Scraper with Tesseract

February 25, 2022

```
import os
import re
import requests
from bs4 import BeautifulSoup as bs
from selenium import webdriver
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.common.desired_capabilities import DesiredCapabilities
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.webdriver.common.by import By
import simplejson as json
from datetime import datetime
```

/Users/azicon/anaconda3/lib/python3.7/site-packages/requests/__init__.py:91: RequestsDependencyWarning: urllib3 (1.26.8) or chardet (3.0.4) doesn't match a supported version!

RequestsDependencyWarning)

0.1 Get website response and font data

```
except:
        driver.refresh();
now = datetime.now().strftime("%d-%m-%Y %H:%M:%S") # get exact datetime at the
\rightarrow time of scrape
os.mkdir("logs/" + now)
driver.get_screenshot_as_file("logs/" + now + "/screenshot.png") # save_
 →screenshot to sanity check later
logs_raw = driver.get_log("performance")
logs = [json.loads(lr["message"])["message"] for lr in logs_raw]
def log_filter(log_):
   return (
        # is an actual response
        log ["method"] == "Network.responseReceived"
        # and json
        and "json" in log_["params"]["response"]["mimeType"]
    )
responses = []
for log in filter(log_filter, logs):
    request_id = log["params"]["requestId"]
    resp_url = log["params"]["response"]["url"]
    print(f"Caught {resp url}")
    response = driver.execute_cdp_cmd("Network.getResponseBody", {"requestId": ___
 →request_id})
    responses.append(response)
```

 $\label{lem:caught} Caught \ https://piaofang.maoyan.com/dashboard-ajax/movie?orderType=0\&uuid=8889a5d f-571f-45aa-9e7c-803e93565323&timeStamp=1645843657776&User-Agent=TW96aWxsYS81LjAgKE1hY21udG9zaDsgSW50ZWwgTWFjIE9TIFggMTBfMTVfNykgQXBwbGVXZWJLaXQvNTM3LjM2IChLSFRNTCwgbG1rZSBHZWNrbykgQ2hyb211Lzk3LjAuNDY5Mi450SBTYWZhcmkvNTM3LjM2&index=820&channelId=40009&sVersion=2&signKey=6ddead5b4a3f722d46590a745bbc6101$

```
[3]: # Get this instance's font file from backend server
body0 = json.loads(responses[0]['body'])
movieList = body0['movieList']['list']
date = body0['calendar']['today']
font_url = body0['fontStyle'].split('"')[-2]

# Get reference fonts from the file tree
from fontTools.ttLib import TTFont
headers = {
```

```
"User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_13_4)

AppleWebKit/537.36 (KHTML, like Gecko) "

"Chrome/66.0.3359.139 Safari/537.36 "

}

woff_url = 'http:' + font_url
response_woff = requests.get(woff_url, headers=headers).content

print("Woff retrieval succuessful: " + str(len(response_woff) > 0))

with open('temp/fonts.woff', 'wb') as f:
    f.write(response_woff)
```

Woff retrieval succuessful: True

0.2 Getting digits from the font data using pytesseract

```
[4]: from fontTools.ttLib import TTFont
   from PIL import ImageFont, Image, ImageDraw, ImageOps
   import pytesseract
   import cv2
   import numpy as np
   import random
   def uniToHex(uni):
       return "&#x" + uni[3:].lower()
   def uni_2_png_stream(txt: str, font: str, img_size=512, font_size=0.7,_
     →invert=False):
       img = Image.new('1', (img_size, img_size), 255)
       draw = ImageDraw.Draw(img)
       font = ImageFont.truetype(font, int(img_size * font_size))
       txt = chr(txt)
       x, y = draw.textsize(txt, font=font)
       draw.text(((img_size - x) // 2, (img_size - y) // 2), txt, font=font,__
     \rightarrowfill=0)
        if invert:
            img = img.convert('L')
            img = ImageOps.invert(img)
            img = img.convert('1')
        #img.save(txt + '.png')
       return img
   def predict_neural(unicode, fontFile):
        image = uni_2_png_stream(int(unicode[3:], 16), fontFile, img_size=28,__
     →font_size=0.5, invert=True)
```

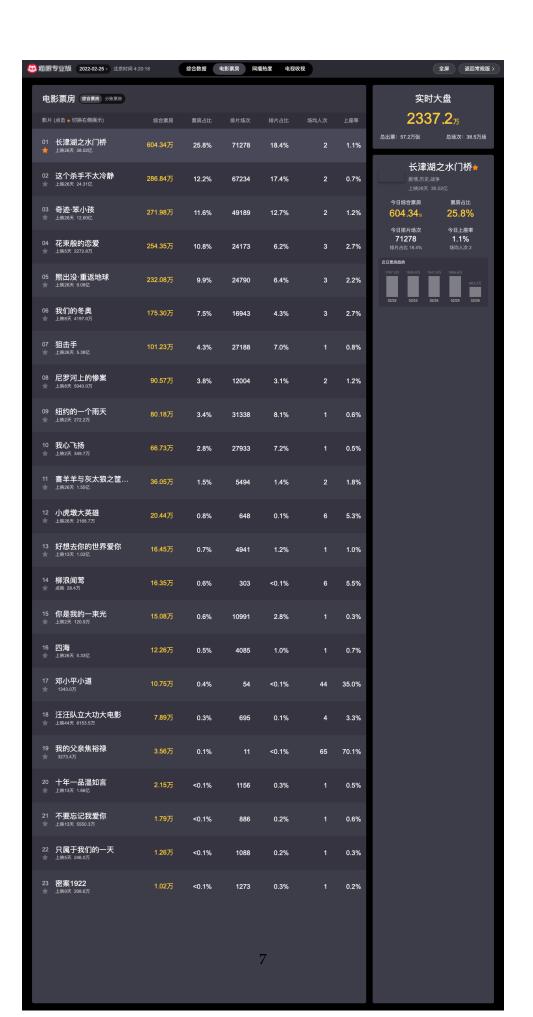
```
image.save(str(unicodeToInt[unicode]) + '_neuro.png')
       matrix_form = np.array(image)
       weighted predictions = np.ndarray.flatten(neural network.run(matrix form))
       most_possible = np.argmax(weighted_predictions)
       return most_possible
   def predict_tesseract(unicode, fontFile, fontSize=0.5):
        image = uni_2_png_stream(int(unicode[3:], 16), fontFile, img_size=1024,__
    →font size=fontSize)
        image.save('logs/' + str(now) + '/' + str(unicode) + '.png')
       text = pytesseract.image_to_string(image, lang="eng", config="--psm 10_u
     →outputbase digits -c tessedit_char_whitelist=0123456789")
       return text
   def predict_tesseract_definite(unicode, fontFile):
       result, size = '', 1
       while not result and size >= 0:
            result = predict_tesseract(x, filename, fontSize=size)
            size -= 0.01
       return result
[5]: # Map contours to numbers - the prediction phase may be very slow
   filename = 'temp/fonts.woff'
   f = TTFont(filename)
   hexToInt = {}
   for x in f.getGlyphNames()[1:-1]:
       predict = predict_tesseract_definite(x, filename)
       hexToInt[uniToHex(x)] = int(predict)
   hexToInt
   2 extra bytes in post.stringData array
[5]: {'&#xe1ae': 4,
     '&#xe63f': 1,
```

0.3 Parsing the data into pandas dataframe

```
[6]: import pandas as pd
    df = pd.DataFrame.from_records(movieList)
[7]: unitLookup = {'': 100, '': 1000, '': 1*10**8}
    #converts the weird character to a float
    def convertToFloat(string):
        spCharLst = string.split(';')
        result = ''
        for i in spCharLst:
            if len(i) > 7: #has a dot in front
                result += '.' + str(hexToInt[i[1:]])
            elif len(i) == 7: #in case of bad parsing
                result += str(hexToInt[i])
        return float(result)
    #helper function for converting the entire block to a single int
    def convertDictToInt(dictionary):
        return int(convertToFloat(dictionary['num']) *__
     →unitLookup[dictionary['unit']])
[8]: df['boxSplitUnit'] = df['boxSplitUnit'].apply(convertDictToInt)
    df['splitBoxSplitUnit'] = df['splitBoxSplitUnit'].apply(convertDictToInt)
    df['movieInfo'] = df['movieInfo'].apply(lambda x : x['movieName'])
                                                           movieInfo showCount \
       avgSeatView avgShowView boxRate boxSplitUnit
[8]:
              1.1%
                                  25.8%
                                               6043400
                                                                    71278
    0
              0.7%
                                  12.2%
                                                                   67234
    1
                                               2868399
              1.2%
                              2
                                  11.6%
    2
                                               2719800
                                                                      49189
    3
              2.7%
                              3
                                  10.8%
                                               2543500
                                                                     24173
    4
              2.2%
                              3
                                   9.9%
                                               2320800
                                                            ů
                                                                    24790
              2.7%
                              3
                                   7.5%
    5
                                               1753000
                                                                      16943
              0.8%
                                   4.3%
    6
                              1
                                               1012300
                                                                        27188
    7
              1.2%
                              2
                                   3.8%
                                                                    12004
                                                905699
              0.6%
                                   3.4%
                                                                    31338
    8
                              1
                                                801800
    9
              0.5%
                              1
                                   2.8%
                                                667300
                                                                       27933
              1.8%
                              2
                                   1.5%
    10
                                                360500
                                                                5494
              5.3%
    11
                              6
                                   0.8%
                                                204400
                                                                       648
              1.0%
                                   0.7%
                                                                   4941
    12
                              1
                                                164500
    13
              5.5%
                              6
                                   0.6%
                                                163500
                                                                         303
    14
              0.3%
                                   0.6%
                                                                    10991
                              1
                                                150800
    15
              0.7%
                              1
                                   0.5%
                                                122600
                                                                          4085
    16
             35.0%
                             44
                                   0.4%
                                                107500
                                                                         54
              3.3%
                                   0.3%
    17
                              4
                                                 78900
                                                                    695
    18
             70.1%
                             65
                                   0.1%
                                                 35600
                                                                       11
    19
              0.5%
                                  <0.1%
                                                 21500
                                                                     1156
```

```
886
     20
               0.6%
                                  <0.1%
                                                17900
                              1
     21
               0.3%
                                  <0.1%
                                                                  1088
                              1
                                                12600
               0.2%
                                  <0.1%
     22
                              1
                                                10200
                                                             1922
                                                                        1273
       showCountRate splitBoxRate
                                   splitBoxSplitUnit sumBoxDesc sumSplitBoxDesc
                                                          38.02
     0
                18.4%
                             26.7%
                                              5758099
                                                                         35.08
                17.4%
                                                                         22.23
     1
                             12.2%
                                              2635899
                                                          24.31
     2
                12.7%
                             11.5%
                                                                         11.52
                                              2479200
                                                          12.60
     3
                 6.2%
                             10.5%
                                                         2272.8
                                                                        2002.5
                                              2259200
     4
                 6.4%
                              9.9%
                                                           9.08
                                                                          8.33
                                              2137100
                 4.3%
     5
                              7.4%
                                                         4197.0
                                                                        3847.0
                                              1591300
     6
                7.0%
                              4.2%
                                               923100
                                                           5.38
                                                                          4.94
     7
                 3.1%
                              3.8%
                                               818600
                                                         5040.0
                                                                        4540.9
                              3.3%
     8
                8.1%
                                               721500
                                                          272.2
                                                                         245.2
     9
                 7.2%
                              2.8%
                                                          349.7
                                                                         316.6
                                               608400
                                                                          1.41
     10
                 1.4%
                              1.5%
                                               328800
                                                           1.55
                 0.1%
                              0.7%
                                                         2168.7
                                                                        2000.6
     11
                                               163600
     12
                 1.2%
                              0.6%
                                                           1.02
                                                                        9283.5
                                               149700
                <0.1%
                              0.7%
                                                           28.4
     13
                                               157300
                                                                          26.9
     14
                 2.8%
                              0.6%
                                               135700
                                                          120.9
                                                                         109.9
     15
                 1.0%
                              0.5%
                                                           5.33
                                                                          4.82
                                               112899
     16
                <0.1%
                              0.4%
                                                99700
                                                         1343.0
                                                                        1336.3
     17
                 0.1%
                              0.3%
                                                70400
                                                         8153.5
                                                                        7218.9
                                                32900
     18
                <0.1%
                              0.1%
                                                         3273.4
                                                                        3255.1
     19
                 0.3%
                             <0.1%
                                                19700
                                                           1.66
                                                                          1.51
     20
                 0.2%
                             <0.1%
                                                16100
                                                         5550.3
                                                                        5021.3
                             <0.1%
                                                                         224.4
     21
                 0.2%
                                                11399
                                                          246.0
     22
                 0.3%
                             <0.1%
                                                 9200
                                                          208.6
                                                                         188.1
 [9]: #
     body0['movieList']['nationBoxInfo']
 [9]: {'nationBoxSplitUnit': {'num': '.',
       'unit': ''},
      'nationSplitBoxSplitUnit': {'num': '.',
       'unit': ''},
      'showCountDesc': '38.5',
      'title': '',
      'viewCountDesc': '57.2'}
[10]: # Comparing with the screenshot earlier
     from IPython.display import Image as displayImage
     displayImage(filename="logs/" + now + "/screenshot.png")
```

[10]:



[]:[