luo	Tasks updates
2020-11-10	1.0	Nayef	Created the initial version

Copyright © 2021 IBM Corporation. All rights reserved.