

# Python & Data - Week 9

打包下載

## 堂上作業

一、目的：請找一個網站嘗試做 Web Scraping。運用 Week 8 提過的方法去判斷是否適合。

二、形態：二人一組完成

三、提交（展示）內容

- 1. 網址
- 2. 具有網頁頁面、網址
- 3. 希望截取的內容及其 HTML 片段
- 4. 用途四、提交日期：未定

## 本期內容

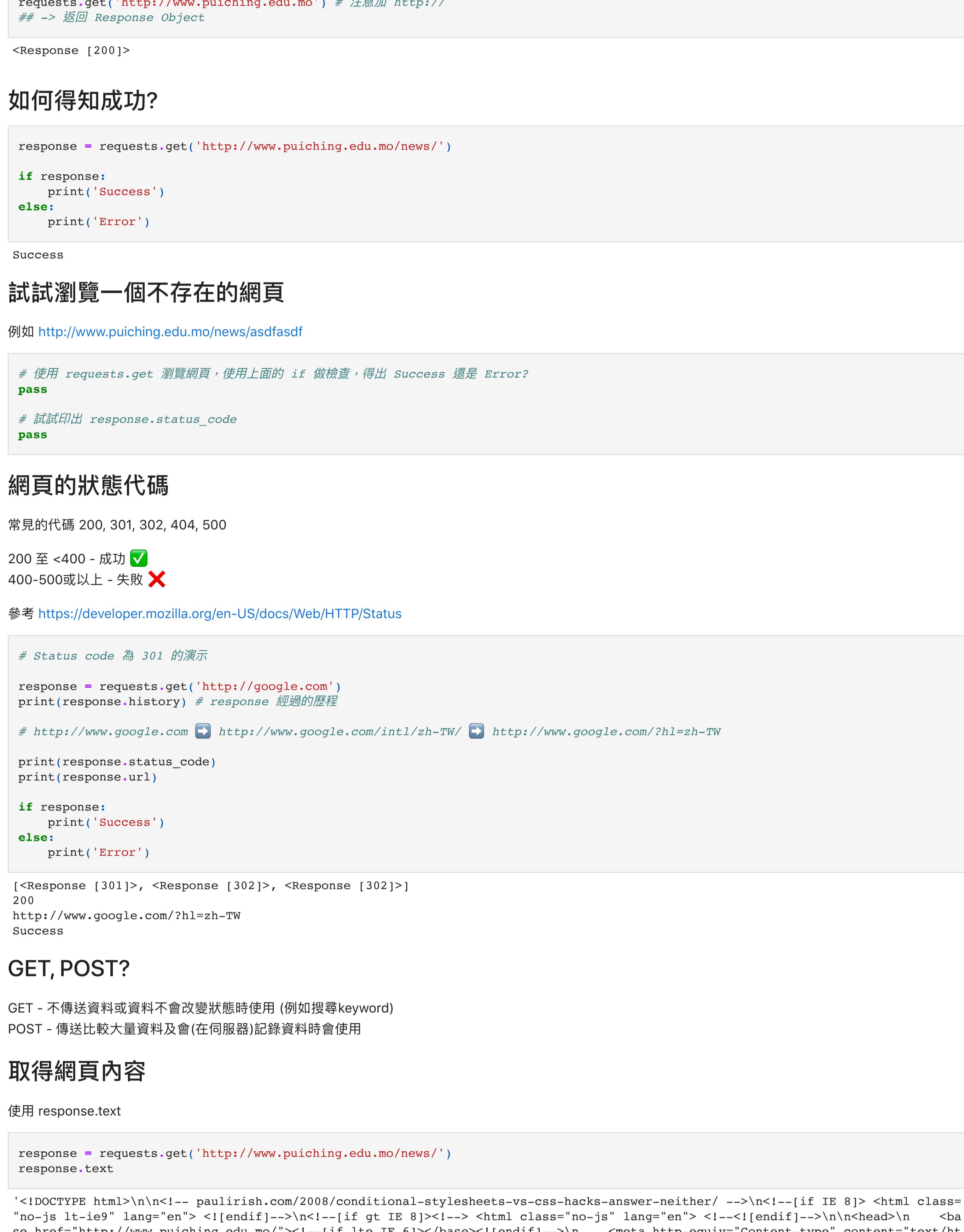
取得內容（瀏覽網頁）進行處理（解析內容）

## 主要可使用兩種技術

1. Requests + Beautiful Soup
  - 用於靜態產生(HTML)的網頁
2. Selenium + Web Driver
  - 用於動態產生(JavaScript)的網頁

## 靜態產生(Static) vs 動態產生(Dynamic)網頁

演示



## Requests 和 BeautifulSoup 是甚麼

1. Requests 取得原始內容
2. BeautifulSoup 進一步解析內容

## 使用 Requests

```
In [1]: # 安裝 Requests
import requests
```

```
In [2]: pip install requests
```

```
DEPRECATION: Configuring installation scheme with distutils config files is deprecated and will no longer work in the near future
If you are using a Homebrew or Linuxbrew Python, please see discussion at https://github.com/Homebrew/homebrew-core/issues/7662
```

```
Requirement already satisfied: requests in /usr/local/lib/python3.9/site-packages (2.27.1)
Requirement already satisfied: charset-normalizer<2.0.0,!=2.0.0rc1,> in /usr/local/lib/python3.9/site-packages (from requests) (2.0.12)
Requirement already satisfied: certifi<2021.10.6,!=2021.10.6rc1,> in /usr/local/lib/python3.9/site-packages (from requests) (2021.10.6)
Requirement already satisfied: idna<3.2.0,!=3.2.0rc1,> in /usr/local/lib/python3.9/site-packages (from requests) (1.26.9)
Requirement already satisfied: ieha<2.0.0,!=2.0.0rc1,> in /usr/local/lib/python3.9/site-packages (from requests) (2.0.0)
WARNING: You are using pip version 21.2.4; however, version 22.0.4 is available.
You should consider upgrading via https://user.local/opt/python3.9/bin/python3.9 -m pip install --upgrade pip.
```

```
Note: you may need to restart the kernel to use updated packages.
```

## 使用 requests "瀏覽" 網頁

```
In [3]: import requests
response = requests.get('http://www.puiching.edu.mo') # 注意加 http://
```

```
#--> 返回 Response Object
<Response [200]>
```

```
Out[3]: Success
```

## 試試瀏覽一個不存在的網頁

例如 <http://www.puiching.edu.mo/news/asdfasdf>

```
In [4]: # 使用 requests.get 激覽網頁，使用上面的 if 做檢查，得出 Success 這是 Error?
pass
```

```
# 請印出 response.status_code
pass
```

```
In [5]: # Status code 為 301 的演示
response = requests.get('http://www.google.com')
print(response.history) # response 經過的歷程
```

```
# http://www.google.com  http://www.google.com/intl/zh-TW  http://www.google.com/?hl=zh-TW
```

```
print(response.status_code)
```

```
if response:
    print('Success')
else:
    print('Error')
```

```
Success
```

## GET, POST?

GET - 不傳送資料或資料不會改變狀態時使用 (例如搜尋keyword)

POST - 傳送比較大量資料及會在伺服器記錄資料時會使用

## 取得網頁內容

使用 response.text

```
In [6]: response = requests.get('http://www.puiching.edu.mo/news/')
```

```
response.text
```

```
Out[6]: <no-dt-type ie9> lang=<nl>-- paulirish.com/2008/conditional-style-sheets-vs-css-hacks-answer-neither--><n!--[if IE 8]> lang=<nl>--><n!--[if !ie9]> lang=<nl>--><n!-- Set the viewport width to device width for mobile --> n <meta name="viewport" content="width=100%" /> n <meta property="og:title" content="最新消息" /> n <title>澳門培正中學 - 最新消息</title> n <!-- Included CSS Files --> n <!--link rel="stylesheet" href="themes/pcms/layout.css" type="text/css" /> n <script src="themes/pcms/javascripts/modernizr.foundation.js"></script> n <link rel="stylesheet" type="text/css" href="http://www.puiching.edu.mo/multilingual/css/langselector.css?e=1377865692" /> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/style.css" type="text/css" /> n <script src="themes/pcms/style.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footer.css" type="text/css" /> n <script src="themes/pcms/footer.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.css" type="text/css" /> n <script src="themes/pcms/footermodernizer.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.foundation.modernizr.css?e=1377865692" type="text/javascript"></script> n <!-- Included CSS File --> n <!--link rel="stylesheet" href="themes/pcms
```