# PYTHON CRAWLER

From Beginner To Intermediate

# Self Introduction

## Cheng-Yi, Yu

erinus.startup@gmail.com

- **LifePlus Inc.**

  Technical Manager

- **Paganini Plus Inc.**

  Senior Software Developer

- **Freelancer**

  ~ 10 Years

# DAY 1

# Python

- **Entry**

  if \_\_name\_\_ == '\_\_main\_\_':

      # do something

- **Method**

  def main ():

      # do something

- **Package**

  import [package]

  import [package] as [alias]

- **Format**

  '%s' % ([parameters ...])

demo01.py

demo02.py

demo03.py

# Python

- **If ... Else ...**

  if [condition]:

      # do something

  else:

      # do something

- **For Loop**

  for item in list:

      # do something

- **Array Slice**

  array[start:end]

# Python

- **Array Creation From For Loop**

  – **Object**

  [item.attr for item in array]

  – **Dictionary**

  [item[key] for item in array]

# Python

- **In**
  - **String**

    if <str> in <str>:

  - **Array**

    if item in list:

  - **Dictionary**

    if <str:key> in <dict>:

# Installation

- **Ubuntu Fonts**

  https://design.ubuntu.com/font/

- **Source Han Sans Fonts**

  https://github.com/adobe-fonts/source-han-sans

- **Visual Studio Code And Extensions**

  https://code.visualstudio.com/

# Installation

- ## Cmder
  http://cmder.net/

- ## Python 3.6
  https://www.python.org/

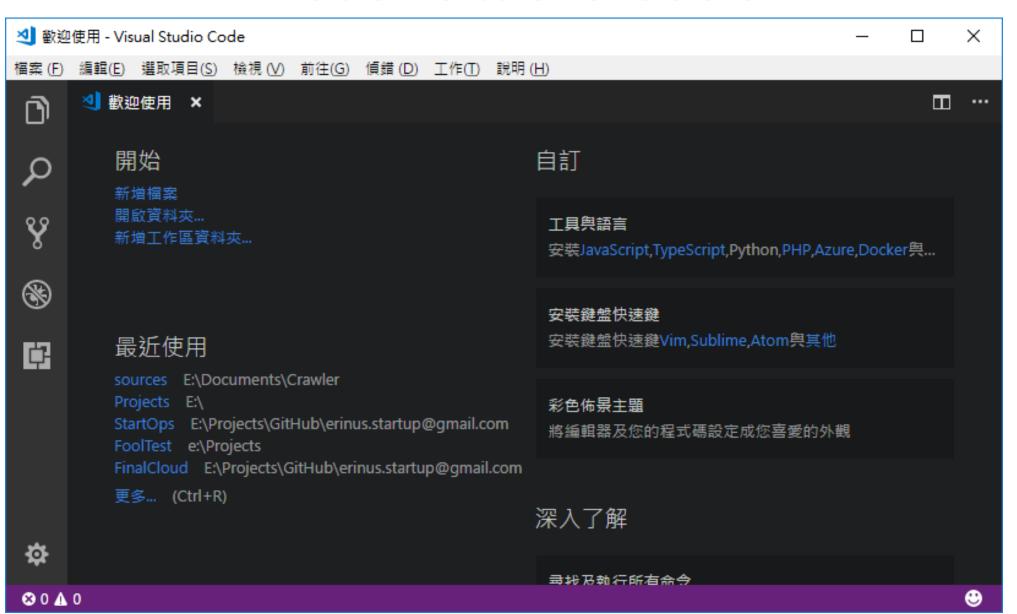- ## Python Packages
  pip install requests

  pip install pyquery
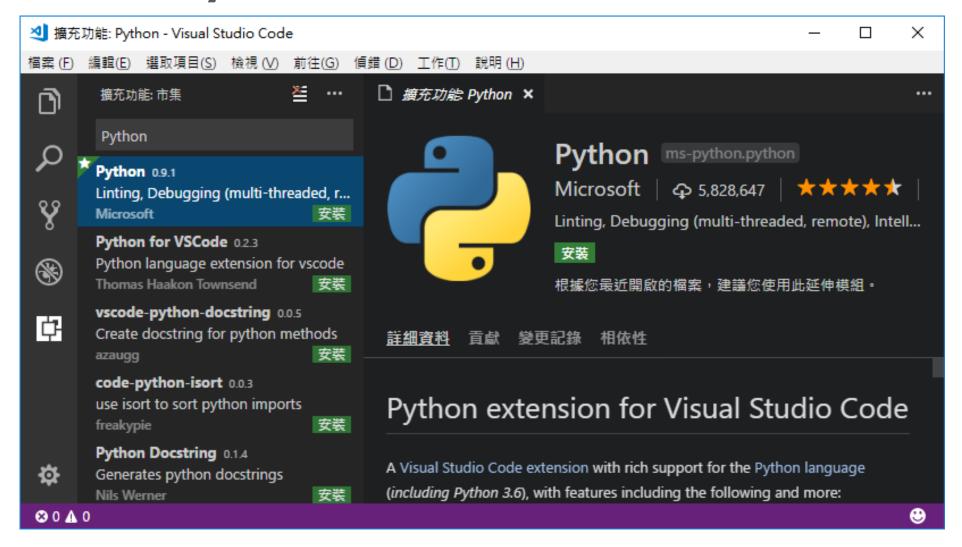
  pip install beautifulsoup4
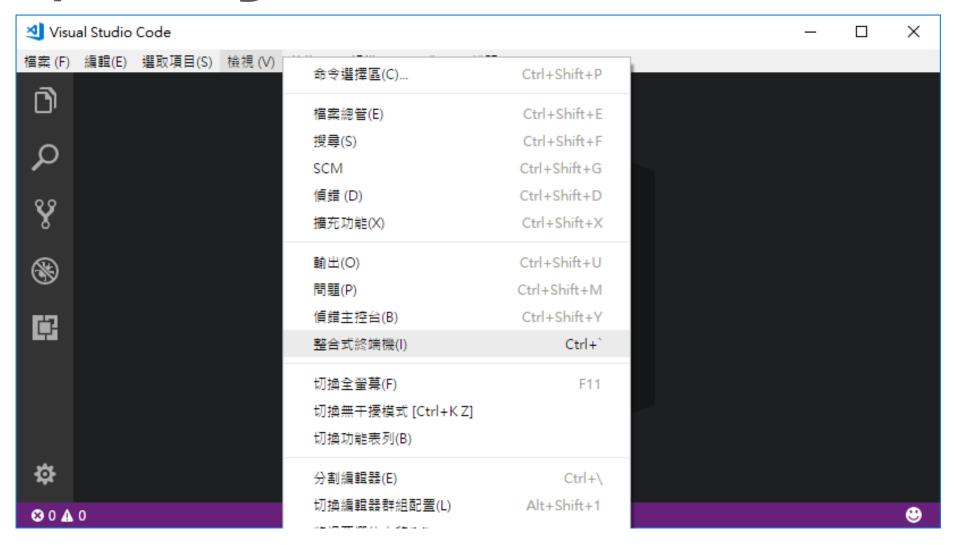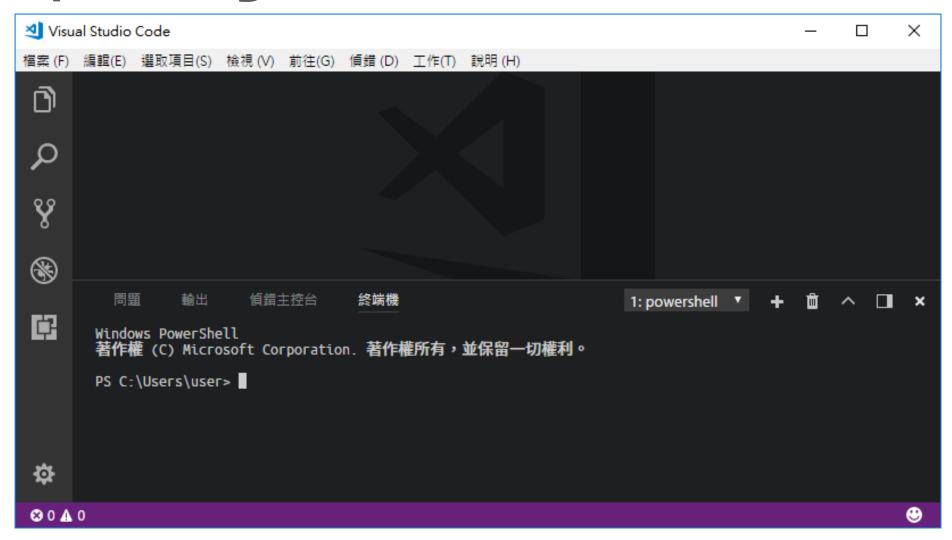
  pip install js2py

  pip install selenium

# Visual Studio Code

# Visual Studio Code

- **Install Python Extensions**

# Visual Studio Code

- **Open Integrated Terminal**

# Visual Studio Code

- **Open Integrated Terminal**

# Built-in

- **Json**

  import json

  - **json.loads**

    json.loads(&lt;str&gt;)

  - **json.dumps**

    json.dumps(&lt;dict&gt;)

demo05.py

demo06.py

# Built-in

- **Xml**

  import xml.etree.ElementTree as ET

  - **Load From File**

    tree = ET.ElementTree(file=<str:filepath>)

    tree = ET.parse(<str:filepath>)

    root = tree.getroot()

  - **Load From String**

    root = ET.fromstring(<str>)

# Built-in

- **Xml**
  - **Child Nodes**

    **Only One Level**

    for node in root:

        # do something

  - **XPath**

    **Multiple Levels**

    nodes = root.findall(&lt;str:expression&gt;)

    for node in nodes:

        # do something

demo07.py

demo08.py

# Built-in

- ## Url

  import urllib.parse as UP

  - ### urlparse

    result = UP.urlparse(<str:url>)

  - ### urlunparse

    url = UP.urlunparse(<ParseResult>)

  - ### quote

    str = UP.quote(<str:unquoted>)

  - ### unquote

    str = UP.unquote(<str:quoted>)

demo09.py

demo10.py

# Built-in

- ## Regular Expression

import re

- ### re.search

  **Find First Match**

  match = re.search(<str:pattern>, <str:text>)

  match.group(<int:index>)

- ### re.findall

  **Find All Matches**

  finds = re.findall(<str:pattern>, <str:text>)

  for find in finds:

     # do something

# Built-in

- **Regular Expression**

  – **re.split**

    **Split By Pattern**

    re.split(<str:pattern>, <str:text>)

  – **re.sub**

    **Replace By Pattern**

    re.sub(<str:pattern>, <str:replace>, <str:text>)

# Built-in

- **Regular Expression**
  - **Expressions**
    1. **Range**

       [Start-End]

       [0-9], [a-z], [A-Z], [a-zA-Z], [0-9a-zA-Z], ...

# Built-in

- **Regular Expression**
  - **Expressions**
    1. **Zero Or More Times**

       **\***

    2. **One Or More Times**

       **+**

    3. **Zero Or One Time**

       **?**

# Built-in

- **Regular Expression**
  - **Expressions**
    1. **Numbers**

       **\d** = [0-9]

    2. **Words**

       **\w** = [a-zA-Z0-9] (ANSI)

       **\w** = [a-zA-Z0-9] + Non-ANSI Characters (UTF-8)

    3. **Spaces, Tabs, ...**

       **\s**

# Built-in

- **Regular Expression**

  – **Expressions**

    1. **Start With**

       **^**

    2. **End With**

       **$**

# Built-in

- ## **Regular Expression**
  - ### **Expressions**
    1. **Named Group**

       **(?P<name>expr)**

       (?P<country>\+\d+)-(?P<phone>\d+)

# AnalySIS

- **Chrome Developer Tools**
  - **Elements**

    **See Elements In DOM**

    Id, Class, Attribute, ...
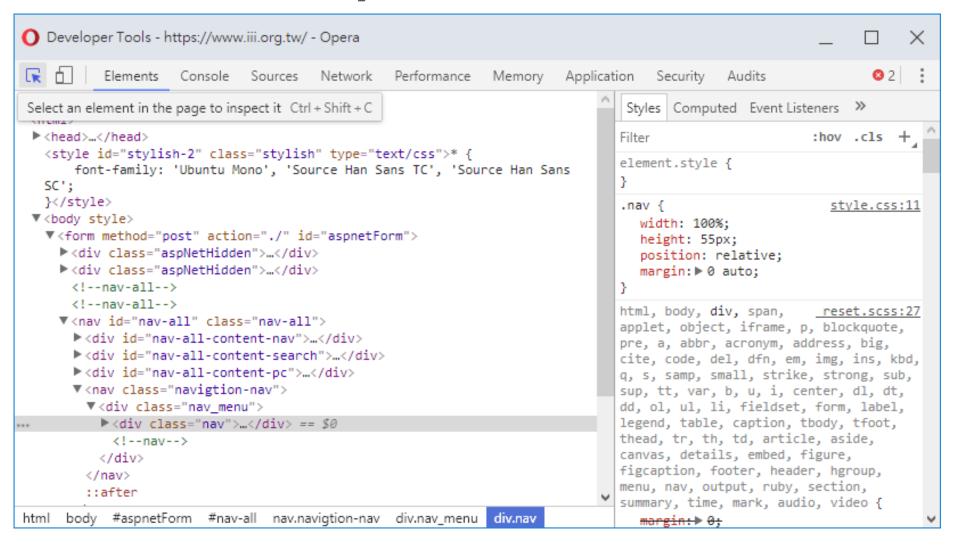
  - **Network**

    **See Requests, Responses**

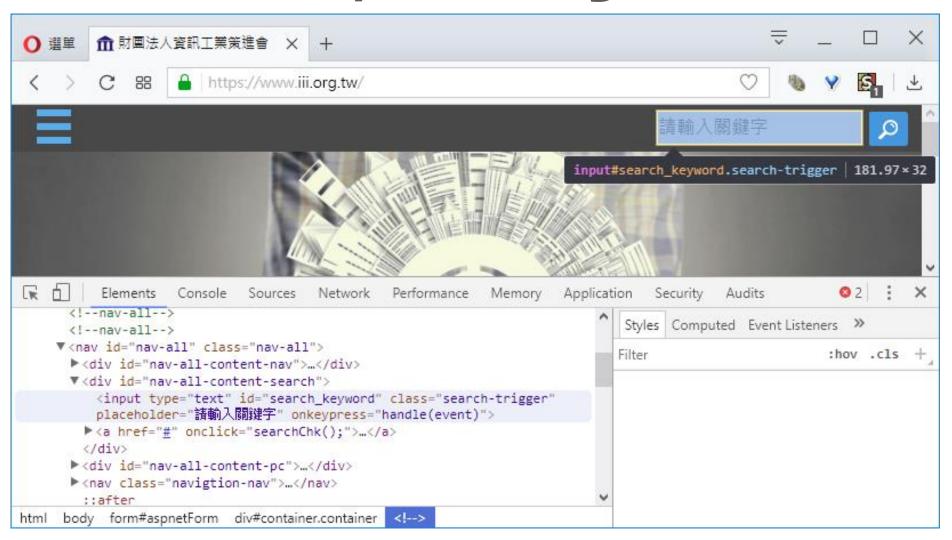    Urls, Methods, Headers, Cookies, Bodies, ...
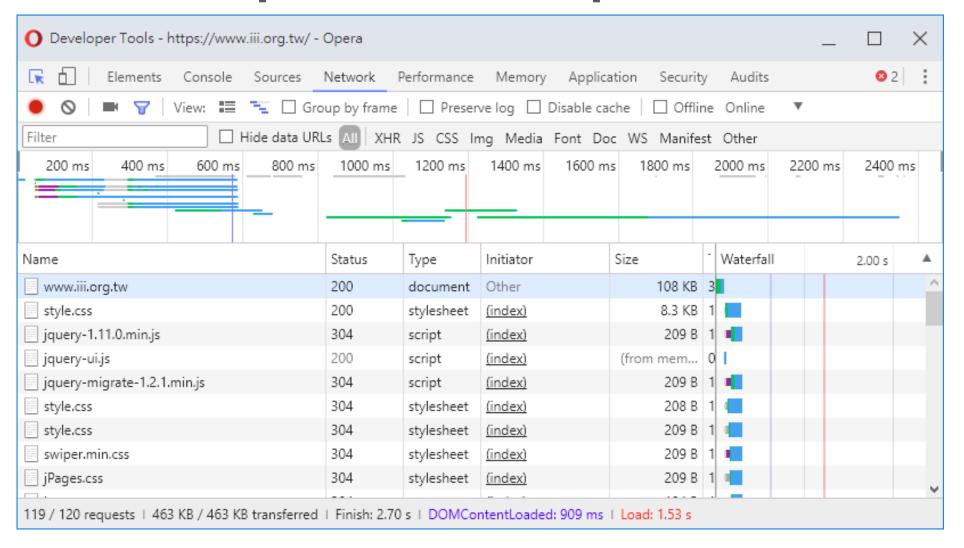
# Elements

- ## Find Element by Mouse Pointer
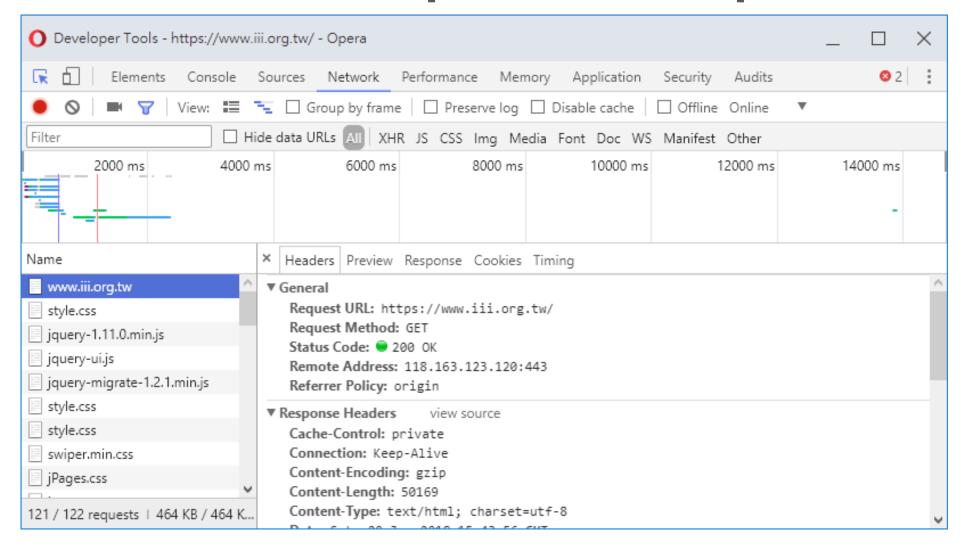
# Elements

- **Find Element by HTML Tag**

# Networks

- **See All Requests And Responses**

# Networks

- **See Details Of Request And Response**

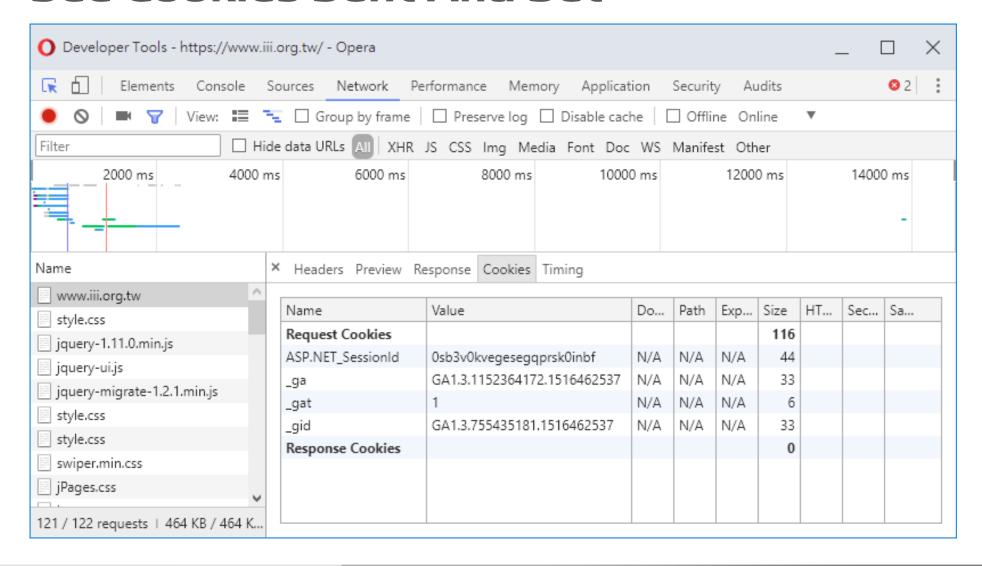# Networks

- **See Response Content**

# Networks

- **See Cookies Sent And Set**

# Documents

- **Requests**

  http://docs.python-requests.org/

- **PyQuery**

  https://pythonhosted.org/pyquery/

- **Beautiful Soup 4**

  https://www.crummy.com/software/BeautifulSoup/bs4/doc/

- **Js2Py**

  https://github.com/PiotrDabkowski/Js2Py

# Packages

- ## Requests

import requests

- ## Request

  1. ## Method

  **GET, POST, ...**

  response = requests.get(\<str:url\>)

  response = requests.post(\<str:url\>, data=\<str:body\>)

  response = requests.post(\<str:url\>, data=\<dict:body\>)

  2. ## Session

  session = requests.Session()

  response = session.get(\<str:url\>)

# Packages

- **Requests**

  import requests

  - **Request**

    - **Headers**

      response = requests.get(<str:url>, headers=<dict>)

    - **Cookies**

      response = requests.get(<str:url>, cookies=<dict>)

# Packages

- **Requests**
  - **Response**
    1. **Status Code**

       response.status_code

    2. **Headers**

       response.headers[<str:name>]

    3. **Cookies**

       response.cookies[<str:name>]

# Packages

- # Requests

  demo14.py

  - ## Response

    1. **Binary Content**

       response.content

    2. **Text Content**

       response.text

    3. **Json Content**

       response.json()

# Packages

- ## PyQuery

import pyquery

- ### Load From String

  d = pyquery.PyQuery(<str:html>)

- ### Load From Url

  d = pyquery.PyQuery(url=<str:url>)

- ### Load From File

  d = pyquery.PyQuery(filename=<str:filepath>)

# Packages

- **PyQuery**
  - **Find**

    p = d(\<str:expression>)

  - **Element To HTML**

    p.html()

  - **Extract Text From Element**

    p.text()

  - **Get Value From Element Attribute**

    val = p.attr[\<str:name>]

# Packages

- ## Beautiful Soup 4

import bs4

  – ### Load From String

  d = bs4.BeautifulSoup(<str:html>, 'html.parser')

# Packages

- **Beautiful Soup 4**
  - **Find**

    p = d.find_all(\<str:tag\>, \<attr-key\>=\<attr-val\>, …)

    p = d.find_all(\<regex\>, \<attr-key\>=\<attr-val\>, …)

    p = d.find_all(\<array\>, \<attr-key\>=\<attr-val\>, …)

    p = d.find(\<str:tag\>, \<attr-key\>=\<attr-val\>, …)

    p = d.find(\<regex\>, \<attr-key\>=\<attr-val\>, …)

    p = d.find(\<array\>, \<attr-key\>=\<attr-val\>, …)

    p = d.select(\<str:expression\>)

    p = d.select_one(\<str:expression\>)

demo16.py

demo17.py

# Packages

- **Beautiful Soup 4**
  - **Extract Text From Element**

    p.get_text()

  - **Get Value From Element Attribute**

    p.get(<str:name>)

demo16.py

demo17.py

# Packages

- **Js2Py**

import js2py

  – **Eval**

    js2py.eval_js(<str:code>)

    res = js2py.eval_js('var o = <str:js>; o')

# DAY 2

# WORKSHOP

- **Apple Daily**

  https://tw.appledaily.com/

  – **Realtime News**

    https://tw.appledaily.com/new/realtime

- **Facebook Page**
  - **Cookies**
  - **Feed**

# DAY 3

# SELENIUM

- **Download ChromeDriver**

  https://sites.google.com/a/chromium.org/chromedriver/

# SELENIUM

- ## Import

  import selenium.webdriver

- ## Initialize

  option = elenium.webdriver.ChromeOptions()

- ## Start

  driver = selenium.webdriver.Chrome(chrome_options=option)

- ## Browse

  driver.get(<str:url>)

- ## Close

  driver.quit()

# SELENIUM

- **Source**

  driver.page_source

# SELENIUM

- **Find One**
  - **find_element_by_id**
  - **find_element_by_name**
  - **find_element_by_tag_name**
  - **find_element_by_class_name**
  - **find_element_by_css_selector**

# SELENIUM

- **Find Multiple**
    - **find_elements_by_name**
    - **find_elements_by_tag_name**
    - **find_elements_by_class_name**
    - **find_elements_by_css_selector**

# SELENIUM

- **Actions**
  - **send_keys**
  - **click**

- **Facebook Page**
  - **Login**
  - **Feed**