# CSCE-638, Programming Assignment #1 SpamLord

Li-Hen Chen 928003907

**1. How to compile and run your code.**

We program the code in Python3 and execute the program under the same directory of the SpamLord.py file.


Eg.

In the terminal just type:


python SpamLord.py data dev/dev/ data dev/devGOLD


**2. Results and Analysis**

For extracting email, although there are some simple cases like "huangrh@cse.tamu.edu" or "huangrh at cse dot tamu dot edu", there are some complex cases such as "<script type="text/javascript">obfuscate('cse.tamu.edu','huangrh')</script>". For extracting phone number, there are some cases we have to deal with such as "(979) 862-2908 ", "979-862-2908" and " 979 862 2908".


- Email Extraction

We should find out the possible email expression first in order to solve the problem and below is the basic email format I found.

    <td class="value">ouster (followed by &ldquo;@cs.stanford.edu&rdquo;)</td>

    <td class="value">teresa.lynn (followed by "@stanford.edu")</td>

    <dd>        <em>melissa&#x40;graphics.stanford.edu</em>

    <address>engler WHERE stanford DOM edu</address>

    email: pal at cs stanford edu,

    d-l-w-h-@-s-t-a-n-f-o-r-d-.-e-d-u

    <dd>        <em>ada&#x40;graphics.stanford.edu</em>

    Email: uma at cs dot stanford dot edu

    hager at cs dot jhu dot edu

    (Fedora) Server at cs.stanford.edu Port 80

    <script> obfuscate('stanford.edu','jurafsky'); </script>

Except for the last example, we can use regular expression to divide each part. First part is characters previous "@". There may by dash between characters so we just use a regular expression to match it. In the second part, there are some cases with "followed by" so we also use ? to filter whether the sentence contains "followed by". In third part, there are some cases "@, @ , at , where ,&#x40;," so we decide to set different regular expression to judge which one is contained in the sentence. For the remaining parts we use same idea to filter the cases like "., ;, dot , dt , DOM" and ".edu, .com, -e-d-u". For the last example "obfuscate('stanford.edu','jurafky');", we just have to set a special regular expression to deal with it.

- Phone Number Extraction

Compared to the email address, the phone number patter is relatively easy. At first part, we set a regular expression to handle to "(123)" and "123". Then we just have to deal with whether the remaining part contains dash and replace it to the normal phone number.

- Test results on dev dataset

True Positive: 59

False Negative: 0

False Positive: 0

**3. Any known bugs, problems, or limitations of your program**

Although we did perfect on dev dataset, I did not handle the cases such as
" 979 862 2908". As results, there may be some errors in the test
dataset. I believe the examples I listed have been dealt with perfectly.