


```
In [15]: # Return string within given tag
         soup.title.string
```

```
Out[15]: 'List of state and union territory capitals in India - Wikipedia'
```

```
In [17]: #We know that, we can tag a link using tag "<a>".
         #So, we should go with option soup.a and it should return the links available
         soup.a
```

```
Out[17]: <a id="top"></a>
```

```
In [21]: #Above, you can see that, we have only one output. Now to extract all the link

all_link =soup.findAll("a")
all_link
```

```
Out[21]: <a id="top"></a>,
  <a class="mw-jump-link" href="#mw-head">Jump to navigation</a>,
  <a class="mw-jump-link" href="#p-search">Jump to search</a>,
  <a href="/wiki/States_and_union_territories_of_India" title="States and un
ion territories of India">States and union <br/> territories of India</a>,
  <a class="image" href="/wiki/File:Flag_of_India.svg"></a>,
  <a href="/wiki/List_of_states_and_union_territories_of_India_by_area" titl
e="List of states and union territories of India by area">Area</a>,
  <a href="/wiki/List_of_states_and_union_territories_of_India_by_populatio
n" title="List of states and union territories of India by population">Popu
lation</a>,
  <a href="/wiki/List_of_Indian_states_and_union_territories_by_GDP" title
="List of Indian states and union territories by GDP">GDP</a>,
  <a href="/wiki/List_of_Indian_states_and_union_territories_by_GDP">GDP</a>
```

```
In [23]: #href is present in 'a tag'
#The href attribute specifies the URL of the page the link goes to.
#If the href attribute is not present, the <a> tag is not a hyperlink

for link in all_link:
    print(link.get('href'))
```

```
None
#mw-head
#p-search
/wiki/States_and_union_territories_of_India
/wiki/File:Flag_of_India.svg
/wiki/List_of_states_and_union_territories_of_India_by_area
/wiki/List_of_states_and_union_territories_of_India_by_population
/wiki/List_of_Indian_states_and_union_territories_by_GDP
/wiki/List_of_Indian_states_and_union_territories_by_GDP_per_capita
/wiki/ISO_3166-2:IN
None
/wiki/List_of_Indian_states_by_Child_Nutrition
/wiki/List_of_states_and_union_territories_of_India_by_crime_rate
/wiki/List_of_states_and_union_territories_of_India_by_households_having_electricity
/wiki/List_of_states_and_union_territories_of_India_by_fertility_rate
/wiki/Forest_cover_by_state_in_India
/wiki/Ease_of_doing_business_ranking_of_states_of_India
/wiki/List_of_Indian_states_and_territories_by_highest_point
```

```
In [26]: #Find the right table: As we are seeking a table to extract information about
#we should identify the right table first. Let's write the command to extract

all_tables=soup.find_all('table')
all_tables
```

```
Out[26]: [<table class="vertical-navbox nowraplinks" style="float:right;clear:right;
width:22.0em;margin:0 0 1.0em 1.0em;background:#f9f9f9;border:1px solid #aaa;padding:0.2em;border-spacing:0.4em 0;text-align:center;line-height:1.4em;
font-size:88%"><tbody><tr><th style="padding:0.2em 0.4em 0.2em;font-size:145%;line-height:1.2em"><a href="/wiki/States_and_union_territories_of_India"
title="States and union territories of India">States and union territories of India</a> <br/> ordered by</th></tr><tr><td style="padding:0.2em 0.4em"><div class="center"><div class="floatnone"><a class="image" href="/wiki/File:Flag_of_India.svg"></a></div></div></td></tr><tr><td class="hlist" style="padding:0 0.1em 0.4em">
<ul><li><a href="/wiki/List_of_states_and_union_territories_of_India_by_area" title="List of states and union territories of India by area">Area</a>
</li>
```

In [41]: *#Now to identify the right table, we will use attribute "class" of table and u*
#In chrome, you can check the class name by right click on the required table
#Inspect element -> Copy the class name OR go through the output of above comm

```
soup.find_all('table', class_="wikitable sortable plainrowheaders")
```

Out[41]: `<table class="wikitable sortable plainrowheaders">`
`<tbody><tr>`
`<th scope="col">No.`
`</th>`
`<th scope="col">State or
union territory`
`</th>`
`<th scope="col">Administrative capital`
`</th>`
`<th scope="col">Legislative capital`
`</th>`
`<th scope="col">Judicial capital`
`</th>`
`<th scope="col">Year of establishment`
`</th>`
`<th scope="col">Former capital`
`</th></tr>`
`<tr>`
`<td>1`
`</td>`

In [42]: *#Extract the information to DataFrame: Here, we need to iterate through each r*
#to a variable and append it to a List. Let's first look at the HTML structure
#(I am not going to extract information for table heading <th>)

```
right_table = soup.find('table' , {"class" : "wikitable sortable plainrowheade
right_table
```

Out[42]: `<table class="wikitable sortable plainrowheaders">`
`<tbody><tr>`
`<th scope="col">No.`
`</th>`
`<th scope="col">State or
union territory`
`</th>`
`<th scope="col">Administrative capital`
`</th>`
`<th scope="col">Legislative capital`
`</th>`
`<th scope="col">Judicial capital`
`</th>`
`<th scope="col">Year of establishment`
`</th>`
`<th scope="col">Former capital`
`</th></tr>`
`<tr>`
`<td>1`
`</td>`

In [47]: *#Above, you can notice that second element of <tr> is within tag <th> not <td>
#Now to access value of each element, we will use "find(text=True)" option with*

```
#Generate Lists
A=[]
B=[]
C=[]
D=[]
E=[]
F=[]
G=[]
for row in right_table.findAll("tr"):
    cells = row.findAll('td')
    states=row.findAll('th') #To store second column data
    if len(cells)==6: #Only extract table body not heading
        A.append(cells[0].find(text=True))
        B.append(states[0].find(text=True))
        C.append(cells[1].find(text=True))
        D.append(cells[2].find(text=True))
        E.append(cells[3].find(text=True))
        F.append(cells[4].find(text=True))
        G.append(cells[5].find(text=True))
```

In [48]: *#import pandas to convert List to data frame*

```
import pandas as pd
df=pd.DataFrame(A,columns=['Number'])
df['State/UT']=B
df['Admin_Capital']=C
df['Legislative_Capital']=D
df['Judiciary_Capital']=E
df['Year_Capital']=F
df['Former_Capital']=G
df
```

Out[48]:

	Number	State/UT	Admin_Capital	Legislative_Capital	Judiciary_Capital	Year_Capital
0	1	Andaman and Nicobar Islands	Port Blair	—	Kolkata	1956
1	2	Andhra Pradesh	Hyderabad	Amaravati	Amaravati	1956
2	3	Arunachal Pradesh	Itanagar	Itanagar	Guwahati	1956
3	4	Assam	Dispur	Guwahati	Guwahati	1956
4	5	Bihar	Patna	Patna	Patna	1956
5	6	Chandigarh	Chandigarh	—	Chandigarh	1966
6	7	Chhattisgarh	Raipur	Raipur	Bilaspur	2000
7	8	Dadra and Nagar Haveli	Silvassa	—	Mumbai	1954

In []:

In []:

In []:

In []:

In []:

In []:

In []: