

Might be useful to see how we can get the last page for this website.

Keep in mind every website is designed differently and the code below will not work with other websites. The way I found the last page number was to find the text representing it in the HTML code.

span class="pageNumbersInfo">Страница 1 OF 25</span

Translated: span class="pageNumbersInfo">Страница 1 OF 25</span

It is interesting to mention that this website has 25 pages of properties for each municipality and each page has 40 property listings.

```
In [16]: def transform_page_number(soup):  
    """  
    This simple function is used to find the last page number for each municipality.  
    Assumes that the page number is a 2 digit number  
    """  
  
    # We are looking for 'span' tag from class 'pageNumbersInfo'  
    last_page = soup.find('span', class_ = 'pageNumbersInfo').text.strip()[-2:]  
    last_page_number.append(last_page)  
  
    return last_page_number  
  
# we create empty list to which we append info  
property_list = []  
last_page_number = []  
  
search = extract(i)  
print(transform_page_number(search))  
  
['25']
```

In order to get results from multiple municipalities we need to alter the code for the 'extract' function. The code for 'transform' function stays the same, it is not the best code. We can see repetition occurring.

The 'modified_extract' function takes another argument called 'municipality' that is the code of the municipality from the URL for it. We will use it instead of 'extract' from now on.

```
In [17]: def modified_extract(page, municipality):  
    """  
    This function is used to change the number of the page from the link.  
    So after we finish with collecting the data from page 1, we can continue to page 2  
    of that webpage and so on.  
    """  
  
    # https://router-network.com/tools/what-is-my-user-agent -> to find the 'user agent'  
    # or google -> 'what is my user agent'  
    # 'user agent' is a mediator between the user and the internet it holds technical information  
    # about the device (the computer used) and the software  
    # the 'user agent' is unique for every person on the internet  
    headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 6.3; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/60.0.3112.101 Safari/537.36'}  
  
    # f-string in which we can change the number of the page  
    # note that 'page' in the f-string is the parameter of the function  
    # in addition to that now we can input different strings that represent municipalities  
    url = f'https://www.imot.bg/poc/1/mot.cq?act=detail&link={municipality}&f={page}'  
  
    r = requests.get(url, headers)  
  
    soup = BeautifulSoup(r.content, 'html.parser')  
  
    return soup
```

```
In [18]: # downloads content from different locations/municipalities  
  
def transform(soup):  
    """  
    This function is used to find and collect information.  
    """  
  
    # We can find all 'div' tags that store the price of each property listing  
    # in that way we can get all 40 listings  
    price_div = soup.find_all('div', class_ = 'price')  
  
    # Each listing is associated with a particular style that actually shows if the listing  
    # is 'Top', 'Vip' or just ordinary listing  
    style_0 = "margin-bottom:0px; border-top:#990000 1px solid; background:url(../images/pictures/top_bg.gif);  
    style_1 = "margin-bottom:0px; border-top:#990000 1px solid; background:url(../images/pictures/vip_bg.gif);  
    style_2 = "margin-bottom:0px; border-top:#990000 1px solid;"  
  
    # We locate all tags 'table' that contain all the information we can scrape and we associate  
    # the 'table' tag with the different styles  
    table_0 = soup.find_all('table', style=style_0)  
    table_1 = soup.find_all('table', style=style_1)  
    table_2 = soup.find_all('table', style=style_2)  
  
    # The following three for loops go through each 'table' tag and scrape information from it  
    # The each for loop stores the information into a dictionary. Each dictionary is then appended  
    # to a list that contains all the information:  
    for i in table_0:  
        property_info = soup.find("a", class_ = 'photoLink').img.get('alt').strip()  
        price = i.find('div', class_ = 'price').text.strip()  
        rooms = i.find('a', class_ = 'lnk1').text.strip()  
        location = i.find('a', class_ = 'lnk2').text.strip()  
        info = i.find('td', width="520", colspan="3", height="50", style="padding-left:4px").text.strip()  
  
        # dictionary to store the collected information  
        property_dictionary_0 = {  
            'property_info': property_info,  
            'price': price,  
            'rooms': rooms,  
            'location': location,  
            'info': info  
        }  
        property_list.append(property_dictionary_0)  
  
    for i in table_1:  
        property_info = soup.find("a", class_ = 'photoLink').img.get('alt').strip()  
        price = i.find('div', class_ = 'price').text.strip()  
        rooms = i.find('a', class_ = 'lnk1').text.strip()  
        location = i.find('a', class_ = 'lnk2').text.strip()  
        info = i.find('td', width="520", colspan="3", height="50", style="padding-left:4px").text.strip()  
  
        # dictionary to store the collected information  
        property_dictionary_1 = {  
            'property_info': property_info,  
            'price': price,  
            'rooms': rooms,  
            'location': location,  
            'info': info  
        }  
        property_list.append(property_dictionary_1)  
  
    for i in table_2:  
        property_info = soup.find("a", class_ = 'photoLink').img.get('alt').strip()  
        price = i.find('div', class_ = 'price').text.strip()  
        rooms = i.find('a', class_ = 'lnk1').text.strip()  
        location = i.find('a', class_ = 'lnk2').text.strip()  
        info = i.find('td', width="520", colspan="3", height="50", style="padding-left:4px").text.strip()  
  
        # dictionary to store the collected information  
        property_dictionary_2 = {  
            'property_info': property_info,  
            'price': price,  
            'rooms': rooms,  
            'location': location,  
            'info': info  
        }  
        property_list.append(property_dictionary_2)  
  
    return len(table_0), len(table_1), len(table_2), len(price_div)  
  
# we create empty list to which we append the information gathered from the three for loops  
property_list = []  
  
# These are two list variables that store strings for the municipality and  
# the actual name of the municipality  
municipality = ['709q9q', '709q7a']  
municipality_translation = ['Veliko Tarnovo', 'Sofia City']  
  
# For loop that does the web scraping  
for i in range(len(municipality)):  
    # = municipality[i]  
    for j in range(1,3):  
        print(f'Page number: {j} - {municipality_translation[i]}')  
        search = modified_extract(j,k)  
        transform(search)  
  
print(len(property_list))  
  
Page number: 1 - Veliko Tarnovo  
Page number: 2 - Veliko Tarnovo  
Page number: 1 - Sofia City  
Page number: 2 - Sofia City  
160
```

What is important is that we have managed to extract all the property listings of two pages for two municipalities. That means that we can scrape all the property listings if we just add the code for each municipality to the list 'municipality'. In addition to that we can see that for the four pages that we scraped the sum of listings is equal to 160, which means that all is working as it should.

What is left to do is to make the code for the 'transform' function more compact and to reduce the amount of code. In that way we will reduce the repetition and improve the quality of the code.

```
In [19]: def modified_transform(soup):  
    """  
    # finds the total number of pages  
    # assumes that the page number is a 2 digit number  
    # span class="PageNumbersInfo">Страница 1 OF 25</span> - Page 1 of 25  
    """  
  
    last_page = soup.find('span', class_ = 'pageNumbersInfo').text.strip()[-2:]  
    last_page_number.append(int(last_page))  
  
    # We can find all 'div' tags that store the price of each property listing  
    # in that way we can get all 40 listings  
    price_div = soup.find_all('div', class_ = 'price')  
  
    # Each listing is associated with a particular style that actually shows if the listing  
    # is 'Top', 'Vip' or just ordinary listing  
    styles = ["margin-bottom:0px; border-top:#990000 1px solid; background:url(../images/pictures/top_bg.gif);",  
    "margin-bottom:0px; border-top:#990000 1px solid; background:url(../images/pictures/vip_bg.gif);",  
    "margin-bottom:0px; border-top:#990000 1px solid;"]  
  
    # The 'table' tag that contain the information do it have class, the only way to distinguish  
    # one from another is the 'style' attribute that they have  
    # We can create a for loop that will collect automatically all the information instead repeating  
    # the process three times for each table style and if something is changed we can easily  
    # amend where it is needed.  
    for i in range(len(styles)):  
        table_x = soup.find_all('table', style=styles[i])  
  
        for j in table_x:  
            property_info = soup.find("a", class_ = 'photoLink').img.get('alt').strip()  
            price = j.find('div', class_ = 'price').text.strip()  
            rooms = j.find('a', class_ = 'lnk1').text.strip()  
            location = j.find('a', class_ = 'lnk2').text.strip()  
            info = j.find('td', width="520", colspan="3", height="50", style="padding-left:4px").text.strip()  
  
            # dictionary to store the collected information  
            property_dictionary = {  
                'property_info': property_info,  
                'price': price,  
                'rooms': rooms,  
                'location': location,  
                'info': info  
            }  
            property_list.append(property_dictionary)  
  
    return len(price_div), last_page_number[0], len(property_info)
```

```
# we create empty list to which we append info  
property_list = []  
  
# variable to store the the last page number  
last_page_number = []  
  
# FOR INFORMATION  
# The URL code for each municipality is different every now and again.  
# This is just an example on how to collect them.  
# Keep in mind that those codes have expired and will give error,  
# because it can not find the page number if the municipality code from the URL is changed.  
# municipality = ['709q7h', '70k4qj', '70yb7j', '70jay6', '70k56o', '70k4qh']  
# municipality_translation = ['Veliko Tarnovo', 'Sofia City', 'Sofia Municipality',  
#                             'Lovetch', 'Burgas', 'Varna']  
  
# You can see that at the moment 'Veliko Tarnovo' code is '71zaba' and not '70q7yh'.  
municipality = ['709q9q']  
municipality_translation = ['Veliko Tarnovo']  
  
# In order to get the last page number we need to run 'modified_extract' function once  
# so we can get the value of 'last_page_number' variable and empty the 'property_list' variable  
# For information - it turns out that this site has only 25 pages with 40 properties  
# per page for each municipality  
search = modified_extract(1, '709q9q')  
modified_transform(search)  
print(transform(search))  
property_list = []
```

The following cell contains two identical in what they do for loops. The difference between them is that the first one finds the last page number, where as the second one has the number hard-coded.

```
In [20]: property_list = []  
  
for i in range(len(municipality)):  
    for j in range(1, int(last_page_number[0])):  
        print(f'Page number: {j} - {municipality_translation[i]}')  
        search = modified_extract(j, municipality[i])  
        modified_transform(search)  
  
# This is the same for-loop as the one from above, it is just using out knowledge that the last page number is  
# for j in range(1, 26, 1):  
#     print(f'Page number: {j} - {municipality_translation[i]}')  
#     search = modified_extract(j, municipality[i])  
#     transform(search)  
  
print(len(property_list))  
print(last_page_number)  
  
Page number: 0 - Veliko Tarnovo  
Page number: 1 - Veliko Tarnovo  
Page number: 2 - Veliko Tarnovo  
Page number: 3 - Veliko Tarnovo  
Page number: 4 - Veliko Tarnovo  
Page number: 5 - Veliko Tarnovo  
Page number: 6 - Veliko Tarnovo  
Page number: 7 - Veliko Tarnovo  
Page number: 8 - Veliko Tarnovo  
Page number: 9 - Veliko Tarnovo  
Page number: 10 - Veliko Tarnovo  
Page number: 11 - Veliko Tarnovo  
Page number: 12 - Veliko Tarnovo  
Page number: 13 - Veliko Tarnovo  
Page number: 14 - Veliko Tarnovo  
Page number: 15 - Veliko Tarnovo  
Page number: 16 - Veliko Tarnovo  
Page number: 17 - Veliko Tarnovo  
Page number: 18 - Veliko Tarnovo  
Page number: 19 - Veliko Tarnovo  
Page number: 20 - Veliko Tarnovo  
Page number: 21 - Veliko Tarnovo  
Page number: 22 - Veliko Tarnovo  
Page number: 23 - Veliko Tarnovo  
Page number: 24 - Veliko Tarnovo  
1000  
[25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25, 25]
```

We got 25 pages, which is the total amount and 25 times 40 is equal to 1000 and that is the exact number of property listings that we got.

We can just add more municipalities and to the 'municipality' list and it will collect all the data for us.

Then we can save the dataframe as shown before as CSV file and read the CSV into Pandas to start cleaning and analyzing the data.

Lastly, I want to show that it can work for different municipalities, but I will hard code it to take the first 2 pages of each so it will show that it can do the job as it should.

```
In [21]: municipality = ['709q9q', '709q5q', '709q8q']  
municipality_translation = ['Sofia City', 'Veliko Tarnovo', 'Varna']  
  
last_page_number = []  
property_list = []  
  
for i in range(len(municipality)):  
    for j in range(0, int(last_page_number[0])):  
        # for j in range(0, 2):  
        print(f'Page number: {j} - {municipality_translation[i]}')  
        search = modified_extract(j, municipality[i])  
        modified_transform(search)  
  
print(len(property_list))  
print(last_page_number)  
  
Page number: 0 - Sofia City  
Page number: 1 - Sofia City  
Page number: 0 - Veliko Tarnovo  
Page number: 1 - Veliko Tarnovo  
Page number: 0 - Varna  
Page number: 1 - Varna  
240  
[25, 25, 25, 25, 25, 25]
```

We can still check the results and make sure that everything is working fine. 6 * 40 = 240. That is obvious, but let's check what are the municipality names.

```
In [22]: # This function is pretty much the same as 'basic_translation' the only difference is that  
# this one doesn't have the pronunciation written down, which we do not need.  
  
def basic_translation_list(text_for_translation, target_language):  
    """  
    Basic function introducing 'google_trans_new'.  
    """  
    # We instantiate 'google_translator' object called 'translator'  
    translator = google_translator(url_suffix='qg', timeout=5)  
  
    # In the brackets we put the text we need translated first, then we define 'lang_tgt' parameter  
    # that is used to tell to what language we need to translate the original text  
    translate = translator.translate(text_for_translation,  
                                    lang_tgt = target_language)  
  
    return print(translate)  
  
def dictionary_into_dataframe(dictionary):  
    """  
    This function takes as an argument a dictionary and creates a DataFrame using Pandas.  
    The function returns as a result the municipalities.  
    """  
  
    # We instantiate a DataFrame object  
    dataframe = pd.DataFrame(dictionary)  
  
    # We take only the 'location' column for translation  
    dataframe_unique_locations = dataframe['location'].unique()  
  
    # The actual translation  
    translation = basic_translation_list(dataframe_unique_locations, 'en')  
  
    return translation
```

```
In [23]: dictionary_into_dataframe(property_list)
```

If we take a second and look at the result we have from the cell below we can see that the municipalities listed below are as follows: 'city of Sofia', 'Veliko Tarnovo region' and 'Varna Region'. That means that the Web scraper is doing exactly what we want.

The main task that we have specified in the beginning of this Jupyter Notebook is now accomplished. That doesn't mean that we still do not have work to do for the data that we have extracted, but this is for another Jupyter Notebook.

Thank you very much for your attention and time.

Kind regards,

Mladen