

# 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<sup>1</sup> 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<sup>2</sup>. 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.

---

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

<sup>2</sup>[https://sphinxcontrib-napoleon.readthedocs.io/en/latest/example\\_google.html](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 <sup>3</sup>.

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
```

---

<sup>3</sup><https://requests.readthedocs.io/en/master/user/install/>

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

```
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=1-5f841066-0c0ea5084573be8a230feaea"
  },
  "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/Studio_Ghibli)
- [https://en.wikipedia.org/wiki/Star\\_Wars](https://en.wikipedia.org/wiki/Star_Wars)
- [https://en.wikipedia.org/wiki/Dungeons\\_%26\\_Dragons](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:

```
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 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.

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).

Combine these functions with your functions from previous tasks in order to test them on the following websites:

- [https://en.wikipedia.org/wiki/Nobel\\_Prize](https://en.wikipedia.org/wiki/Nobel_Prize)
- <https://en.wikipedia.org/wiki/Bundesliga>
- [https://en.wikipedia.org/wiki/2019%E2%80%932020\\_FIS\\_Alpine\\_Ski\\_World\\_Cup](https://en.wikipedia.org/wiki/2019%E2%80%932020_FIS_Alpine_Ski_World_Cup)

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 BeautifulSoup 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" and "years" 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
```

The dates returned should be formatted like this `year/month/day`. The function should be able to find dates even if they are written as strings like `March 12, 2013`. The function should furthermore be able to recognize, if there is no day and month given, to just return the year. We also want to mention that there is a regional ambiguity with regards to date formatting. For instance, `9/10 2011`, could be either October, 9th or September, 10th. You can assume one of them, but mention which one you chose in your `README.md`. 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` method taught in the lecture. This will give us a list of all the strings matching your regular expression.

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

*Steps for creating the script*

- visit the following websites
- [https://en.wikipedia.org/wiki/Linus\\_Pauling](https://en.wikipedia.org/wiki/Linus_Pauling)
- [https://en.wikipedia.org/wiki/Rafael\\_Nadal](https://en.wikipedia.org/wiki/Rafael_Nadal)
- [https://en.wikipedia.org/wiki/J.\\_K.\\_Rowling](https://en.wikipedia.org/wiki/J._K._Rowling)

- [https://en.wikipedia.org/wiki/Richard\\_Feynman](https://en.wikipedia.org/wiki/Richard_Feynman)

- [https://en.wikipedia.org/wiki/Hans\\_Rosling](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 :

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.

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

[illegible]

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

```
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 two 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).

#	Event	Date	Venue	Type	Winner	Second	Third	Details
1783	1	18 October 2019	<span><span></span></span> Garmisch	GS 430				
1784	2	14 November 2019	<span><span></span></span> Lech/Zürs	PG 407				
		28 November 2019	<span><span></span></span> Lake Louise	DH 406				
		29 November 2019	<span><span></span></span> Lake Louise	SG 405				
		4 December 2019	<span><span></span></span> Beaver Creek	SG 404				North American Tour cancelled due to the coronavirus pandemic; most likely to be rescheduled in 2020-21
		5 December 2019	<span><span></span></span> Beaver Creek	DH 403				
		6 December 2019	<span><span></span></span> Beaver Creek	GS 402				
1785	3	5 December 2019	<span><span></span></span> Val d'Isère	GS 401				
1786	4	6 December 2019	<span><span></span></span> Val d'Isère	GS 400				
1787	5	12 December 2019	<span><span></span></span> Val d'Isère	DH 400				
1788	6	13 December 2019	<span><span></span></span> Val d'Isère	SG 400				
1789	7	18 December 2019	<span><span></span></span> Val Gardena/Gröden	SG 400				
1790	8	19 December 2019	<span><span></span></span> Val Gardena/Gröden	DH 400				
1791	9	20 December 2019	<span><span></span></span> Åre	GS 400				
1792	10	21 December 2019	<span><span></span></span> Åre	SL 400				
1793	11	22 December 2019	<span><span></span></span> Madonna di Campiglio	SL 400				
1794	12	28 December 2019	<span><span></span></span> Madonna di Campiglio	DH 400				
1795	13	29 December 2019	<span><span></span></span> Bormio	SG 400				
1796	14	6 January 2020	<span><span></span></span> Zagreb	SL 400				
1797	15	8 January 2020	<span><span></span></span> Zagreb	GS 400				
1798	16	9 January 2020	<span><span></span></span> Adelboden	GS 400				
1799	17	10 January 2020	<span><span></span></span> Adelboden	SL 400				
1800	18	15 January 2020	<span><span></span></span> Adelboden	DH 400				
1801	19	16 January 2020	<span><span></span></span> Hengelo	DH 400				
1802	20	17 January 2020	<span><span></span></span> Hengelo	SL 400				
1803	21	22 January 2020	<span><span></span></span> Hengelo	SG 400				
1804	22	23 January 2020	<span><span></span></span> Kitzbühel	DH 400				
1805	23	24 January 2020	<span><span></span></span> Kitzbühel	SL 400				
1806	24	26 January 2020	<span><span></span></span> Schladming	SL 400				
1807	25	30 January 2020	<span><span></span></span> Chamorro	SL 400				
1808	26	31 January 2020	<span><span></span></span> Chamorro	SL 400				
1809	27	5 February 2020	<span><span></span></span> Garmisch-Partenkirchen	SG 400				
1810	28	6 February 2020	<span><span></span></span> Garmisch-Partenkirchen	DH 400				
<b>FIS Alpine World Ski Championships 2021 (8-21 February)</b>								
1811	29	27 February 2021	<span><span></span></span> Bansko	GS 400				
1812	30	28 February 2021	<span><span></span></span> Bansko	GS 400				
1813	31	6 March 2021	<span><span></span></span> Cortina	DH 400				
1814	32	7 March 2021	<span><span></span></span> Cortina	SG 400				
1815	33	13 March 2021	<span><span></span></span> Kvitfjell	GS 400				
1816	34	14 March 2021	<span><span></span></span> Kvitfjell	SL 400				
1817	35	17 March 2021	<span><span></span></span> Kvitfjell	DH 400				
1818	36	19 March 2021	<span><span></span></span> Kvitfjell	SG 400				

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 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 <sup>4</sup> formatting.

Check that your script also works when accessing [https://en.wikipedia.org/wiki/2020%E2%80%9321\\_FIS\\_Alpine\\_Ski\\_World\\_Cup](https://en.wikipedia.org/wiki/2020%E2%80%9321_FIS_Alpine_Ski_World_Cup).

#### Files to Deliver in this Subtask

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.

<sup>4</sup><https://github.com/adam-p/markdown-here/wiki/Markdown-Cheatsheet#tables>



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](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.

There you will extract the names of the teams which made it to the conference semifinals (This is shown in the section "Bracket").<sup>5</sup>. 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.

---

<sup>5</sup>Hint: It should be 8.

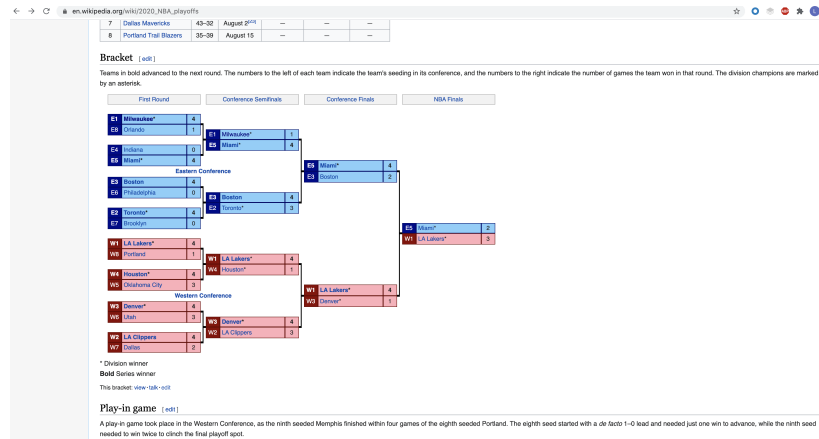


Figure 4: Bracket Section.

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

**2019–20 Milwaukee Bucks season**

From Wikipedia, the free encyclopedia

The **2019–20 Milwaukee Bucks season** was the 52nd season of the franchise in the National Basketball Association (NBA).<sup>[1]</sup> The Bucks entered the season following a playoff defeat in six games from the Toronto Raptors in the Eastern Conference Finals.

On February 20, 2020, the Bucks became the first team to clinch a playoff berth after the Chicago Bulls defeated the Washington Wizards.

The season was suspended by the league officials following the games of March 11<sup>[2]</sup> after it was reported that Rudy Gobert tested positive for COVID-19.<sup>[3]</sup>

The Bucks returned to play on July 31 against the Boston Celtics, at the ESPN Wide World of Sports Complex in Orlando, Florida.<sup>[4]</sup>

When the format of the Bubble announced, Bucks won the Central Division on June 4 and tied Detroit Pistons' record of 9 Central Division titles. After a 23-point comeback win over the Miami Heat, the Bucks clinched the top seed in the Eastern Conference for the second consecutive season. The Bucks faced the Orlando Magic in the first round.

In a significant development following the shooting of Jacob Blake in Kenosha, Wisconsin, the Bucks struck in solidarity with Black Lives Matter protests, refusing to play Game 5 of their series against the Magic on August 26.

On September 8, the Milwaukee Bucks season came to an end when the 5th seeded Miami Heat eliminated them in a five-game upset in the conference semifinals.

**Contents** [ show ]

- Draft picks
- Roster
- Standings
- 2.1 Division
- 3.2 Conference
- Game log
  - 1 Preseason
  - 2 Regular season
  - 3 Playoffs
- Transactions
  - 5.1 Overview
  - 5.2 Trades
  - 5.3 Free agency
  - 5.3.1 Additions
  - 5.3.2 Subtractions
- References

**Draft picks** [ edit ]

Main article: 2019 NBA draft

Round	Pick	Player	Position(s)	Nationality	School
1	30	Kenneth Farrier Jr.	Shooting Guard	<span><span></span></span> United States	USC

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:

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

← → en.wikipedia.org/wiki/2019-20\_Milwaukee\_Bucks\_season

### Roster

Roster listing

Milwaukee Bucks roster									
Players								Coaches	
Pos.	No.	Name	Height	Weight	DOB (YYYY-MM-DD)	From	Head coach		
F	34	Antetokounmpo, Giannis	6 ft 11 in (2.11 m)	242 lb (110 kg)	1994-12-06	Greece	• Mike Budenholzer		
F	43	Antetokounmpo, Thanasis	6 ft 6 in (1.98 m)	219 lb (99 kg)	1990-07-18	Greece	Assistant coach(es)		
G	6	Bledsoe, Eric	6 ft 5 in (1.96 m)	214 lb (97 kg)	1969-12-09	Kentucky	• Vin Baker		
G/F	23	Brown, Sterling	6 ft 5 in (1.96 m)	219 lb (99 kg)	1995-02-10	SMU	• Chad Foster		
G	24	Connaughton, Pat	6 ft 5 in (1.96 m)	209 lb (95 kg)	1993-01-06	Notre Dame	• Darvin Ham		
G	0	Dinwiddie, Dante	6 ft 4 in (1.93 m)	203 lb (92 kg)	1987-01-31	Villanova	• Charles Lee		
G	3	Hill, George	6 ft 3 in (1.91 m)	188 lb (85 kg)	1988-05-04	IUPUI	• Josh Longstaff		
F	7	Iyemori, Chase	6 ft 9 in (2.06 m)	235 lb (107 kg)	1987-05-15	Turkey	• Patrick St. Andrews		
G/F	26	Komer, Kyle	6 ft 7 in (2.01 m)	212 lb (96 kg)	1981-03-17	Cwrighton	• Ben Sullivan		
C	11	Lopez, Brook	7 ft 0 in (2.13 m)	282 lb (128 kg)	1988-04-01	Stanford	Legend		
C	42	Lopez, Robin	7 ft 0 in (2.13 m)	281 lb (127 kg)	1988-04-01	Stanford	• (C) Team captain		
G	15	Mason, Frank (TW)	5 ft 11 in (1.80 m)	180 lb (86 kg)	1994-04-03	Kansas	• (GP) Unsigned draft pick		
G	9	Mathews, Wesley	6 ft 4 in (1.93 m)	230 lb (105 kg)	1995-10-14	Marquette	• (FA) Free agent		
F	22	Middleton, Khris	6 ft 7 in (2.01 m)	202 lb (92 kg)	1991-08-12	Texas A&M	• (S) Suspended		
G/F	13	Reynolds, Cameron (TW)	6 ft 7 in (2.01 m)	235 lb (107 kg)	1995-02-07	Tulane	• (OL) On assignment to G League affiliate		
F	20	Williams, Marvin	6 ft 8 in (2.03 m)	237 lb (108 kg)	1986-06-19	North Carolina	• (TW) Two-way affiliate player		
F	5	Wilson, D. J.	6 ft 10 in (2.08 m)	231 lb (105 kg)	1996-02-19	Michigan	• Injured		

Standings

Division

Central Division	W	L	PCT	GB	Home	Road	Div	GP
<b>x - Milwaukee Bucks</b>	56	17	.767	0.0	30-8	26-12	13-1	73
x - Indiana Pacers	45	28	.616	11.0	25-11	20-17	8-7	73
Chicago Bulls	22	43	.338	30.0	14-20	8-23	7-9	65
Detroit Pistons	20	46	.303	32.5	11-22	9-24	5-10	66
Cleveland Cavaliers	19	46	.292	33.0	11-25	8-21	4-10	65

Figure 6: Roster Milwaukee Bucks.

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.

← → en.wikipedia.org/wiki/Giannis\_Antetokounmpo

### Legend

GP	Games played	GS	Games started	MPG	Minutes per game
FG%	Field goal percentage	SP%	Shooting percentage	FT%	Free throw percentage
REB	Rebounds per game	APG	Assists per game	STL	Steals per game
BPG	Blocks per game	PPG	Points per game	BPG	Blocks per game
				BPG	Blocks per game

### NBA

#### Regular season

Year	Team	GP	GS	MPG	FG%	SP%	FT%	REB	APG	BPG	PPG
2013-14	Milwaukee	77	23	24.0	.414	.247	.663	4.4	1.9	.8	6.8
2014-15	Milwaukee	81	71	31.4	.491	.199	.741	6.7	2.6	.9	10.0
2015-16	Milwaukee	80	79	35.3	.506	.257	.724	7.7	4.3	1.2	14.0
2016-17	Milwaukee	80	80	35.6	.522	.272	.770	8.7	5.4	1.6	19.0
2017-18	Milwaukee	75	75	36.7	.529	.307	.780	10.0	4.8	1.5	26.0
2018-19	Milwaukee	72	72	32.9	.576	.256	.729	12.5	5.8	1.3	27.7
2019-20	Milwaukee	65	65	30.4	.553	.304	.653	13.8	5.6	1.0	29.5
Career		528	483	32.5	.526	.284	.722	8.9	4.3	1.2	13.0
All-Star		4	4	28.8	.653	.231	.667	8.8	3.0	1.3	1.0

#### Playoffs

Year	Team	GP	GS	MPG	FG%	SP%	FT%	REB	APG	BPG	PPG
2015	Milwaukee	6	6	33.5	.366	.000	.739	7.0	2.7	.5	11.5
2017	Milwaukee	6	6	40.8	.536	.400	.543	9.5	4.0	2.2	17.0
2018	Milwaukee	7	7	40.0	.570	.386	.691	9.6	6.3	1.4	23.7
2019	Milwaukee	15	15	34.3	.492	.307	.637	12.3	4.9	1.1	28.5
2020	Milwaukee	9	9	30.8	.559	.325	.580	13.8	5.7	.7	26.7
Career		43	43	35.2	.514	.325	.627	11.0	4.8	1.1	15.0

### Awards and honours

- 2x NBA Most Valuable Player: 2019, 2020
- NBA Defensive Player of the Year: 2020
- 4x NBA All-Star: 2017, 2018, 2019, 2020
- NBA Most Improved Player: 2017
- 4x All-NBA Selection
- All-NBA First Team: 2019, 2020
- All-NBA Second Team: 2017, 2018

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<sup>6</sup>

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](https://en.wikipedia.org/wiki/2020_NBA_playoffs)
- 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`

---

<sup>6</sup>Since we take the 3 best players of each team, we will end up comparing the abilities of 24 players in total.

## 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](https://en.wikipedia.org/wiki/Parque_18_de_marzo_de_1938) to [https://en.wikipedia.org/wiki/Bill\\_Mundell](https://en.wikipedia.org/wiki/Bill_Mundell) using only urls in Wikipedia articles. The script should work with any wikipedia url. The websites given will be your test case. In order to assign a winner, we will evaluate this for 2 undisclosed urls. The most efficient (and correct) script will win a prize :)

**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!