Lab 1:

Resources Requred:
Linux Host – Tested on Ubuntu
Windows / MAC Host – For OpenDNS Client

OpenDNS Account – http://www.opendns.com
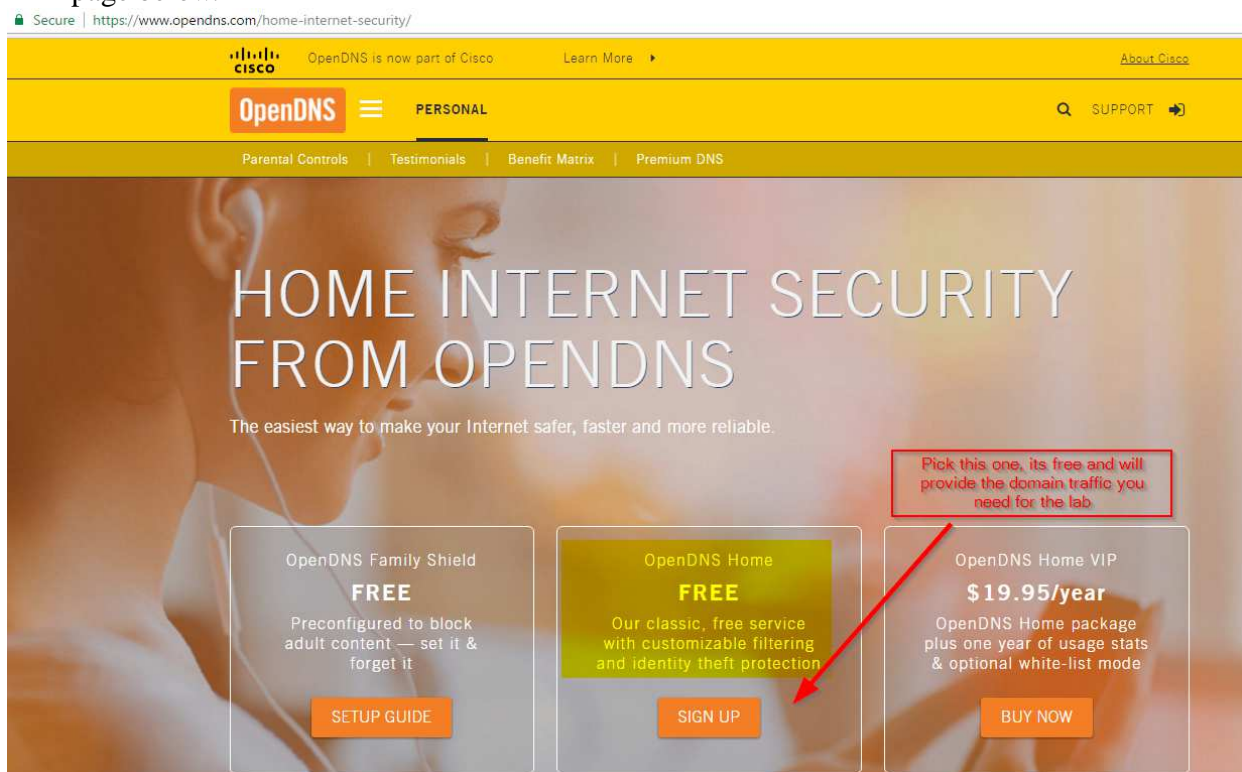Jupyter Notebook  - http://jupyter.org/

A Vmware Virtual Machine is available for download on http://www.MLresearchLab.com. If you want
to have everything pre-built with the sample data
- If you do not have VMware, VMware Player is free or you can convert it to Virtual Box which
is also free

Goal: Utilize machine learning (ML) random forest classifier to analyze OpenDNS
(www.opendns.com) logs of user domain requests to determine if a user has been visiting work related
sites.

1. Create an OpenDNS account by visiting www.opendns.com and click personal to reach the web
page below.



2. Click the orange "Sign Up" button and follow the instructions to create an accont.

3. Once you have the account created, download the client

## Download the OpenDNS Updater

We recommend that you use our **client-side software** to keep your dynamic IP updated for your network.
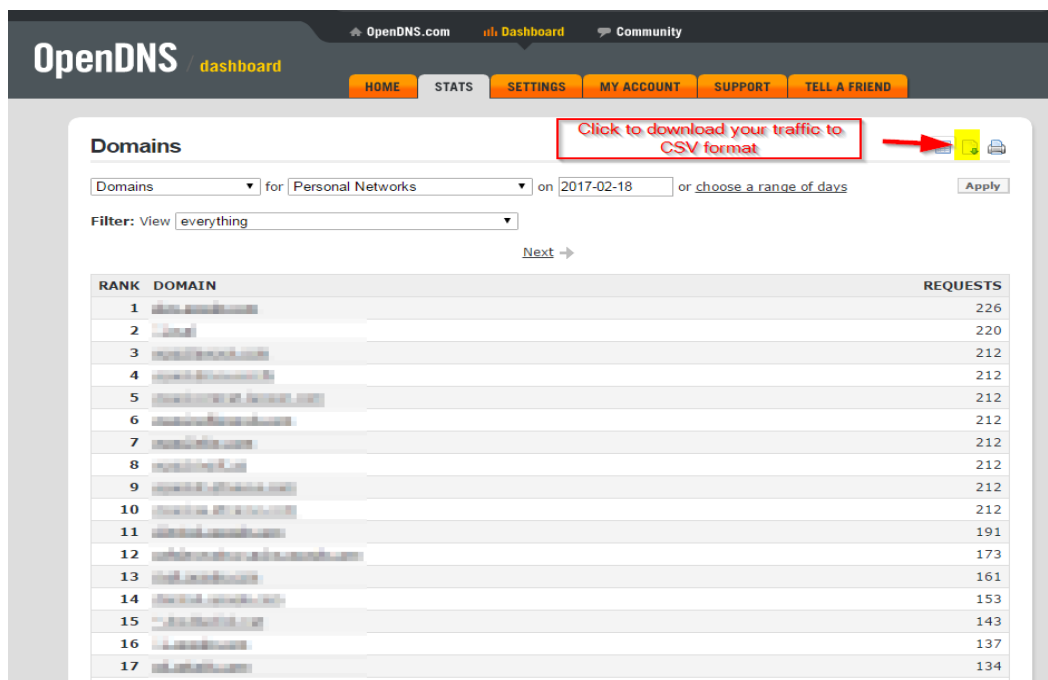
Windows Client          Mac Client

4. Install the client on your test machine or Virtual Machine (VM).

5. Once the client is installed it will ask for the same credentials you signed up under, so enter those in and you should see the screen shot below.



6. Make sure the client under "Using OpenDNS" says "Yes" which means now your DNS traffic is being forwarded to the OpenDNS servers.

7. Browse random websites to generate traffic on the host that has the OpenDNS client, its best to visit at least 50 websites or more

8. Once you have visited 50 or more websites log back into OpenDNS you should see traffic in your dashboard. (If there is no traffic it might take some time to populate so try back a few hours later or a day later)

9. If you have traffic in your dashboard you are now ready to begin to download the traffic as data for the ML lab. Click "Stats" the orange tab in the center of the dashboard then "Domains" on the left pane and you should see something similar to the screenshot below.

10. To download the DNS traffic data you will look for a button of a file download in the top right hand corner of the "Domains" Dashboard. (See screenshot from previous step)

11. Upon Success you should end up with a comma separated value (CSV) file on your computer, go ahead and open it with your favorite editor



12. In your editor and remove the following columns
"Rank", "Total", "Blacklisted", "Blocked by Category", "Blocked as Botnet", "Blocked as Malware", "Blocked as Phishing", and "Resolved by SmartCache" as shown below
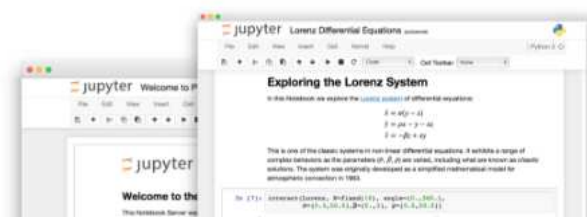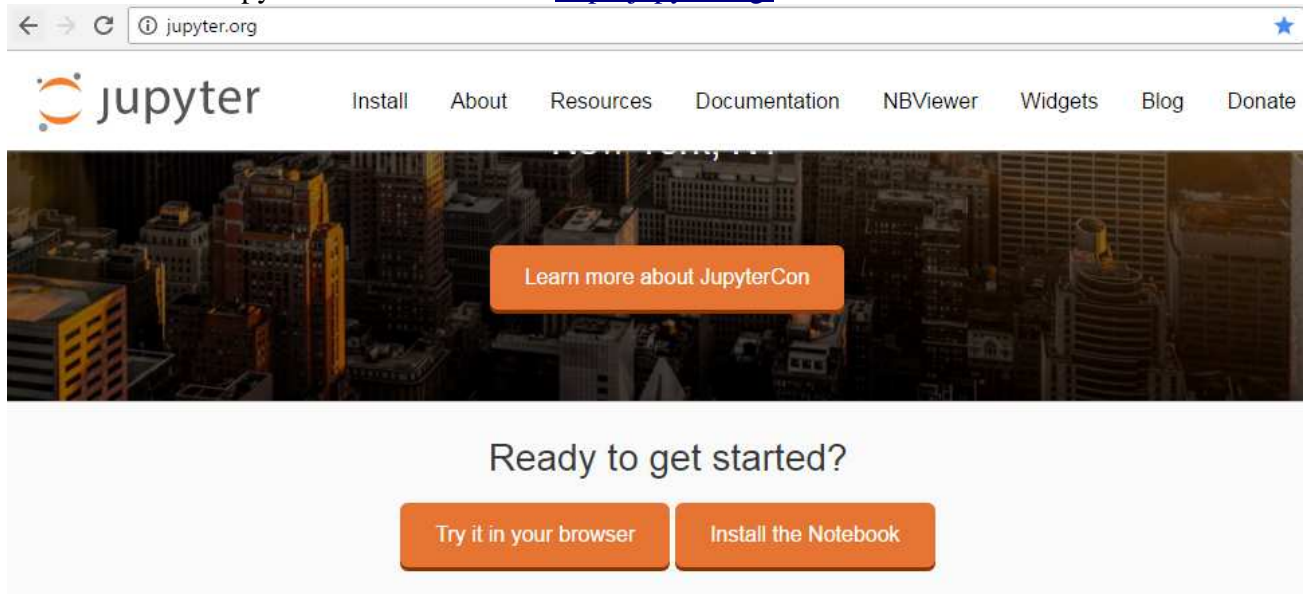
| Domain | Academic Fraud | Adult Themes | Adware | Alcohol | Anime/Manga/Webcomic | Auctions | Automotive | Blogs |
|---|---|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

13. Remember the purpose of this lab is not to determine if the domains your user visited is malicious or not (it is beyond the scope of this lab). Rather that if they visit sites that are not Work Related, so you should end up with the "Domains" and "Categories of websites classifications"(e.g. Academic Fund, Adult Themes, Adware) in your CSV file.

14. To prepare your data, remove rows that have a "0" in all "Categories of websites classifications". These rows do not help the ML and should be removed because they do not have a "1" in any of fields.

15. Now the data inside the CSV needs to be trained for the ML engine so you need to create a new column called "WorkRelated" at the end of the CSV file. So you should now have"Domains", "Categories of websites classifications", and "WorkRelated" as columns inside the CSV file.

16. You must now research the domain that is listed in the "Domain" column and if you determine that domain is work related put a "1" in the "WorkRelated" column else wise put a "0" in the "WorkRelated" column.

17. You should now have a "1" or "0" in the "WorkRelated" column for every row. Once confirmed, remove the "Domain" column because the "Domain" does not help the ML make a decision, the "Categories of website classifications" and "WorkRelated" columns help the ML make the decisions.

18. Your CSV file should only have the "Categories of website classifications" and "WorkRelated" columns and each row should have a "1" in at least one of the "Categories of website

classifications" and a "1" or "0" in "WorkRelated" column.

19. Now the data is ready to be ingested for ML.

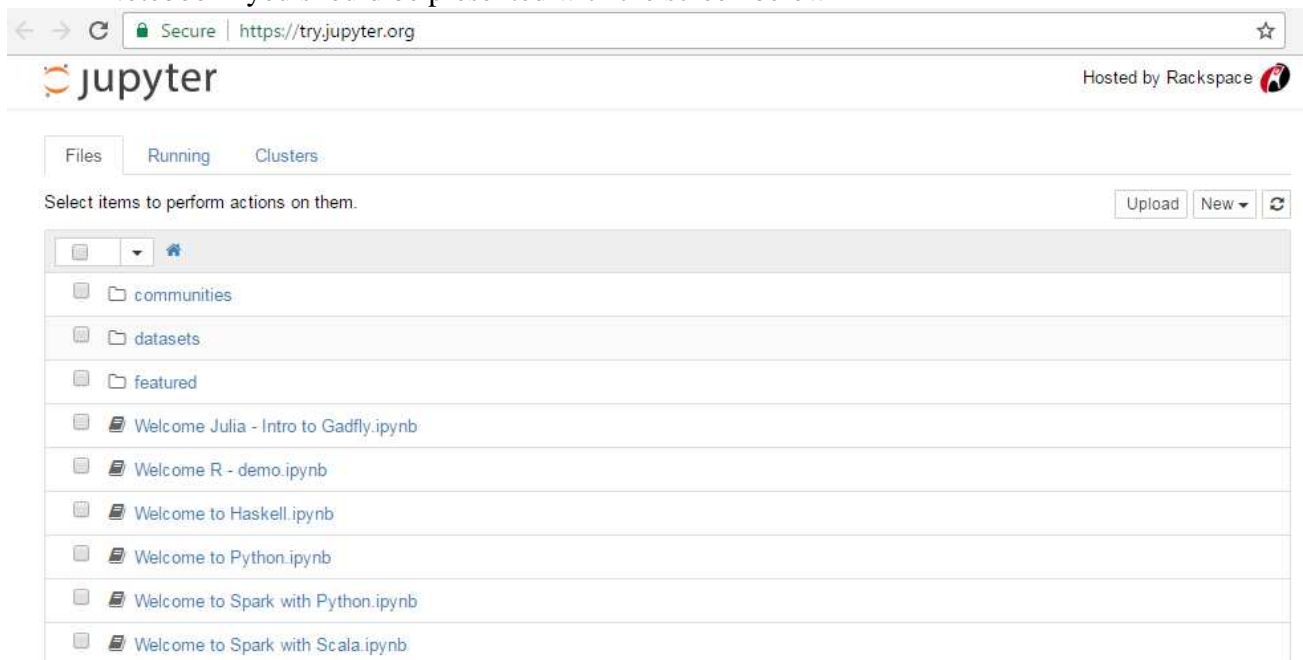20. Visit the Jupyter Notebook website: http://jupyter.org/



21. You have two choices you can choose to try it in your browser and for the lab exercise you can upload the data to the online site OR install Jupyter Notebook locally on your computer.

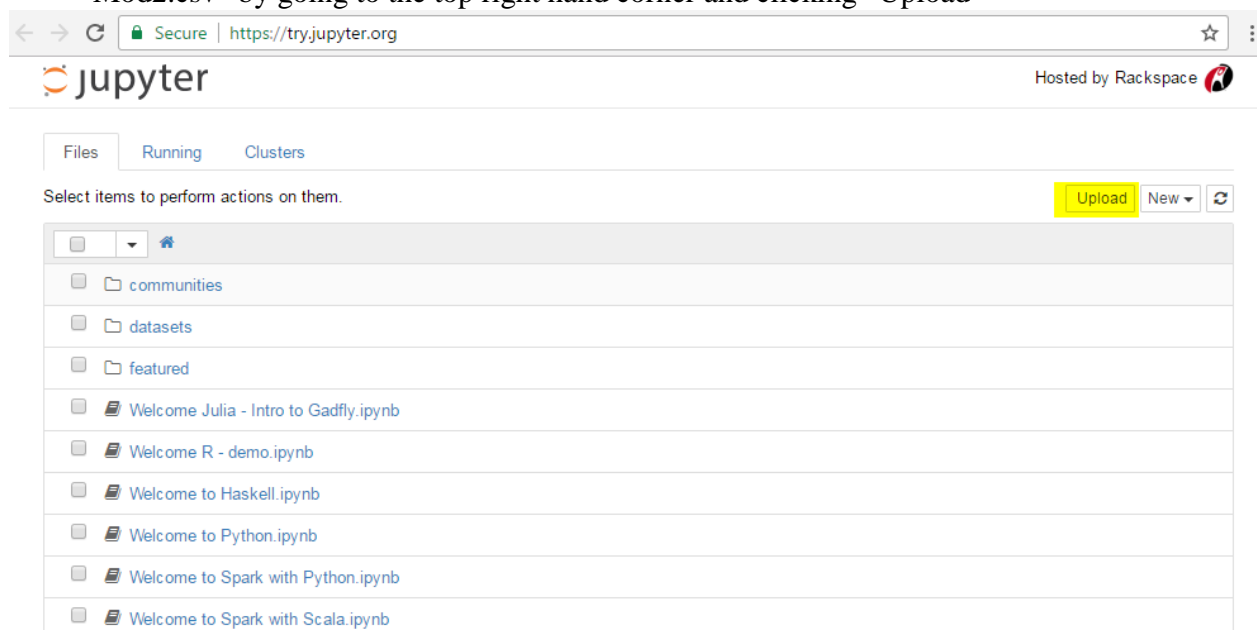It is purely personal preference and if you want to upload your data onto a 3$^{rd}$ party site

The instructions below work for both "Try it in your browser" or "Install the Notebook"
For those that use "Install the Notebook" replace the language "Jupyter website" with your localhost

22. If you click "Try it in your browser" or you started Jupyter from a local install from "Install the Notebook" you should be presented with the screen below



23. First thing is to upload your data to the Jupyter website or use the test sample data provided "Mod2.csv" by going to the top right hand corner and clicking "Upload"

24. The Jupyter website will ask for the file location, once supplied you should see the file ready to be upload as seen in the screenshot below



25. Your data should now be upload and you can verify by seeing if the file shows up in the listing as shown in the screenshot below

26. Now that the data has been uploaded, the next step is to upload the actual ML code so on the top right hand corner click "New" and click "Python 2"



27. You should be presented with a blank notebook as shown below

28. Open the "Lab1-Code.py" code from the Github link on www.MLresearchlab.com and Copy the entire contents and paste into your notebook as shown below



```
In [ ]:  # Created by Dickson Kwong & Kevin Figueroa
         # MLresearchLab.com
         # Lab 1 Random Forest
         # Python 2.7 run inside of Jupyter Notebook
         # [Question to ask ML]
         # Determine based on traffic from OpenDNS.com if a user has been primarily visting work related sites or not

         #MIT License

         #Copyright (c) 2017 Dickson Kwong & Kevin Figueroa

         #Permission is hereby granted, free of charge, to any person obtaining a copy
         #of this software and associated documentation files (the "Software"), to deal
         #in the Software without restriction, including without limitation the rights
         #to use, copy, modify, merge, publish, distribute, sublicense, and/or sell
         #copies of the Software, and to permit persons to whom the Software is
         #furnished to do so, subject to the following conditions:

         #The above copyright notice and this permission notice shall be included in all
         #copies or substantial portions of the Software.

         #THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR
         #IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY,
         #FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE
         #AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER
         #LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM,
         #OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE
         #SOFTWARE.

         # Pythons scientific libraries
         import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         from sklearn.ensemble import RandomForestClassifier
         # You need this import division for python 2.7 to display percentages correctly
         from __future__ import division

         # The line below will only work in Jupyter Notebook -> purpose: it shows plots inside Jupyter Notebook
```

29. Run the entire cell by either clicking the button on the top center of the screen or using the menu dropdown "Cell" → "Run Cell Below" or you can use the keyboard shortcut "Shift-Enter". If your data file is correctly formatted and prepared you should see results as shown below with no errors.

30. Now the results are limited because it only shows the first result and the last result what you want to do is start splitting lines of code into "Cells" so that you can see the individual results of each line of code. You split code into cells by click "Edit" → "Split Cells"

31. Once you start splitting code into "Cells" you start to see individual results such as below when you run that "Cell"



32. If you encounter problems or are unable to run the code please do not hesitate to email "Info@MLresearchLab.com" or on Twitter ML Research Lab@ML _Research_Lab

There is Vmware image (Link on www.MLresearchLab.com) of a local install with Jupyter installed with the sample data all ready to go if you do not want to upload your data online and you want to have a quick environment to learn from.

Questions / Comments / Concerns please do not hesitate to reach out to us we want to help the security community to grow and understand ML

Thank you for trying out the lab and more labs to come!