Introduction
oooo

Code
oooooooo

Analysis
ooooooooo

Other
ooo

# NIL College Athlete Web Scraping Project

James Hamlette

The University of Nebraska at Omaha

ECON 8320: Graduate Tools for Data Analysis Term Project
May 11, 2023

Introduction
oooo

Code
oooooooo

Analysis
ooooooooo

Other
ooo

# Outline

Introduction
●○○○

Code
○○○○○○○○

Analysis
○○○○○○○○○

Other
○○○

# Introduction

## NIL Athlete Data Scraping

Project and Goal

- Scrape NIL endorsement/sponsor deals and related data
- Few websites list NIL deals and less provide usable data
- Quickly growing data and constantly changing

Using On3 and NIL College Athletes Websites

- https://www.on3.com/nil/rankings/player/nil-100/
- https://nilcollegeathletes.com/athletes

Introduction
○○●○

Code
○○○○○○○○○

Analysis
○○○○○○○○○○

Other
○○○

# Websites

Athletes | Universities | Companies ⌄ | Deals

**Athletes with Sponsorships and Endorsements**

List of student athletes getting sponsored and who they are represented by.

| NAME | SPONSORS | UNIVERSITY | SPORT | |
|------|----------|------------|-------|--|
| Jashon Hubbard | 614 Chiropractic<br>Barstool Sports<br>CRMD Ice Cream<br>Celsius<br>Chill Cryotherapy<br>Essentia Water<br>Ez Fresh Meals<br>Go Puff<br>Liquid I.V.<br>Max Effort Muscle<br>Playa Bowls<br>Rewild Yoga | The Ohio State University | Wrestling | More info › |
| Jonathan Shuskey | Barstool Sports<br>Liquid I.V.<br>Swing Juice<br>Talco Industrial Chemicals<br>The Winston Collection | Christian Brothers University | Golf | More info › |
| Collin Gillespie | Barstool Sports<br>Outback Steakhouse | Villanova University | Basketball | More info › |
| Buddy Boeheim | Enduraphin<br>Three Wishes Cereal | Syracuse University | Basketball | More info › |
| Armando Bacot | Jimmy's Seafood | University of North Carolina at Chapel Hill | Basketball | More info › |
| Aaron McLaughlin | Barstool Sports | North Carolina State University | Football | More info › |

Figure: NIL Athletes Website

Introduction
○○○●

Code
○○○○○○○○○

Analysis
○○○○○○○○○○

Other
○○○

# Websites Cont.



Figure: On3 Website

Introduction
○○○○

Code
●○○○○○○○

Analysis
○○○○○○○○○○

Other
○○○

Code

# Nil Deals Code

```
1  ###########ALL NIL College Athlete Website Scraping, Cleaning, and
       Analyzing Code################
2
3  #Load relevant libraries
4  import requests
5  from bs4 import BeautifulSoup
6  import numpy as np
7  import pandas as pd
8  import re
9  from urllib.parse import urljoin
10
11 #Define function that takes a url link and creates columns of data
       as follows
12 def collectNames(startURL):
13     myPage = requests.get(startURL)
14     parsed = BeautifulSoup(myPage.text)
15
16     #Start with the names of the athletes via tag "a"
17     a = parsed.find_all('td', class_="px-2 md:px-6 py-4 whitespace-
       nowrap text-sm font-medium text-gray-900")
18     n=[i.a.text.strip() for i in a]
19
20     #Append names to "ndata"
21     ndata=[]
22     for x in n:
23         ndata.append(x)
24     ndata=pd.DataFrame(ndata, columns=['Name'])
```

Introduction
○○○○

Code
○○●○○○○○

Analysis
○○○○○○○○○○

Other
○○○

## Nil Deals Code Cont.

```
1   #Append Sponsor to "t3"
2   t=[]
3   for i in a:
4       try:
5           d=i.find_next_sibling()
6           t.append(d.text)
7       except:
8           t.append("not listed")
9   t2=list(t)
10  t2new= [item.strip().replace('\n','') for item in t2]
11  t3=pd.DataFrame(t2new, columns=['Sponsors'])
12
13  #Append University and Sport to "datab"
14  b= parsed.find_all('span', class_="truncate")
15  blist= []
16  for i in b:
17      blist.append(i.text)
18  blist2=list(blist)
19  #Values are in succeeding positions, so create lists for every
       other to split
20  left = []
21  right = []
22  for i, j in enumerate(blist2):
23      if i%2==0:
24          left.append(j.strip())
25      else:
26          right.append(j.strip())
```

## Nil Deals Code part 3

```
1
2      #zip the two lists of Universities and Sports Together
3      b3 = list(zip(left, right))
4      b4=[list(i) for i in b3]
5      datab= pd.DataFrame(b4, columns=['University', 'Sport'])
6
7      #Inner join of ndata and datab
8      all_data= ndata.join([datab,t3])
9
10     #Parse through all remaining pages if there is one, then
        concatenate using recursive function
11     try:
12         nextPage= urljoin( 'https://nilcollegeathletes.com', parsed.
        find('div', class_="-mt-px flex w-0 flex-1 justify-end").a['href
        '])
13     except:
14         nextPage=None
15     if nextPage:
16         return pd.concat([all_data,collectNames(nextPage)], axis=0)
17     else:
18         return all_data
19 ##############################################END OF FUNCTION
        ##############################################
20 main_data=collectNames('https://nilcollegeathletes.com/athletes')
21 main_data
22 main_data.to_csv("main_data.csv")
```

Introduction
oooo

Code
ooooo●ooo

Analysis
oooooooooo

Other
ooo

# ON3 Top 100 Scraping

```
1  #################### ON3 TOP 100 OF ALL ATHLETES SCRAPING, CLEANING
      , ANALYSIS####################
2
3  #import relevant libraries
4  import numpy as np
5  import pandas as pd
6  import plotly.express as px
7  import requests
8  from bs4 import BeautifulSoup
9
10 #Define a scraping function to take a url link value
11 def onThree(scrapeurl):
12     myPage_ = requests.get("https://www.on3.com/nil/rankings/player/
        nil-100/")
13     soup = BeautifulSoup(myPage_.text)
14
15     #Scrape Names and append to dataframe
16     oo= soup.find_all('a', class_="MuiTypography-root MuiLink-root
        MuiLink-underlineNone NilPlayerRankingItem_name__nzSp9
        MuiTypography-h5 MuiTypography-colorPrimary")
17     oo=list(oo)
18
19     ood=[]
20     for o in oo:
21         oo2=o.text
22         ood.append(oo2)
23     ood=pd.DataFrame(ood, columns=['Name'])
```

## ON3 Top 100 Scraping Cont.

```
1  #Scrape the text/string number of followers
2      for o in oo:
3          uu= soup.find_all('p', class_="MuiTypography-root
       NilPlayerRankingItem_followersNumber__ifWQr MuiTypography-body1
       MuiTypography-colorTextPrimary")
4          uu=list(uu)
5      uud=[]
6      for u in uu:
7          try:
8              uu2=u.text
9              uud.append(uu2)
10         except:
11             uud.append("blank")
12     uud=pd.DataFrame(uud, columns=['Followers'])
13     #Scape the String of NIL Valuation
14     vv= soup.find_all('p', class_="MuiTypography-root
       NilPlayerRankingItem_valuationCurrency__oSkvo MuiTypography-
       body1 MuiTypography-colorTextPrimary")
15     vvd=[]
16     for v in vv:
17         vvs=v.text
18         vvd.append(vvs)
19     vvd=pd.DataFrame(vvd, columns=['Valuation'])
20     #Join the datasets
21     full_data=ood.join([uud,vvd])
22     #No further pages so just return
23     return full_data
```

Introduction
○○○○

Code
○○○○○○○●○○

Analysis
○○○○○○○○○○

Other
○○○

# Data Snapshot

```
1   collectNames('https://nilcollegeathletes.com/athletes')
```

|   | Name | Univeristy | Sport | Sponsors |
|---|------|-----------|-------|----------|
| 0 | Jashon Hubbard | The Ohio State University | Wrestling | 614 Chiropractic ... |
| 1 | Jonathan Shuskey | Christian Brothers University | Golf | Barstool Sports ... |
| 2 | Collin Gillespie | Villanova University | Basketball | Barstool Sports ... |
| 3 | Buddy Boeheim | Syracuse University | Basketball | Enduraphin ... |
| 4 | Armando Bacot | University of North Carolina at Chapel Hill | Basketball | Jimmy's Seafood |
| ... | ... | ... | ... | ... |

|   | Name | Followers | Valuation |
|---|------|-----------|-----------|
| 0 | Bronny James | 12.8M | $7.4M |
| 1 | Livvy Dunne | 11.3M | $3.5M |
| 2 | Arch Manning | 255K | $3.2M |
| 3 | Caleb Williams | 277K | $2.7M |
| 4 | Travis Hunter | 1.3M | $1.7M |
| ... | ... | ... | ... |
| 95 | Hailey Van Lith | 837K | $520K |
| 96 | Devin Leary | 37K | $519K |
| 97 | Flory Bidunga | 11.6K | $519K |
| 98 | Aaron Bradshaw | 8.4K | $516K |
| 99 | Dwight McGlothern | 31K | $513K |

Introduction
○○○○

Code
○○○○○○○●

Analysis
○○○○○○○○○○
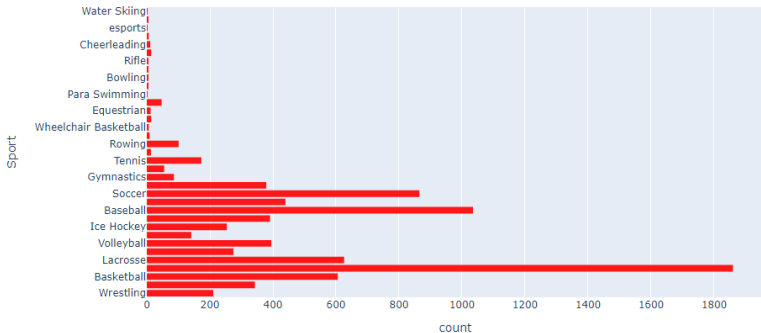
Other
○○○

# ON3 Top 100 Cleaning

```
1  ##Begin Extraction using defined function
2  tophundred=onThree("https://www.on3.com/nil/rankings/player/nil-100/
       ")
3  tophundred.to_csv("on3top100.csv")
4  #######################################CLEANING BELOW
       ########################################
5  #Valuation has "$", so remove to allow for quantitative analysis
6  tophundred['Valuation']=tophundred['Valuation'].str.replace('$','')
7  ###Values for thousands, millions etc are as "5K"; write function to
       convert to numeric values "5,000"
8  def value_change(num):
9      if num[-1:]=='K':
10         return float(num[:-1]) * 10**3
11     elif num[-1:]=='M':
12         return float(num[:-1]) * 10**6
13     elif num[-1:]=='B':
14         return float(num[:-1]) * 10**9
15     else:
16         num=float(num)
17 #Use value_change function to apply it to both followers and
       valuation columns without replacing original data
18 tophundred['Followers_total']=tophundred['Followers'].apply(
       value_change)
19 tophundred['Valuation_total']=tophundred['Valuation'].apply(
       value_change)
20 #Generate Rank variable based on the index
21 tophundred['Rank']=tophundred.index +1
```

Introduction
○○○○

Code
○○○○○○○○

Analysis
●○○○○○○○○○

Other
○○○

Analysis

Introduction
○○○○

Code
○○○○○○○○○

Analysis
○●○○○○○○○○

Other
○○○

# NIL Deals by Sport

```python
import plotly.express as px
import plotly.graph_objs as go
px.histogram(main_data, y="Sport", color_discrete_sequence=['red'],opacity=0.9, title="Sport Distribution among NIL D
```

Sport Distribution among NIL Deals



Football has the most deals by far; Barstool made up over 88% of NIL deals

Introduction
○○○○

Code
○○○○○○○○○

Analysis
○○●○○○○○○○

Other
○○○

# On3 Top 100 Athletes Valuation

```python
1  import plotly.express as px
2  import pandas as pd
3
4  px.scatter(tophundred, x='Valuation_total', y='Rank', title="Top 100 Athletes NIL Valuation by Rank")
5
```



Top 100 Athletes NIL Valuation by Rank

Introduction
○○○○

Code
○○○○○○○○

Analysis
○○○●○○○○○○

Other
○○○

# On3 Top 100 Athletes Summary

| | | |
|---|---|---|
| 1 | summary | |

| | Followers_total | Valuation_total |
|---|---|---|
| count | 100 | 100 |
| mean | 711,599 | 878,000 |
| std | 1,939,215 | 674,010 |
| min | 3,000 | 474,000 |
| 25% | 27,000 | 541,500 |
| 50% | 64,000 | 715,500 |
| 75% | 252,500 | 907,750 |
| max | 12,900,000 | 5,900,000 |

Introduction
○○○○

Code
○○○○○○○○○

Analysis
○○○○○●○○○○

Other
○○○

# On3 Top 100 Athletes Followers

```
1  import plotly.express as px
2  import pandas as pd
3
4  px.scatter(tophundred, x='Followers_total', y='Rank', title="Top 100 Athletes Social Media Followers by Rank")
5
```

Top 100 Athletes Social Media Followers by Rank

Introduction
○○○○

Code
○○○○○○○○○

Analysis
○○○○○●○○○

Other
○○○

# On3 Top 100 Football Positions

```
1  import plotly.express as px
2
3  px.bar(on3top100_football, x="Position",y="Valuation_total", title="NIL Valuation by Football Position")
```

NIL Valuation by Football Position

Introduction
○○○○

Code
○○○○○○○○○

Analysis
○○○○○○○●○○

Other
○○○

# On3 Top 100 Football Positions Cont.
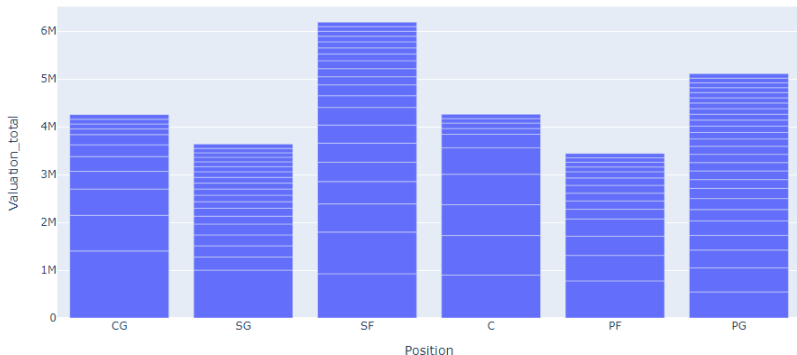
```python
1  table1=on3top100_football['Position'].value_counts()
2  table1
```

```
QB      24
RB      15
WR      12
EDGE    10
OT       8
CB       6
S        6
TE       5
LB       5
DL       5
IOL      4
```

Introduction
oooo

Code
ooooooooo

Analysis
oooooooo●o

Other
ooo

# On3 Top 100 Basketball Positions

```
1  import plotly.express as px
2  px.bar(on3top100_Basketball, x="Position",y="Valuation_total", title="NIL Valuation by Basketball Position")
3
```

NIL Valuation by Basketball Position

Introduction
0000

Code
000000000

Analysis
00000000●

Other
000

# On3 Top 100 Basketball Positions Cont.

```python
1  table2=on3top100_Basketball['Position'].value_counts()
2  table2
```

```
PG    26
SF    20
SG    19
PF    14
CG    11
C     10
```

Introduction
○○○○

Code
○○○○○○○○

Analysis
○○○○○○○○○○

Other
●○○

# Other

Introduction
oooo

Code
oooooooo

Analysis
ooooooooo

Other
o●o

## Not Working

- Twitter links proved difficult to scrape from NIL
- Values on NIL website were not readily available
- Load More Button is not very compatible with Beautiful Soup
- University and Individual Names using hyphens or apostrophes were unable to be recognized without correction

## Further Ideas

More Quantitative Values

- difficult to do thorough analysis without more than just values
- could look to pull further athlete data from roster websites or NCAA database to match with NIL valuation deals

This is a consistently updating list of deals so only pulling new values would be the most helpful moving forward (especially for On3)