# Mandatory Assignment 5
## Regular Expressions & Web Scraping (30 + 15 Points/40 + 5 Points)

### University of Oslo - IN3110/IN4110

### Fall 2020

Your solution to this mandatory assignment needs to be placed in the directory `assignment5` in your Github repository. The repository has to contain a `README.md` file with information on how to run your scripts, required dependencies and packages (versions the code ran with) as well as how to install them. Documentation on how to run your examples[1] is required. Furthermore, your code needs to be well commented and documented. **All** functions need to have docstrings explaining what the function does, how it should be used, an explanation of the parameters and return value(s) (including types). We **highly** recommend you use a well-established docstring style such as the Google style docstrings[2]. However, you are free to choose your own docstring style. We expect your code to be well formatted and readable. Coding style and documentation will be part of the point evaluation for **all tasks** in this assignment.

**Files to Deliver for this Task**

All `.py`-files will be delivered in the same `assignment5` folder. Each subtask will create an output folder with all outputs as stated in the subtask.

*Files required for all tasks*

- `README.md` including dependencies and their installation and commands you used for running your examples and creating your output files

**Note:** Regular expressions can be tricky to write, and they can get hairy very quickly. In addition to being difficult to write, regular expressions are even harder to read. Documentation for your regular expressions should also be given. You are also required to supply files that can be used to demonstrate your solutions.

- You are required to supply everything necessary to demonstrate your code, i.e. you should include the code examples you use when testing your work, or something equivalent.

---

[1] Not only important to get all of the credits, but it is a really important and not understated habit to be established :)

[2] https://sphinxcontrib-napoleon.readthedocs.io/en/latest/example_google.html

## 5.1 Sending url requests (2 Points)

In order to parse html and become the regex master of the year, you need to access the html body first. This can be done using the python request module [3].

The module can be installed as follows:

```
pip install requests
```

To retrieve data from any destination, the most common approach is the **get** Request. This is the most used HTTP method. In our example, we will use the get request to fetch data from the website of choice. The result will be printed out as html data. The request module can be imported and used as shown below.

```
#import the module
import requests as req

#grabbing the content of https://en.wikipedia.org/wiki/URL
resp = req.get("https://en.wikipedia.org/wiki/URL")

#get() method returns a response object
print(resp.text)
<!DOCTYPE html>
<html class="client-nojs" lang="en" dir="ltr">
<head>
<meta charset="UTF-8"/>
<title>URL - Wikipedia</title>

#shortened for the purpose of viewing
```

Fantastic, now we can actually display the data in html!

Often you are interested in a specific type of data. With GET requests, you can send parameters in the form of query strings. How does that work? If you were to construct the URL by hand, the data you are interested in would be given as key/value pairs in the URL after a question mark, e.g. httpbin.org/get?key=val. Using get requests, you can provide these arguments as a dictionary of strings, using the params keyword argument. That means, if you wanted to pass key1=value1 and key2=value2 to httpbin.org/get the code would in theory look like this:

```
import requests as req

params = {'key1': 'value1', 'key2': 'value2'}
r = req.get('https://httpbin.org/get', params=params)

# the url would be encoded accordingly
print(r.url)
https://httpbin.org/get?key2=value2&key1=value1
```

---

[3]https://requests.readthedocs.io/en/master/user/install/

For a real example you can look at the following snippet:

```python
import requests as req


params = {'user_name': 'admin', 'password': 'password'}
r = req.get('http://httpbin.org/get', params=params)

print(r.url)
http://httpbin.org/get?user_name=admin&password=password

# print the.text content of the request
print(r.text)

{
  "args": {
    "password": "password",
    "user_name": "admin"
  },
  "headers": {
    "Accept": "*/*",
    "Accept-Encoding": "gzip, deflate",
    "Host": "httpbin.org",
    "User-Agent": "python-requests/2.22.0",
    "X-Amzn-Trace-Id": "Root=1-5f841066-0
    c0ea5084573be8a230feaea"
  },
  "origin": "51.175.234.27",
  "url": "http://httpbin.org/get?user_name=admin&password=
    password"
}
```

Your task will be to create the function `get_html` which makes a url request of a given website. You should test your function on the following websites:

- `https://en.wikipedia.org/wiki/Studio_Ghibli`

- `https://en.wikipedia.org/wiki/Star_Wars`

- `https://en.wikipedia.org/wiki/Dungeons_%26_Dragons`

Furthermore you should extend the function to optionally take in parameters to be passed to the get function. There should be an optional argument allowing to specify that the response url and text will be saved to a text file with a specified name `optionalargument.txt`. If the optional argument is not set, the response is just returned.

The function should be able to be used like this:

```python
get_html(url, params, output)
```

```
# if specified , the response txt and url get printed to a .
    txt − file with the name output.txt
```

The extended function with parameters should be tested for these websites:

- `https://en.wikipedia.org/w/index.php`

  - with parameters title=Main_Page and action=info

- `https://en.wikipedia.org/w/index.php`

  - with parameters title=Hurricane_Gonzalo and oldid=983056166

**Files to Deliver in this Subtask**

Create a folder `requesting_urls` where you put all output files from our test runs created for solving this subtask.

*Files Required in this Subtask*

- `requesting_urls.py`

- Five files inside the `requesting_urls` folder containing the function output for the five websites example websites.

## 5.2 Regex for filtering URLs (5 Points)

In this task, you will be making functions for finding urls in a body of html using regex. Create a script named `filter_urls.py` that includes a function `find_urls` that receives a string of html and returns a list of all urls found in the text. You will be using the `re.findall` method taught in the lecture. This will give us a list of all the strings matching your regular expression. To be strict, we only consider urls that are anchor[4] hyperlinks with the `<a>` HTML tag. The function should ignore fragment identifiers. That is, links that start with a hash (#) symbol. If the url points to a different page, but a specific fragment/section using the hash symbol, the fragment part of the url should be stripped before it is returned by the function.

**Example:** The string `https://www.example.com/somepage#someidentifier` becomes `https://www.example.com/somepage`

In order to handle relative urls, the function could optionally receive a base url.

Further, make a function `find_articles` that calls `find_urls` and returns only urls to Wikipedia articles. This function should also use regex and **be able to handle any chosen language** (`no.wikipedia.org`, `en.wikipedia.org` etc).

**Note:** Non wikipedia pages can be ignored. Combine these functions with your functions from previous tasks in order to test them on the following websites:

---

[4]name of the HTML-element, with the HTML-tag "a"

- `https://en.wikipedia.org/wiki/Nobel_Prize`

- `https://en.wikipedia.org/wiki/Bundesliga`

- `https://en.wikipedia.org/wiki/2019%E2%80%9320_FIS_Alpine_Ski_World_Cup`

**Note 1:** You should only return normal Wikipedia articles, so no special namespace article (or files). While you could exclude just the specific reserved namespaces [5], it is sufficient to exclude all articles with a colon. This will exclude some actual articles (for example `https://en.wikipedia.org/wiki/Avengers:_Endgame`, but we won't deduce points for it.

**Note 2:** You should only find urls that are in the href attribute of the a tag. It is *not always* the first attribute which means matching just ¡a href= is not enough, but you can assume some things about the use of special html characters since regex cannot parse html alone.

**Note 3:** You need to include relative urls, best be combined with the base url. Examples of relative urls: `/wiki/Executive_director` and `//en.wikipedia.org/wiki/Wikipedia:Contact_us`

**Note 4¿** The outputs for both first `fins_urls` and `fins_articles` should be added to your outputs.

The found urls should be saved to a .txt file, which name can be defined by an optional argument. **Files to Deliver in this Subtask**

Create a folder `filter_urls` where you put all output files from our test runs created for solving this subtask.

*Files Required in this Subtask*

- `filter_urls.py`

- Three files inside the `filter_urls` folder containing a list of the urls returned for each of the three example websites.

**Note:** You should not use Beautiful Soup or any HTML parser in this task, only regex.

## 5.3 Regular Expressions for finding Dates (IN3110 OPTIONAL & IN4110 REQUIRED 10 Points)

In this task, you will be making functions for finding "dates" in a body of html using regex. Create a script named `collect_dates.py` that includes a function `find_dates` that receives a string of html and returns a list of all dates found in the text in the following format:

```
1)  1998/10/12
2)  1998/11/04
3)  1999/01/13
```

---

[5]i.e. look here `https://en.wikipedia.org/wiki/Wikipedia:Namespace`

**Which date formats do we need to consider?** You need to consider the wikipedia standardized formats [6]. These are the four formats:

```
DMY:  13 October 2020
MDY:  October 13, 2020
YMD:  2020 October 13
ISO:  2020−10−13
```

**Note:** You also need to include the case, where the month is abbreviated, e.g. Oct for October, which adds these date types to consider.

```
DMY:  13 Oct 2020
MDY:  Oct 13, 2020
YMD:  2020 Oct 13
```

We are aware, that Wikipedia should not support the abbreviated month, but we would like to allow your code to be flexible, in case you are handling other websites in the future

The function should furthermore be able to **recognize, if there is no day**, but only the month and the year. Then it should just return the year and the months. This optional day argument should work for the following formats resulting from the four formats you are working on:

```
DMY:  October 2020
MDY:  October, 2020
```

Since `2020-10` and `2020 October` can result in many false positives, there will **NOT** be points withdrawn if for that one if the missing date version is **NOT** implemented.

For a list of dates found which have two dates with `year/month/day` and one with `year/month` the returned list could look like this:

```
1)  1998/10/12
2)  1998/11
3)  1999/01/13
```

This is consequently the complete list of formats to consider:

```
DMY:  13 October 2020 AND October 2020
MDY:  October 13, 2020 AND October, 2020
YMD:  2020 October 13
ISO:  2020−10−13
```

The dates **returned** should be formatted like this `year/month/day`.

Moreover, there should be an optional argument allowing to specify that the resulting list will be saved to a text file with a specified name `optionalargument.txt`.

You will be using the `re.findall` and `re.sub` method taught in the lecture. This will give us a list of all the strings matching your regular expressio and allow to change the formatting.

Here is the list of websites, you need to deliver a report in form of a text-file for:

*Steps for creating the script*

---

[6]`https://en.wikipedia.org/wiki/Template:Date`

- visit the following websites

- `https://en.wikipedia.org/wiki/Linus_Pauling`

- `https://en.wikipedia.org/wiki/Rafael_Nadal`

- `https://en.wikipedia.org/wiki/J._K._Rowling`

- `https://en.wikipedia.org/wiki/Richard_Feynman`

- `https://en.wikipedia.org/wiki/Hans_Rosling`

Save your implementation to `collect_dates.py`.

**Files to Deliver in this Subtask**

Create a folder `filter_dates_regex` where you put all output files from our test runs created for solving this subtask.

*Files Required in this Subtask*

- `collect_dates.py`

- Five files inside the `filter_dates_regex` folder containing a list of found dates in sorted order, for each website given.

**Note:** You should not use Beautiful Soup or any HTML parser in this task, only regex.

## 5.4 Making your Life easier with Soup for filtering datetime objects (IN3110 & IN4110 8 Points)

Need to plan your watch parties for the upcoming skiing season? Why not send out all the dates and times to your friends that are as well into (watching) skiing? In this task you will use the Python tool Beautiful soup to parse data. You will extract the datetime objects representing the event "date" , a regular expression fetching the "venue" and the discipline "type".

You can install Beautiful Soup via:

```
pip install beautifulsoup4
```

**Note:** Make sure that you are running python3 and your pip points to the right python version as well. If you are running into issues we highly recommend setting up a `virtualenv` or a `conda env`.

Before starting with the fun-coding part, it is recommended to take some time to get familiar with the html tags of the website we are extracting information from.

**Note:** It will really help for this task if you just take a few minutes to familiarize yourself with the page setup.

We recommend to navigate to the website using your favorite browser, right click and 'Inspect'. Get familiar with the elements used. Alternatively, right click on the web page and view the page source. Since we are interested in a table, search for the table class :

Figure 1: Table class in the html. You get to this view using "View Source Page" mode.

Okay, so we would like to access the table of soups now. To do so we need to first grab the HTML from the web page, so that we have something to parse through. Then we will extract the soup table.

```python
from bs4 import BeautifulSoup
import requests as req

url = "https://en.wikipedia.org/wiki/List_of_soups"

# request url as we have already learned - yeah!
request = req.get(url)

soup = BeautifulSoup(request.text, "html.parser")

#check the title of the wikipedia page
print(soup.title)

# get the soup table
soup_table = soup.find('table', {"class":'wikitable sortable'})
```

Great, now we need to know how to navigate through the table. To figure out which column holds which information, it is recommended to look at the header. Therefore, we will use the html <th> tag, which defines the header cell

in an html table. On a sidenote, a html table has to kinds of cells: header cells created with `<th>` elements and data cells created with `<td>` elements.

To iterate through the data in the table, we first need to access the table rows. Table rows are defined by the `<tr>` tag. A `<tr>` element contains `<td>` or `<th>` elements. The table headers will be extracted via `<th>` elements,

Write a script `time_planner.py` that first requests the url from `https://en.wikipedia.org/wiki/2019%E2%80%9320_FIS_Alpine_Ski_World_Cup`. The request can be done with a function created in earlier tasks.

Then parses through the html using `Beautiful Soup` and finds the main table in the calendar section (see figure below).



Figure 2: Main table in calendar section.

Write a function **extract_events** that extracts the data in the date, venue and discipline column. The function will live in the script `time_planner.py`. Since you would like to gamble with your friends beforehand you will create a betting slip. That slip will contain and a table which columns represent the date, venue and discipline as well as a column for the tip of which athlete or country wins that competition. The empty betting slip should be created automatically by your script and saved to a new file named `betting_slip_empty.md` using Markdown [7] formatting.

Check that your script also works when accessing `https://en.wikipedia.org/wiki/2020%E2%80%9321_FIS_Alpine_Ski_World_Cup`.

**Files to Deliver in this Subtask**

---

[7] `https://github.com/adam-p/markdown-here/wiki/Markdown-Cheatsheet#tables`

Create a folder `datetime_filter` where you put all output files from our test runs created for solving this subtask. *Files Required in this Subtask*

- `time_planner.py`

- `datetime_filter/betting_slip_empty.md` - a nicely formatted table for your bets, containing date, venue and discipline, but with an empty "who wins" column.

The resulting betting slip would look like this:

BETTING SLIP

Name:

| DATE | VENUE | DISCIPLINE | Who Wins? |
|------|-------|-----------|-----------|
| 01/01/2021 | nowhere | sitting | |
| | | | |
| | | | |
| | | | |

Figure 3: Betting Slip.

## 5.5 NBA Player Statistics Season 2019/2020 (IN3110 & IN4110 15 Points)

Go Miami Heat - Whooo! Let's imagine you are actually into basketball. After this unpredictable season, you really want dive into the player statistics. Lucky you! In this task you are going to write a script `fetch_player statistics.py` which visits the "2020 NBA playoffs" website on wikipedia (this one: `https://en.wikipedia.org/wiki/2020_NBA_playoffs`) and creates some player statistics. As in task `5.4` will use `BeautifulSoup` for parsing.

As always, you first need to make a url request using your function `get_html`.

Then you want to navigate, using `BeautifulSoup` to the main table in the "Bracket" section shown below.
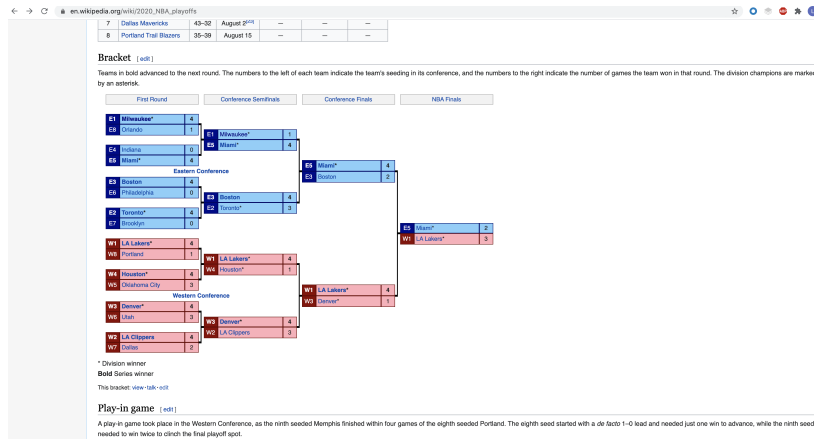
Figure 4: Bracket Section.

There you will extract the names of the teams which made it to the conference semifinals (This is shown in the section "Bracket"). [8]. You will also create a function `extract_url` which will allow you to extract the team urls in this table. You will follow the urls to the wikipedia websites of the corresponding teams in the Bracket section.

---

[8]Hint: It should be 8.

The next figure shows the website you get to when following the Milwaukee Bucks link in the conference semifinals.



Figure 5: Milwaukee Bucks Page.

In the next step you need to identify all the players that played for the specific team in that season. Therefore you need to navigate to the table holding the team roster. For the Milwaukee Bucks the roster section looks like this:



Figure 6: Roster Milwaukee Bucks.

Fantastic, let's start on the player statistic section! In order to filter the

statistics for one specific player you need to extract the player name and the url to this players wikipedia page. You will request the `url` to that player which can be found in the roster. If we follow the link to Giannis Antetokounmpo this would be the part of the webpage showing the player statistics.



Figure 7: Giannis Antetokounmpo.

We are interested in the NBA Regular season table. Therefore we need to navigate to this one. We want to extract for the row **Year 2019-20** and the columns points per game, blocks per game and rebounds per game.

For your statistics you are going to compare the best players from the teams in the conference semifinals. "Best" is defined by the highest count of points per game in our case.

The three best players of each team make it into our comparison pool[9] If statistics for a player is missing it is fine to ignore this player. Since you are only selecting the top 3 from each team you could just give these people the lowest possible score (this way they are not included).

We want to have a look at the best player with respect to points per game, with respect to blocks per game and rebounds per game.

Therefore we will plot the player over the points/blocks/rebounds. You can choose your favorite plotting tool for this. We would like to have players of each team grouped together ( this can be color-coded or with a bracket, or another visualization).

Your implementation should be saved to a file called `fetch_player statistics.py`.
The produced plots should be saved.
*Steps for creating the script*

- visit `https://en.wikipedia.org/wiki/2020_NBA_playoffs`

_____

[9]Since we take the 3 best players of each team, we will end up comparing the abilities of 24 players in total.

- crawl through the team and player wikipedia websites via the internal urls as described above

- define the best 3 players of each team in the conference semifinals based on the "Points per game" in 2019/20

- fetch their statistics in the categories points per game, blocks per game and rebounds per game

- create three plots (players of each team grouped together in the plots):

   - Players over points per game

   - Players over blocks per game

   - Players over rebounds per game

- Save the plots!

**Files to Deliver in this Subtask** Create a folder `NBA_player_statistics` where you put all output files from our test runs created for solving this subtask.
*Files Required in this Subtask*

- `fetch_player statistics.py`

- `NBA_player_statistics/players_over_ppg.png`

- `NBA_player_statistics/players_over_bpg.png`

- `NBA_player_statistics/players_over_rpg.png`


## 5.6  Challenge - Wiki Race with URLs (5 bonus points and a price for the winner!)

Everyone has probably heard of golf. You try to get the ball into the hole with the least amount of hits. Let us play some Wikipedia golf.

Write a script using your previous functions that finds the shortest way (in number of urls to visit) from `https://en.wikipedia.org/wiki/Parque_18_de_marzo_de_1938` to `https://en.wikipedia.org/wiki/Bill_Mundell` using only urls in Wikipedia articles. The script should work with any wikipedia url.

The main objective is to find the shortest path. You can expect to work on the english wiki only for this task. The websites given will be your test case.

**Note:** Since the simplest approach might technically work, but not be the fastest finding the solution, we would like you to still be able to provide a solution for bonus points, therefore you can find the shortest path lengths between these two wikipedia articles: `https://en.wikipedia.org/wiki/Nobel_Prize` and `https://en.wikipedia.org/wiki/Array_data_structure`.

14

In order to assign a winner, we will evaluate this for 2 undisclosed urls. The most efficient (and correct) script will win a prize :)

Your chances for winning are even higher, if your script is fast.

**Note:** You are going to be on english wikipedia en.wiki .

**You have to use python ;)**

**Files to Deliver in this Subtask**

Create a folder `wiki_race_challenge` where you put all output files from our test runs created for solving this subtask.

*Files Required in this Subtask*

- `wiki_race_challenge`.py

- `wiki_race_challenge/shortest_way.txt`

Wohoo - You finished another assignment! Congrats!

**Update October, 13th 2020**

- check Section 5.2 - anchor explanation added, added requirements for urls, identifier example

- check Section 5.3 - added 3 Notes, defined date formats more precisely)

**Update October, 14th 2020**

- check Section 5.3 - defined that yyyy/mm[/dd] , where only the day optional

**Update October, 15th 2020 Q&A**

- **For subtask 5.3, are we meant to convert the date with regular expressions, or can we use a standard library?** - Regex should be used everywhere in task 5.3 and 5.3 - so yes, use regex please.

- **Do we need to use a regex to fetch the venue in 5.4?** - No need to fetch it with regex in 5.4

**Update October, 16th 2020** Updated Task 5.2, 5.3 and 5.6 questions. **Q&A**

- **For subtask 5.3, are we meant to convert the date with regular expressions, or can we use a standard library?** - Regex should be used everywhere in task 5.3 and 5.3 - so yes, use regex please.

- **What's a relative url and what's a base url?** - The url you see when searching for Giannis Antetokounmpo in your browser is the full or absolute url `https://en.wikipedia.org/wiki/Giannis_Antetokounmpo`. An absolute URL is the entire address from the protocol (https), to the domain name including the location within your website (path). A relative url does not use the full web address but only contains the path followed by the domain. Therefore it assumes the link is on the same site and is part of the same root domain. The relative url linking to Giannis Antetokounmpos wiki page looks like this:

```
<a href="/wiki/Giannis_Antetokounmpo"
```

The base url is represented by `https://en.wikipedia.org/` in our example.

- **Do all of the dates in 5.3 be formatted the same? As in either 06 or 6 for june, or is it fine that they are mixed?** The dates should al be formatted in the same style, which is **yyyy\mm\dd**, so june would be represented with 06.

- **(5.3) When sorting should we use regex for that as well or is it fine to use inbuilt functions?** For the sorting of the list you can use a built-in function.

- **In 5.3 it says we should get the dates without days as well. Im wondering if this is needed for ISO dates as well or only the other formats?** Yes, this is required for the iso dates as well. So basically all date types you include in the yyyymmdd version of the function, should be included.

**Update October, 19th 2020** Updated 5.3 completely and simplified challenge. Updated Q& A. **Q&A**

- **5.3 - Should we remove duplicate dates?** - You can keep duplicate dates without us withdrawing points.

- **5.4 - In the video lecture about beautiful soup I saw that you were using the lxml parser. Are we allowed to use it, or do we need to use the standard python html.parser?** - You can always use what you learned in the lecture.

- **5.4 - Are we supposed to focus on the 19-20 or 20-21 season?** -Your implementation should work on both pages, and you just test both, to verify this. The 2019-20 table contains maybe a little bit more interesting information, so I would focus on making this one work first.

- **5.4 - Are we meant to use only soup or soup and regex in 5.4?** - 5.4 is a beautiful soup task , so you can use beautiful soup only :)

- **5.4 - In task 5.5 there's a lot of players with missing information for example:** `https://en.wikipedia.org/wiki/Johnathan_Motley` **misses the 2019-20 season** `https://en.wikipedia.org/wiki/Paul_Watson_(basketball)` **misses a stats table** - If statistics for a player is missing it is fine to ignore this player. Since you are only selecting the top 3 from each team I would say you could just give these people the lowest possible score (this way they are not included).