# Profiling internet users CIS6930 Information Security and Privacy

Jaynab Khatun

UID- U06687764

## Introduction:

The objective of this project is to find out at what time window a user is statically indistinguishable from the number of other user's data across. There are three time window 10s, 227s and 5 minutes and total 54 users.

## Source of data:

The source of data for this project is Cisco NetFlow version 5. Many parameters has been provided like Packet, Octets, Real First Packet, Real End Packet, duration, Doctets/packet and many other variables. There are 54 data sheet for each user for a month long period.

## Requirement:

1. Windows, 64 bit machine
2. Anaconda, Spider to run python3.6
3. Data for 54 users

## Summary of contribution:

In this project, 54 user's data in excel sheet are collected, then modify it with parameter Doctet, Real First Packet, Duration. For ease to handle excel sheet other column has discarded as those are not needed. At same time, calculated the Doctet/Duration and ignore the duration value where it is zero to handle division zero. I have collected epoch data for $1^{st}$ February, $8^{th}$ February from internet which help to find time window calculation in program. Write program in python to divide the time window for 10s, 227 s and 5 minutes which define as 300s, the program name as first.py. While dividing time widow, count the frequency as well so that it will help to find out the total number window in week for particular user. In this way find out all data for week1, week2 and write excel sheet which named as New1 for 10s window, New2 for 227s window, New3 for 300s window.
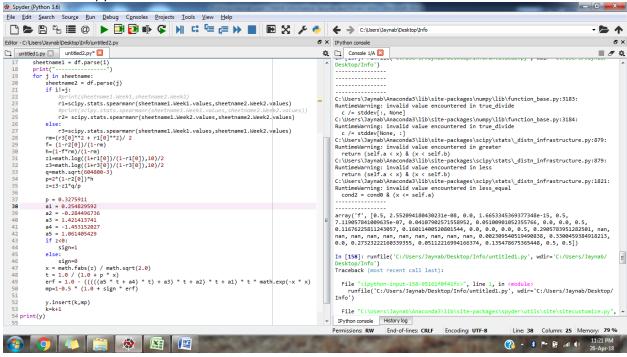
Now the second part of project is to find out Spearman Correlation coefficient and MRR-Test which is known as Z value,I have written another program in python named as second.py. I have used "spearmanr" function which is predefined in Python3.6 for each user and consider week1, week2. At the same time calculate the Z- value and P value for correlation coefficient of user and corresponding week. The collected all P value for each user and put it in a list of array. The image show
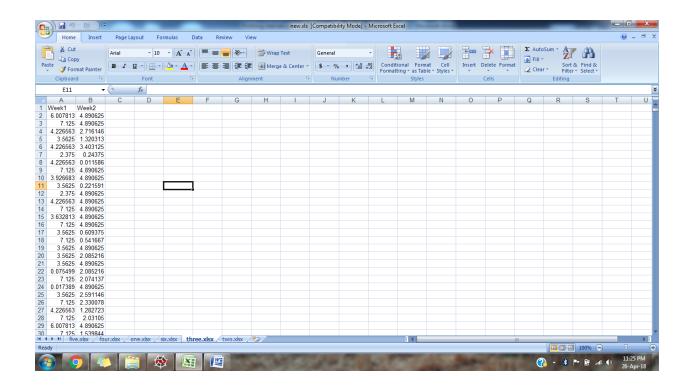
## How to run in Python:

1. In 54 excel sheet, Doctet, Real First packet and Duration is selected and calculated the Doctet/duration. Name as one.xlsx, two.xlsx, three xlsx like this till 54 users.
2. Keep all modified excel sheet in same work folder where python program is created
3. Read all excel sheet and produce the new1, new2 and new3 excel sheet using python program.

4. Now take new1, new2, new3 excel sheet input and find the spearman correlation coefficient and corresponding Z-value, P-value using program second.py. The program is already written in python and read excel sheet. User needs to run to python program separately one by one(first run First.py the run second.py)

The below screen shot has taken for five user and print their corresponding P value using second.py



Below screen shot is for New1 excel sheet created by executing the first.py program

## Analysis:

Form the Table of P value, it is observe that for same user in same week and time window, the P value is 0.5.  Look out for number of average P value in P value matrix for each time window and find the 10s window:

| Week1 & Week2 | One | Two | Three | Four | five |
|---|---|---|---|---|---|
| One | 0.5 | .0255209 | 0.0 | 0.166533 e-15 | 0.0711905 |
| Two | 0.5 | 0.5 | 0.0410790 | 0.0510 | 0.0 |
| Three | 0.2715938 | 0.1167 | 0.5 | 0.16011 | 0.0 |
| Four | 0.5 | 0.0290578 | 0.00230954 | 0.5 | 0.3300459 |
| five | 00 | 0.2732 | 0.511 | 0.5 | 0.5 |

## Conclusion: 
The project has done in Python3.6 and result has shown in the table with P value for each user. Most of P values are greater that  0.05.

## Reference:

[1] Soheil Sarmadi (Ph D Student of University of South Florida)

[2] Sriram Chellappan(Assistant professor of  University of South Florida)