

AMIA 2019 Informatics Summit

Workshop W07: Hands-On Full Life Cycle Data Science Workshop

Thank you for your interest in the workshop. This is a “hands-on” workshop which means that you will need to bring a laptop so that you can follow along, execute code and do short exercises as we work through aspects of a data science project. You can either 1) run it in the cloud or 2) install the environment on your laptop and execute it locally.

1. The advantage to running the environment in the cloud is that there is nothing to install! We will be using the Google Colab environment. One disadvantage is that the environment will shut down after 12 hours of use. That is fine for the purpose of this workshop, but probably is not what you’d want to use for a production data science environment.
2. The advantage to installing it locally is that everything you need for a data science project is available on your computer. The disadvantage is that it is sometimes difficult and time consuming to install all of the elements of the environment. To install it locally, we will use the Anaconda Python distribution which has all of the packages that we will be using. It should install equally well on with Windows or Mac machines.

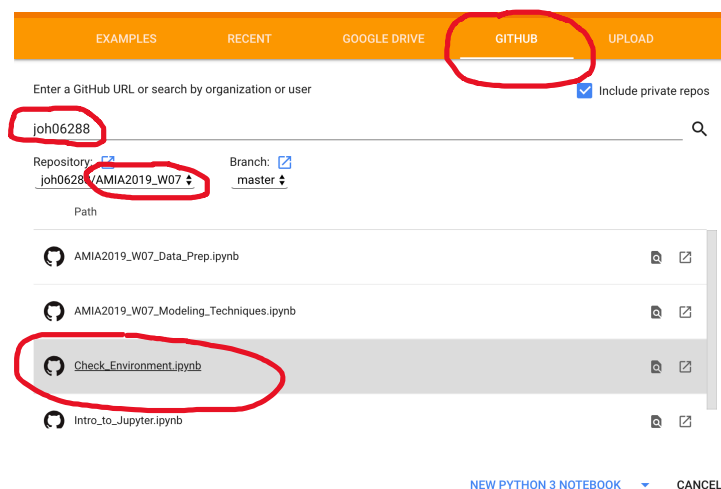
1. Access Jupyter notebooks in the cloud using Google Colab

This is probably the easiest way to get started. You can access Jupyter notebooks in a shared cloud environment using a Google research project called Colab. This has the advantage of letting you try out a notebook without having to install anything on your local computer. But the disadvantage is that the session will timeout after 12 hours. It is not suitable for doing real data science work, but it is a good option for this workshop.

To access the system, go to the following URL in a browser (Chrome works best):

<https://colab.research.google.com/> This will launch a new virtual environment where you can play with the workshop notebooks. Check to make sure the environment is working by doing the following:

1. When you go <https://colab.research.google.com/>, a pop-up screen will ask you which notebook to open. Click the “GITHUB” tab and then put in joh06288 for the Github URL and find the AMIA2019_W07 repository and select the “Check_Environment” notebook:



2. Navigate to the main menu and select “Runtime -> Run all”. You will be warned that the notebook is not from Google and that your runtimes will be reset. Pick “Run anyway” and “Yes” to run the notebook. If

Warning: This notebook was not authored by Google.

This notebook is being loaded from [GitHub](#). It may request access to your data stored with Google, or read data and credentials from other sessions. Please review the source code before executing this notebook. To prevent this notebook reading state from other sessions, you can reset all runtimes.

☒ Reset all runtimes before running

CANCEL RUN ANYWAY

Reset all runtimes

Are you sure you want to reset all runtimes? State of all runtimes, including all local variables and files, will be lost.

CANCEL YES

The screenshot shows a Jupyter Notebook titled 'Check_Environment.ipynb'. The 'Runtime' menu is open, with 'Run all' highlighted. The code cell contains the following Python code:

```
# Check to make sure all the necessary libraries are installed at the proper version level.

import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import matplotlib
import seaborn as sns
from distutils.version import StrictVersion
import IPython

print("numpy version: %6.6s" % np.__version__)
print("pandas version: %6.6s" % pd.__version__)
print("matplotlib version: %6.6s" % matplotlib.__version__)
print("IPython version: %6.6s" % IPython.__version__)
print("seaborn version: %6.6s" % sns.__version__)

if StrictVersion(np.__version__) >= StrictVersion('1.13.0') and \
    StrictVersion(pd.__version__) >= StrictVersion('0.20.0') and \
    StrictVersion(matplotlib.__version__) >= StrictVersion('2.0.0') and \
    StrictVersion(IPython.__version__) >= StrictVersion('5.5.0') and \
    StrictVersion(sns.__version__) >= StrictVersion('0.7.0'):
    print('\nCongratulations, your environment is setup correctly!')
else:
    print('\nEnvironment is NOT setup correctly!')
```

The output of the code cell shows the versions of the libraries:

```
numpy version: 1.14.6
pandas version: 0.22.0
matplotlib version: 3.0.3
IPython version: 5.5.0
seaborn version: 0.7.1

Congratulations, your environment is setup correctly!
```

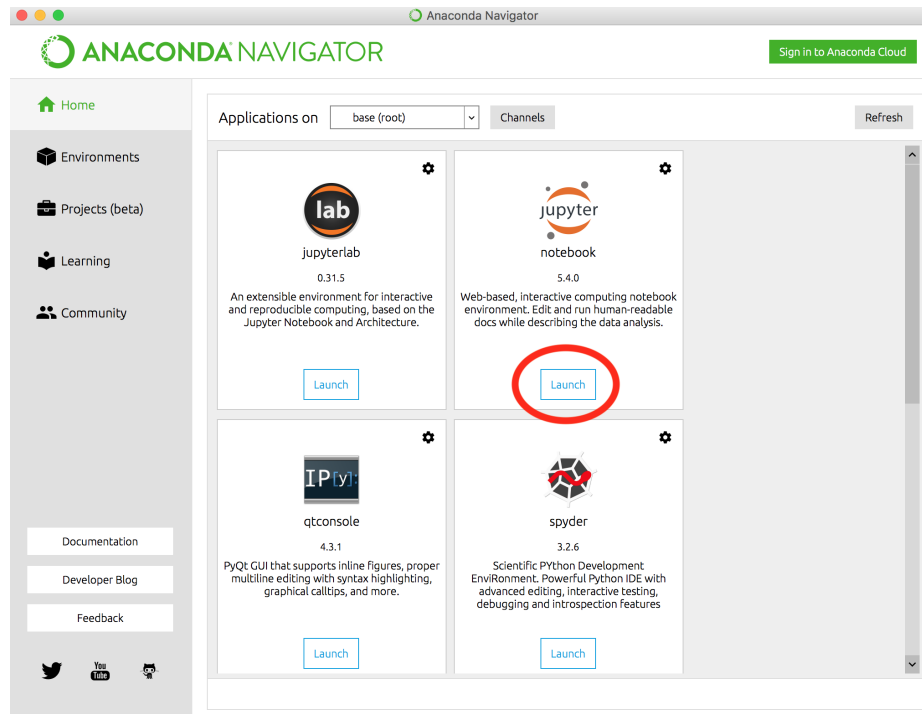
3. The final message should say “Congratulations, your environment is setup correctly!”
4. If you have problems, there will be time at the Workshop to help you troubleshoot. Please come to the workshop meeting room 30 minutes early on the day of the workshop and we can help resolve the problem.
5. If you are interested, you can take a look at the “Intro_to_Jupyter” notebook which is a quick overview of some of the libraries that we will be using during the workshop.

2. Install locally using the Anaconda distribution

If you don't want to use the cloud, follow the instructions to install the full Anaconda Distribution (<https://www.anaconda.com/download>), which includes Python, the Jupyter Notebook system, and other commonly used packages for scientific computing and data science. Make sure you install the latest Python 3 version of the distribution. After a successful install, you can run the following command at the Terminal (Mac/Linux) or Command Prompt (Windows) to start the Jupyter system:

```
jupyter notebook
```

You can also launch the Jupyter Notebook system from the Anaconda Navigator:



Check to make sure the environment is installed correctly by doing the following:

1. In a CMD prompt (Windows) or Terminal prompt (OSX) create a work directory that you will put the Workshop content in.

```
mkdir workdir
```
2. Retrieve the workshop content from the github server using the following command:

```
git clone https://github.com/joh06288/AMIA2019_W07.git
```
3. From the browser window that Jupyter notebook opened for you when it started, navigate to your workdir and go into the AMIA2019_W07 folder.
4. Select the "Check_Environment" notebook. A new browser window will open up.
5. Select "Cell -> Run All" from the Jupyter Notebook menu. The following output should be produced:

The screenshot shows a Jupyter Notebook interface with the title 'Check_Environment (autosaved)'. The notebook has a menu bar with 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. Below the menu bar is a toolbar with icons for running cells, saving, and other actions. The main area of the notebook contains a code cell with the following code:

```
In [2]: 1 # Check the environment
2
3 import pandas as pd
4 import numpy as np
5 from matplotlib import pyplot as plt
6 import matplotlib as mpl
7 import seaborn as sns
8 from distutils.version import StrictVersion
9 import IPython
10
11 print("numpy version: %s" % np.__version__)
12 print("pandas version: %s" % pd.__version__)
13 print("matplotlib version: %6.6s" % mpl.__version__)
14 print("IPython version: %6.6s" % IPython.__version__)
15 print("seaborn version: %6.6s" % sns.__version__)
16
17 if StrictVersion(np.__version__) >= StrictVersion('1.13.0') and \
18    StrictVersion(pd.__version__) >= StrictVersion('0.20.0') and \
19    StrictVersion(mpl.__version__) >= StrictVersion('2.0.0') and \
20    StrictVersion(IPython.__version__) >= StrictVersion('5.5.0') and \
21    StrictVersion(sns.__version__) >= StrictVersion('0.7.0'):
22     print('\nCongratulations, your environment is setup correctly!')
23 else:
24     print('\nEnvironment is NOT setup correctly!')
25
```

The output of the code cell is displayed below the code:

```
numpy version: 1.16.2
pandas version: 0.24.1
matplotlib version: 2.1.2
IPython version: 6.2.1
seaborn version: 0.8.1
Congratulations, your environment is setup correctly!
```

The 'Run All' option in the 'Cell' menu is circled in red, and the final output message is also circled in red.

6. The final message should say “Congratulations, your environment is setup correctly!”
7. If you have problems, there will be time at the Workshop to help you troubleshoot. Please come to the workshop meeting room 30 minutes early on the day of the workshop and we can help resolve the problem.
8. If you are interested, you can take a look at the “Intro_to_Jupyter” notebook which is a quick overview of some of the libraries that we will be using during the workshop.