Discussion Forums

# Week 4

| SUBFORUMS |
|---|
| **All** |
| Assignment: Visualizing Two Datasets |

← Week 4

---

CF  ## How to download data from a website into a CSV file?　　📌 ⌄
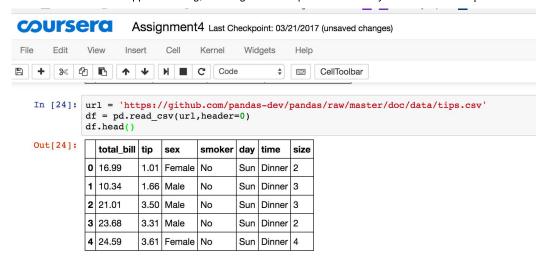
claudio freitas Week 4 · a year ago

Can anyone provide an example?

⇧ 2 Upvotes　　　💬 Reply　　　Follow this discussion

---

**HIGHLIGHTED POST**

SG　Sophie Greene　Teaching Staff　· a year ago · Edited　　⌄

If you have a link to the CSV file, pandas read_csv can be used exactly the same way, instead of the file name, you pass the url

here is an example

```
coursera      Assignment4 Last Checkpoint: 03/21/2017 (unsaved changes)

File    Edit    View    Insert    Cell    Kernel    Widgets    Help

[save] + ✂ ⧉ 🗏 ↑ ↓ ▶ ■ C    Code  ⬍    ⌨ CellToolbar
```

```
In [24]:  url = 'https://github.com/pandas-dev/pandas/raw/master/doc/data/tips.csv'
          df = pd.read_csv(url,header=0)
          df.head()
```

Out[24]:

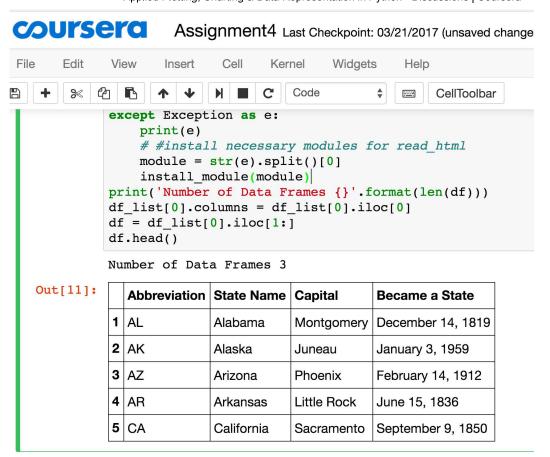|   | total_bill | tip | sex | smoker | day | time | size |
|---|------------|-----|-----|--------|-----|------|------|
| 0 | 16.99 | 1.01 | Female | No | Sun | Dinner | 2 |
| 1 | 10.34 | 1.66 | Male | No | Sun | Dinner | 3 |
| 2 | 21.01 | 3.50 | Male | No | Sun | Dinner | 3 |
| 3 | 23.68 | 3.31 | Male | No | Sun | Dinner | 2 |
| 4 | 24.59 | 3.61 | Female | No | Sun | Dinner | 4 |

If you have a web page containing a table and you want to read that table you can use read_html. You will need to load the missing modules like lxml in the online platform, that said, the changes are not permanent, i.e. the modules you install will be deleted when you log out /restart the Jupyter server

I've written a script that can help load the needed modules automatically

```
1    import pandas as pd
2    from IPython.display import display, HTML
3
4
5    def install_module(module):
6        ! conda install "$module" -y
7        js_cmd = ['IPython.notebook.kernel.restart();',
8                  'IPython.notebook.select(1);',
9                  'IPython.notebook.execute_cell();'
10                 ]
11       js = "<script>{0}</script>".format(' '.join(js_cmd))
12       display(HTML(js))
13
14   url = 'https://simple.wikipedia.org/wiki/List_of_U.S._states'
15   try:
16       df_list = pd.read_html(url)
17   except Exception as e:
18       print(e)
19       # #install necessary modules for read_html
20       module = str(e).split()[0]
21       install_module(module)
22   print('Number of Data Frames {}'.format(len(df_list)))
23   df_list[0].columns = df_list[0].iloc[0]
24   df = df_list[0].iloc[1:]
25   df.head()
```

the script will automatically restart the kernel to propagate the changes and will return the list tables available in the input url page

**coursera**     Assignment4 Last Checkpoint: 03/21/2017 (unsaved change

| File | Edit | View | Insert | Cell | Kernel | Widgets | Help |
|------|------|------|--------|------|--------|---------|------|

💾  ➕  ✂  🗐  🗎  ↑  ↓  ▶|  ■  C     Code     ⬍     ⌨  CellToolbar

```python
    except Exception as e:
        print(e)
        # #install necessary modules for read_html
        module = str(e).split()[0]
        install_module(module)
print('Number of Data Frames {}'.format(len(df)))
df_list[0].columns = df_list[0].iloc[0]
df = df_list[0].iloc[1:]
df.head()
```

Number of Data Frames 3

Out[11]:

|   | Abbreviation | State Name | Capital | Became a State |
|---|--------------|------------|---------|----------------|
| 1 | AL | Alabama | Montgomery | December 14, 1819 |
| 2 | AK | Alaska | Juneau | January 3, 1959 |
| 3 | AZ | Arizona | Phoenix | February 14, 1912 |
| 4 | AR | Arkansas | Little Rock | June 15, 1836 |
| 5 | CA | California | Sacramento | September 9, 1850 |

This way your peer reviewers will be able to run the code in the online platform without issues

If you want to save the html table you read using read-html to datafile you can use

```
1   #csv
2   df.to_csv('filename.csv')
3   #or excel might need module xlwt installed
4   df.to_excel('file.xls')
5   
```

⇧ 7 Upvotes      Jump to post

**Earliest    Top    Most Recent**

⌄

Pierre Masson · a year ago

Hi Sophie,

I have tried to use the script. As a result, I receive the following message

lxml not found, please install it