



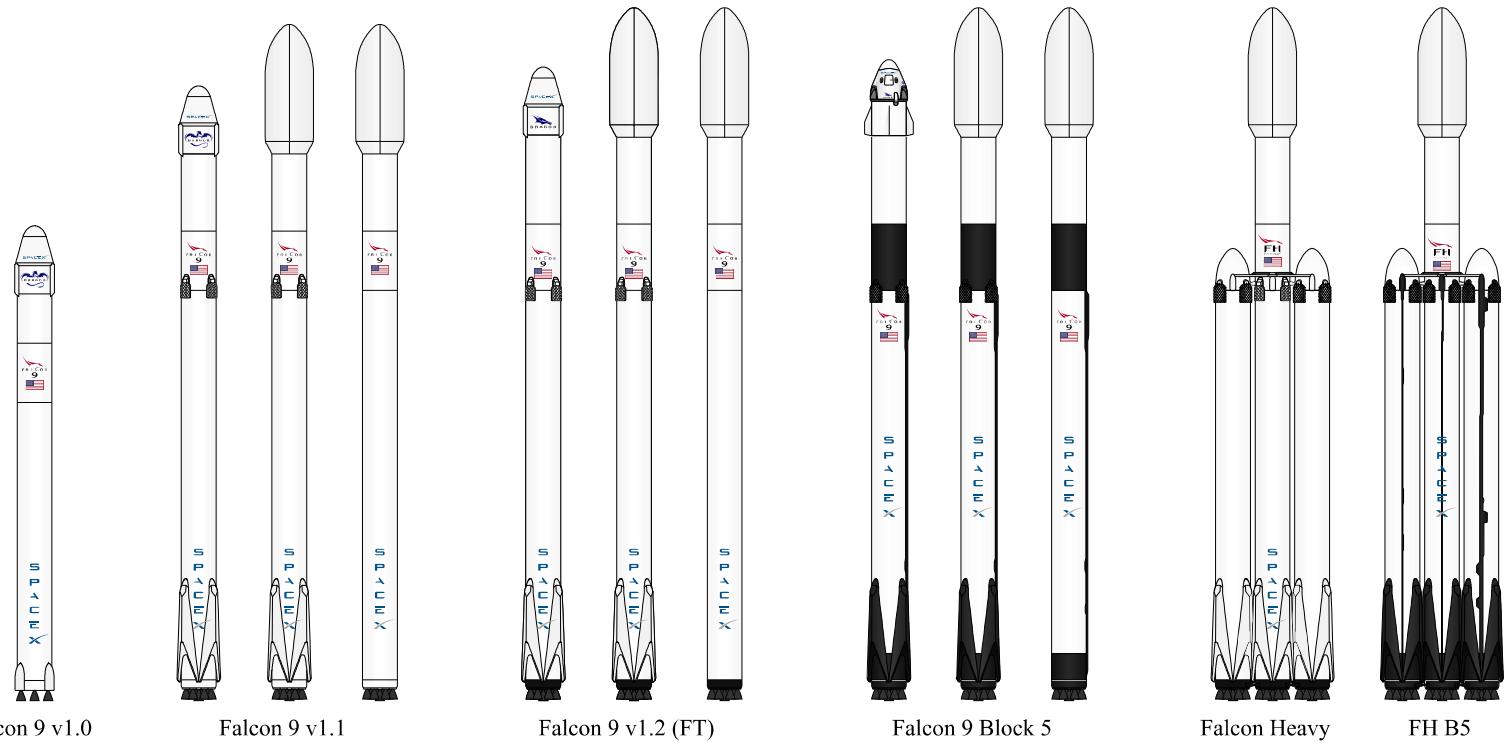
# Space X Falcon 9 First Stage Landing Prediction

## 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)



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:

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 Δ B1049.4	CCAFS, SLC-40	Starlink 2 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[5]</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 Δ 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 <a href="#">Dragon 2</a> abort system after <a href="#">Max Q</a> . The capsule fired its <a href="#">SuperDraco</a> engines, reached an apogee of 40 km (25 mi), deployed parachutes after reentry, and <a href="#">splashed down</a> in the ocean 31 km (19 mi) downrange from the launch site. The test was previously slated to be accomplished with the <a href="#">Crew Dragon Demo-1</a> 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 <a href="#">mass simulator</a> in place of its engine.								
80	29 January 2020, 14:07 <sup>[501]</sup>	F9 B5 Δ B1051.3	CCAFS, SLC-40	Starlink 3 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[5]</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 Δ B1056.4	CCAFS, SLC-40	Starlink 4 v1.0 (60 satellites)	15,600 kg (34,400 lb) <sup>[5]</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 Δ B1059.2	CCAFS, SLC-40	SpaceX CRS-20 (Dragon C112.3 Δ)	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 <a href="#">ESA</a> 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 swap out the second stage instead of replacing the faulty part. <sup>[509]</sup> It was SpaceX's 50th successful landing of a first stage booster, the third flight of the Dragon C112 and the last launch of the cargo <a href="#">Dragon</a> 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 <a href="#">Merlin 1D</a> 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

```
In [1]: !pip3 install beautifulsoup4
!pip3 install requests
```

```
Requirement already satisfied: beautifulsoup4 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (4.11.1)
Requirement already satisfied: soupsieve>1.2 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from beautifulsoup4) (2.3.2.post1)
Requirement already satisfied: requests in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (2.29.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from requests) (3.1.0)
Requirement already satisfied: idna<4,>=2.5 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from requests) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from requests) (1.26.15)
Requirement already satisfied: certifi>=2017.4.17 in /home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (from requests) (2023.5.7)
```

```
In [2]: import sys
```

```
import requests
from bs4 import BeautifulSoup
import re
import unicodedata
import pandas as pd
```

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

```
In [3]: def date_time(table_cells):
```

```
    """
    This function returns the data and time from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    return [data_time.strip() for data_time in list(table_cells.strings)][0:2]
```

```
def booster_version(table_cells):
```

```
    """
    This function returns the booster version from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    out=''.join([booster_version for i,booster_version in enumerate( table_cells.strings) if i%2==0][0:-1])
    return out
```

```
def landing_status(table_cells):
```

```
    """
    This function returns the landing status from the HTML table cell
    Input: the element of a table data cell extracts extra row
    """
    out=[i for i in table_cells.strings][0]
    return out
```

```
def get_mass(table_cells):
```

```

mass=unicodedata.normalize("NFKD", table_cells.text).strip()
if mass:
    mass.find("kg")
    new_mass=mass[0:mass.find("kg")+2]
else:
    new_mass=0
return new_mass

def extract_column_from_header(row):
"""
This function returns the landing status from the HTML table cell
Input: the element of a table data cell extracts extra row
"""
if (row.br):
    row.br.extract()
if row.a:
    row.a.extract()
if row.sup:
    row.sup.extract()

column_name = ' '.join(row.contents)

# Filter the digit and empty names
if not(column_name.strip().isdigit()):
    column_name = column_name.strip()
return column_name

```

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](#) Wikipedia updated on [9th June 2021](#)

In [4]: static\_url = "https://en.wikipedia.org/w/index.php?title=List\_of\_Falcon\_9\_and\_Falcon\_Heavy\_launches&oldid=1027686922"

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

First, let's perform an HTTP GET method to request the Falcon9 Launch HTML page, as an HTTP response.

In [8]: # use requests.get() method with the provided static\_url  
# assign the response to a object  
r = requests.get(static\_url)  
data = r.text  
soup = BeautifulSoup(data, "html.parser")

Create a `BeautifulSoup` object from the HTML response

```
In [10]: # Use BeautifulSoup() to create a BeautifulSoup object from a response text content
print(soup.title)
```

```
<title>List of Falcon 9 and Falcon Heavy launches - Wikipedia</title>
```

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

```
In [11]: # Use soup.title attribute
print(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

```
In [12]: # Use the find_all function in the BeautifulSoup object, with element type `table`
# Assign the result to a list called `html_tables`
html_tables = soup.find_all('table')
```

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

```
In [14]: # Let's print the third table and check its content
first_launch_table = html_tables[2]
#print(first_launch_table)
```

You should able to see the columns 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 Universal Time">UTC</a>)
</th>
<th scope="col"><a href="/wiki/List_of_Falcon_9_first-stage_boosters" title="List of Falcon 9 first-stage boosters">Version,<br/>Booster</a> <sup class="reference" id="cite_ref-booster_11-0"><a href="#cite_note-booster-11">[b]</a></sup>
```

```
</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-12">[c]</a>
</sup>
</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 landing tests">Booster<br/>landing</a>
</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

```
In [15]: column_names = []

# Apply find_all() function with `th` element on first_launch_table
# Iterate each th element and apply the provided extract_column_from_header() to get a column name
# Append the Non-empty column name (`if name is not None and Len(name) > 0`) into a List called column_names
table_headers = first_launch_table.find_all('th')
# print(table_headers)

for j, table_header in enumerate(table_headers):
    name = extract_column_from_header(table_header)
    if name is not None and len(name) > 0:
        column_names.append(name)

print(column_names)
```

['Flight No.', 'Date and time ( )', 'Launch site', 'Payload', 'Payload mass', 'Orbit', 'Customer', 'Launch outcome']

Check the extracted column names

```
In [16]: print(column_names)

['Flight No.', 'Date and time ( )', 'Launch site', 'Payload', 'Payload mass', 'Orbit', 'Customer', 'Launch outcome']
```

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

We will 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

In [17]:

```
launch_dict= dict.fromkeys(column_names)

# Remove an irrelevant column
del launch_dict['Date and time ( )']

# Let's initial the launch_dict with each value to be an empty list
launch_dict['Flight No.'] = []
launch_dict['Launch site'] = []
launch_dict['Payload'] = []
launch_dict['Payload mass'] = []
launch_dict['Orbit'] = []
launch_dict['Customer'] = []
launch_dict['Launch outcome'] = []
# Added some new columns
launch_dict['Version Booster']=[]
launch_dict['Booster landing']=[]
launch_dict['Date']=[]
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:

In [18]:

```
extracted_row = 0
for table_number,table in enumerate(soup.find_all('table',"wikitable plainrowheaders collapsible")):
    for rows in table.find_all("tr"): # get table row
        if rows.th: #check to see if first table heading is a number corresponding to Launch a number
            if rows.th.string:
                flight_number=rows.th.string.strip()
                flag=flight_number.isdigit()
            else:
                flag=False
            row=rows.find_all('td') # #get table element
            if flag: #if it is number save cells in a dictionary
                extracted_row += 1
                # Flight Number value
                # TODO: Append the flight_number into Launch_dict with key `Flight No.`
```

```
launch_dict['Flight No.'].append(flight_number)
# print(flight_number)
datatimelist=date_time(row[0])

# Date value
# TODO: Append the date into Launch_dict with key `Date`
date = datatimelist[0].strip(',')
launch_dict['Date'].append(date)
# print(date)

# Time value
# TODO: Append the time into Launch_dict with key `Time`
time = datatimelist[1]
launch_dict['Time'].append(time)
# print(time)

# Booster version
# TODO: Append the bv into Launch_dict with key `Version Booster`
bv=booster_version(row[1])
if not(bv):
    bv=row[1].a.string
launch_dict['Version Booster'].append(bv)
# print(bv)

# Launch Site
# TODO: Append the bv into launch_dict with key `Launch Site`
launch_site = row[2].a.string
launch_dict['Launch site'].append(launch_site)
# print(launch_site)

# Payload
# TODO: Append the payload into Launch_dict with key `Payload`
payload = row[3].a.string
launch_dict['Payload'].append(payload)
# print(payload)

# Payload Mass
# TODO: Append the payload_mass into Launch_dict with key `Payload mass`
payload_mass = get_mass(row[4])
launch_dict['Payload mass'].append(payload_mass)
# print(payload)

# Orbit
# TODO: Append the orbit into Launch_dict with key `Orbit`
orbit = row[5].a.string
launch_dict['Orbit'].append(orbit)
# print(orbit)
```

```

# Customer
# TODO: Append the customer into Launch_dict with key `Customer`
try:
    customer = row[6].a.string
except:
    customer = "None"

launch_dict['Customer'].append(customer)
# print(customer)

# Launch outcome
# TODO: Append the launch_outcome into Launch_dict with key `Launch outcome`
launch_outcome = list(row[7].strings)[0]
launch_dict['Launch outcome'].append(launch_outcome)
# print(launch_outcome)

# Booster landing
# TODO: Append the Booster_landing into Launch_dict with key `Booster Landing`
booster_landing = landing_status(row[8])
launch_dict['Booster landing'].append(booster_landing)
# print(booster_landing)

# debugging: find length differences
# print("-----")
# curr_len = [len(val) for key, val in Launch_dict.items()]
# print(curr_len)
# print("-----")
print("number of extracted rows: ", extracted_row)

```

number of extracted rows: 121

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

In [19]: `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

Yan Luo

## 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.