

Configuration Manual

MSc Internship
Programme Name

Michael Oluwasegun Akinrele
Student ID: X18109489

School of Computing
National College of Ireland

Supervisor: Mr Vikas Sahni

National College of Ireland
MSc Project Submission Sheet
School of Computing



Student Name: Michael Oluwasegun Akinrele
Student ID: X18109489
Programme: CyberSecurity **Year:** 2019
Module: MSc Internship
Lecturer: Mr Vikas Sahni
Submission Due Date: 12/12/2019
Project Title: Detection of Phishing and Spam Emails Using Ensemble Technique
Word Count: 324 **Page Count:** 6

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

I agree to an electronic copy of my thesis being made publicly available on NORMA the National College of Ireland's Institutional Repository for consultation.

Signature:

Date:

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

MICHAEL OLUWASEGUN AKINRELE
X18109489

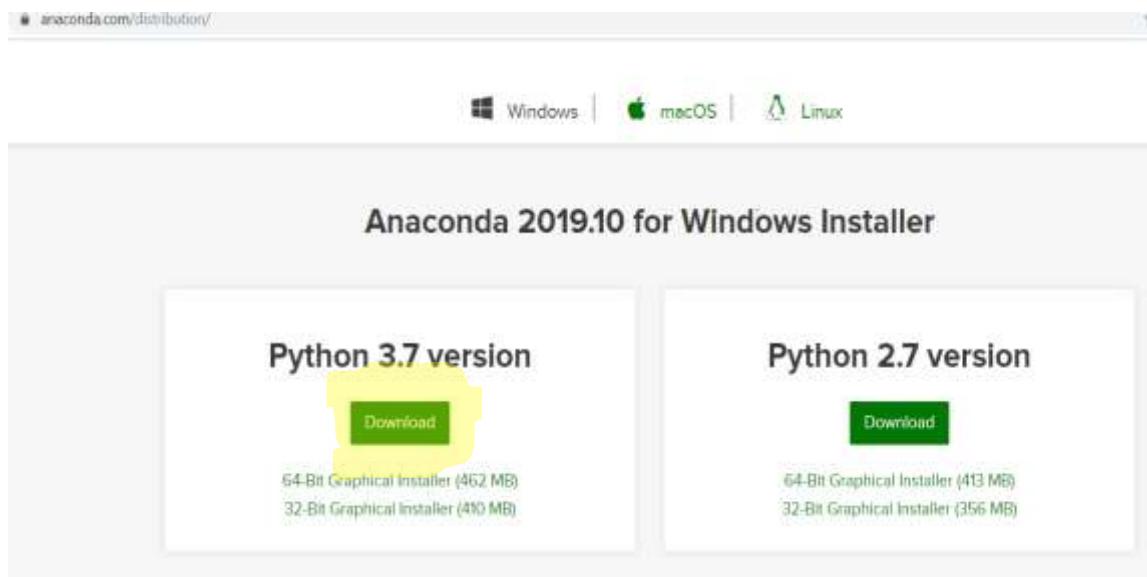
This configuration manual entails the tools capabilities and features used in the course of the project. It gives specific directions on the best replica of the experiment carried out.

The below steps shows how the installation is carried out on Windows Operating System.

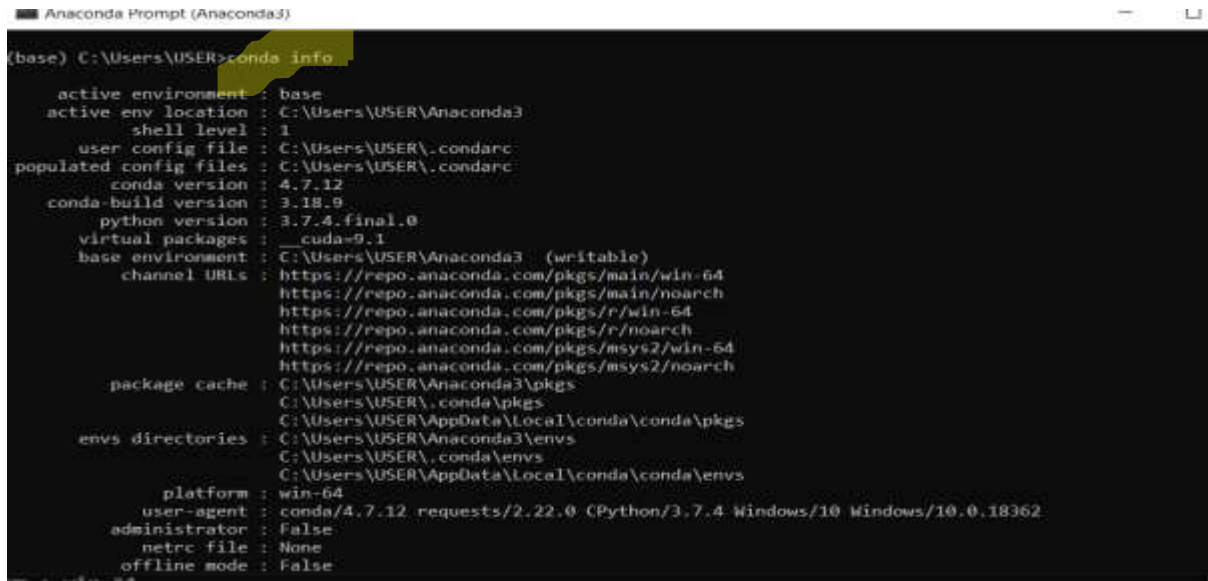
1 Installing Anaconda3

This section shows installation of Anaconda on an operating system of windows version 10. Anaconda3 is a designed package manager which is an open source for actualizing data analysis and machine learning project.

Step 1: Download the latest version of Anaconda installer with <https://anaconda.com/distribution>



Step 2: After successful download, check the file integrity with the highlighted command to be sure that the download is right.



```

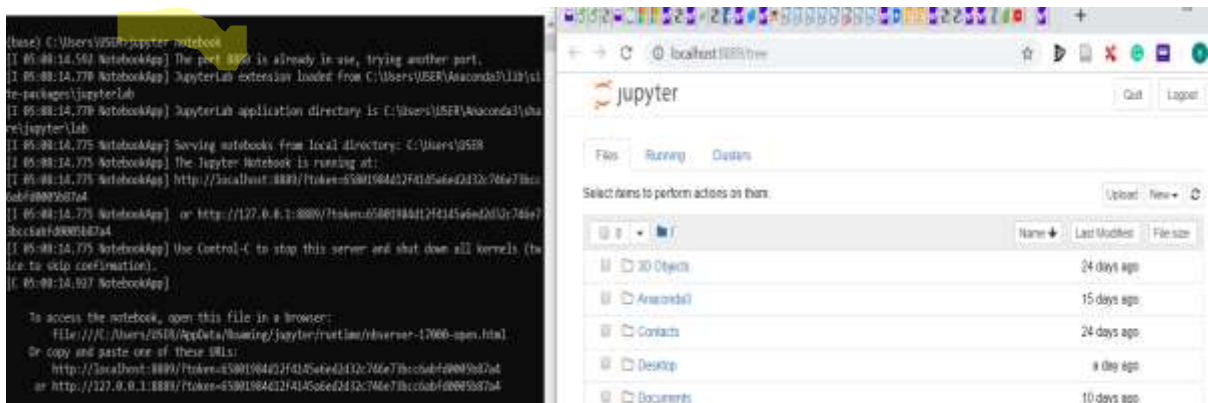
(base) C:\Users\USER>conda info

active environment : base
active env location : C:\Users\USER\Anaconda3
shell level        : 1
user config file   : C:\Users\USER\.condarc
populated config files : C:\Users\USER\.condarc
conda version      : 4.7.12
conda-build version: 3.18.9
python version     : 3.7.4.final.0
virtual packages   : _cuda=9.1
base environment   : C:\Users\USER\Anaconda3 (writable)
channel URLs       : https://repo.anaconda.com/pkgs/main/win-64
                   https://repo.anaconda.com/pkgs/main/noarch
                   https://repo.anaconda.com/pkgs/r/win-64
                   https://repo.anaconda.com/pkgs/r/noarch
                   https://repo.anaconda.com/pkgs/msys2/win-64
                   https://repo.anaconda.com/pkgs/msys2/noarch
package cache     : C:\Users\USER\Anaconda3\pkgs
                   C:\Users\USER\.conda\pkgs
                   C:\Users\USER\AppData\Local\conda\conda\pkgs
envs directories  : C:\Users\USER\Anaconda3\envs
                   C:\Users\USER\.conda\envs
                   C:\Users\USER\AppData\Local\conda\conda\envs
platform          : win-64
user-agent        : conda/4.7.12 requests/2.22.0 CPython/3.7.4 Windows/10 Windows/10.0.18362
administrator    : False
netrc file        : None
offline mode      : False

```

2 Jupyter Notebook IDE Installation using pip

Step 1: From the Anaconda prompt, run Jupyter notebook



```

(base) C:\Users\USER>jupyter notebook
[I 05-08-14.58] NotebookApp] The port 8888 is already in use, trying another port.
[I 05-08-14.778] NotebookApp] JupyterLab extension loaded from C:\Users\USER\Anaconda3\lib\site-packages\jupyterlab
[I 05-08-14.779] NotebookApp] JupyterLab application directory is C:\Users\USER\Anaconda3\share\jupyterlab
[I 05-08-14.775] NotebookApp] Serving notebooks from local directory: C:\Users\USER
[I 05-08-14.775] NotebookApp] The Jupyter Notebook is running at:
[I 05-08-14.775] NotebookApp] http://localhost:8888/?token=65801984d2f4d4566d2d32c766e73cc5abf09095b7d4
[I 05-08-14.775] NotebookApp] or http://127.0.0.1:8888/?token=65801984d2f4d4566d2d32c766e73cc5abf09095b7d4
[I 05-08-14.775] NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).

To access the notebook, open this file in a browser:
File:///C:/Users/USER/AppData/Local/Temp/jupyter-runtime/runtime-17066-open.html
Or copy and paste one of these URLs:
http://localhost:8888/?token=65801984d2f4d4566d2d32c766e73cc5abf09095b7d4
or http://127.0.0.1:8888/?token=65801984d2f4d4566d2d32c766e73cc5abf09095b7d4

```

3 Importing and Extraction of dataset

Step 1: Import necessary python libraries for dataset preprocessing

```
# Necessary Python Libraries
import pprint, os, re, email, csv
from bs4 import BeautifulSoup
from collections import Counter
from urllib.parse import urlparse
from IPY import IP
```

Step 2: Download dataset from the site as follows; (<https://monkey.org/~jose/phishing/>) and Spam Dataset (<https://spamassassin.apache.org/old/publiccorpus/>) and get path to spam, easy_ham, test_emails, and phishing from the Dataset imported.

```
spam_path = "C:\\Users\\USER\\Desktop\\phisbam\\phish-spamming-detection\\datasets\\spam"
easy_ham_path = "C:\\Users\\USER\\Desktop\\phisbam\\phish-spamming-detection\\datasets\\easy_ham"
test_path = "C:\\Users\\USER\\Desktop\\phisbam\\phish-spamming-detection\\datasets\\test_emails"
phishing_2017_path = "C:\\Users\\USER\\Desktop\\phisbam\\phish-spamming-detection\\datasets\\phishing\\phishing_2017"
phishing_2018_path = "C:\\Users\\USER\\Desktop\\phisbam\\phish-spamming-detection\\datasets\\phishing\\phishing_2018"
```

Step 3: Extracting all the 40 features from the dataset

```
def overall_feature_extraction(path, label, mail):
    necessary_fields = extract_necessary_fields(path, mail)
```

Step 4: Extract features of all mails in all paths and create a csv with the target attribute 'labels'

```
# Paths and corresponding labels
all_paths_labels = {spam_path : "Spam",
                    easy_ham_path : "Ham",
                    phishing_2017_path : "Phishing",
                    phishing_2018_path : "Phishing"}
pprint.pprint(all_paths_labels, width = 1)

# Datasets used in the paper
datasets = ["final_dataset_HamSpam", "final_dataset_HamPhishing", "final_dataset_HamSpamPhishing"]
```

Step 5: After successful feature extraction and dataset creation

```
create_features_csv: success
extract_features_into_csv: success
dataset_creation: success
```

4 Analysis of dataset

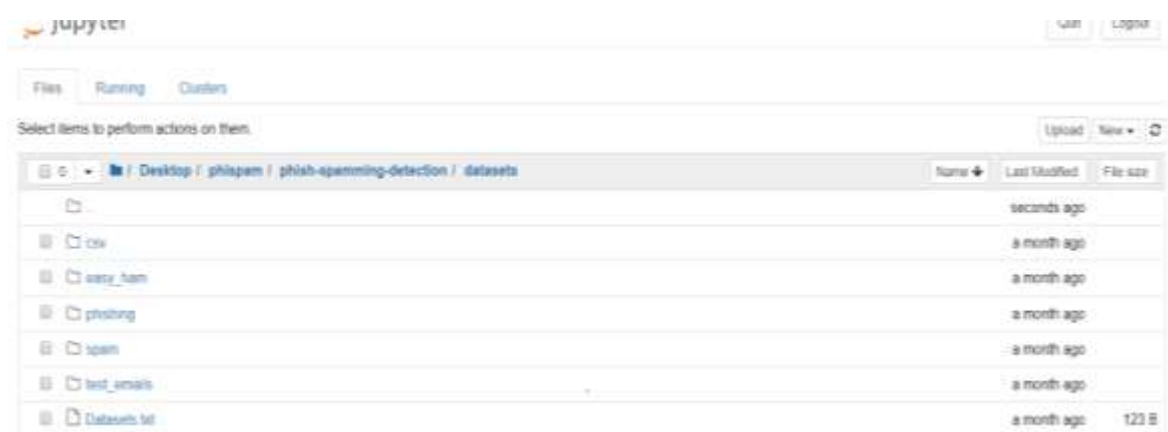
Step 1: The below codes show the imported libraries in Jupyter notebook for analysis

```
In [1]: %matplotlib inline

import matplotlib.pyplot as plt
plt.style.use('ggplot')

import numpy as np
import os, math, sys
from collections import Counter
from IPython.display import display
import subprocess
import pandas as pd
from itertools import groupby
from operator import itemgetter
import timeit
```

Step 2: Insert easy_ham, phishing, spam and test_email files in the same folder, run the preprocessing files to download three csv files; which are test, final_dataset_HamPhishing, final_dataset_HamSpam, and final_dataset_HamSpamPhishing.



Step 3: Find below the view of our preprocessing and data analysis files in Jupyter notebook.



5 Read File and Sklearn Configuration

Step 1: Specify the directory that files are located and read **final_dataset_HamSpamPhishing** file only

```
csv_directory = "C:\\Users\\USER\\Desktop\\phispam\\phish-spamming-detection\\datasets\\csv\\"
files = ['final_dataset_HamPhishing', 'final_dataset_HamSpam', 'final_dataset_HamSpamPhishing']
```

Step 2: Additional command to read head and tail of final dataset

```
global df
def get_data(filename):
    filename = csv_directory + filename + '.csv'
    df = pd.read_csv(filename)
    return df

filename = files[2]
df = get_data(filename)
display(df.head())
display(df.tail())
```

Step 3: Import all the algorithm from Sklearn

```
from sklearn import model_selection

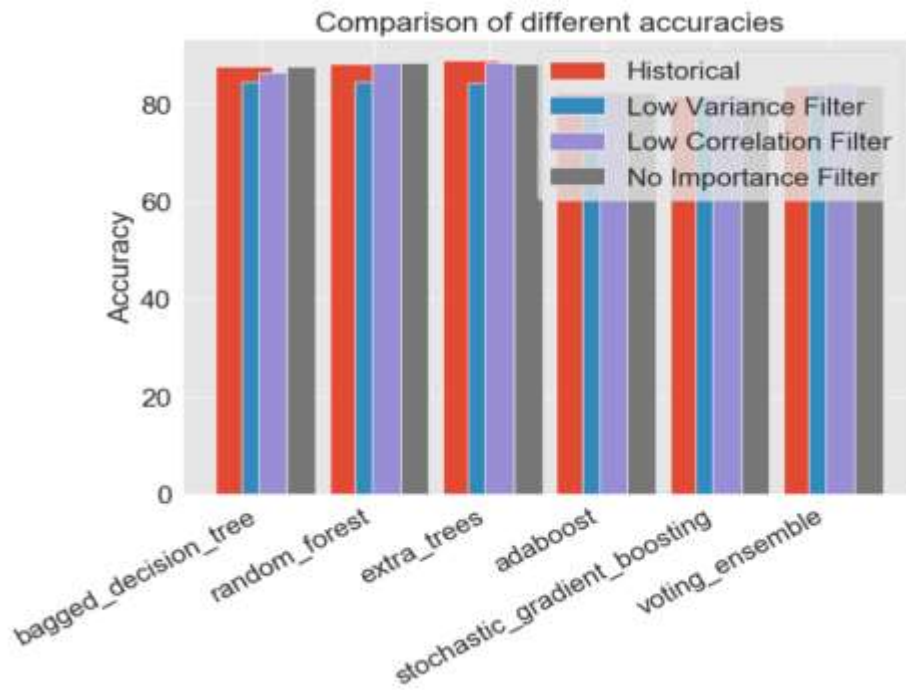
from sklearn.ensemble import BaggingClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import ExtraTreesClassifier

from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import GradientBoostingClassifier

from sklearn.svm import SVC
from sklearn.ensemble import VotingClassifier

from sklearn import metrics
from sklearn.naive_bayes import GaussianNB
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
```

Step 4: Knowledge chart of the classification model



References

- [1] <https://docs.anaconda.com/anaconda-enterprise-4/ae-and-nav/>
- [2] <https://scikit-learn.org/stable/install.html#installing-the-latest-release>