Web Scraping of AmbitionBox website

```
In [1]:
            import numpy as np
            import pandas as pd
            import requests
            from bs4 import BeautifulSoup
In [2]:
         requests.get('https://www.ambitionbox.com/list-of-companies?page=1')
   Out[2]: <Response [403]>
            headers={'User-Agent':'Mozilla/5.0 (Windows NT 6.3; Win 64; x64) Apple
In [3]:
            webpage= requests.get('https://www.ambitionbox.com/list-of-companies?pd
In [4]:
In [5]:
            webpage
   Out[5]: '<!doctype html>\n<html data-n-head-ssr lang="en" data-n-head="%7</pre>
            B%22lang%22:%7B%22ssr%22:%22en%22%7D%7D">\n <head >\n
            arset="UTF-8">\n
                                <meta name="viewport" content="width=device-wi</pre>
            dth,initial-scale=1,minimum-scale=1">\n
                                                         <meta http-equiv="X-UA-</pre>
            Compatible content="IE=edge"> \n
                                                 <title>List of companies - 79
            Øk companies | AmbitionBox</title><meta data-n-head="ssr" name="co</pre>
            pyright" content="2023 AmbitionBox"><meta data-n-head="ssr" name</pre>
            ="revisit-after" content="1 day"><meta data-n-head="ssr" name="app
            lication-name" content="AmbitionBox"><meta data-n-head="ssr" name</pre>
            ="content-language" content="EN"><meta data-n-head="ssr" name="goo
            gle-signin-client id" content="462822053404-hphug4pkahqljh2tc96g35
            at47o4isv2.apps.googleusercontent.com"><meta data-n-head="ssr" pro
            perty="fb:app_id" content="712617688793459"><meta data-n-head="ss</pre>
            r" name="theme-color" content="#ffffff"><meta data-n-head="ssr" na
            me="msapplication-navbutton-color" content="#ffffff"><meta data-n-</pre>
            head="ssr" name="apple-mobile-web-app-status-bar-style" content="d
            efault"><meta data-n-head="ssr" property="fb:admins" content="1000
            01438127755,100000444923785"><meta data-n-head="ssr" property="og:
            site_name" content="AmbitionBox"><meta data-n-head="ssr" property</pre>
```

```
In [6]:
            soup = BeautifulSoup(webpage, 'lxml')
            soup
   Out[6]: <!DOCTYPE html>
            <html data-n-head="%7B%22lang%22:%7B%22ssr%22:%22en%22%7D%7D" data
            -n-head-ssr="" lang="en">
            <head>
            <meta charset="utf-8"/>
            <meta content="width=device-width,initial-scale=1,minimum-scale=1"</pre>
            name="viewport"/>
            <meta content="IE=edge" http-equiv="X-UA-Compatible"/>
            <title>List of companies - 790k companies | AmbitionBox</title><me
            ta content="2023 AmbitionBox" data-n-head="ssr" name="copyright"/>
            <meta content="1 day" data-n-head="ssr" name="revisit-after"/><met</pre>
            a content="AmbitionBox" data-n-head="ssr" name="application-name"/
            ><meta content="EN" data-n-head="ssr" name="content-language"/><me</pre>
            ta content="462822053404-hphug4pkahqljh2tc96g35at47o4isv2.apps.goo
            gleusercontent.com" data-n-head="ssr" name="google-signin-client_i
            d"/><meta content="712617688793459" data-n-head="ssr" property="f
            b:app id"/><meta content="#ffffff" data-n-head="ssr" name="theme-c</pre>
            olor"/><meta content="#ffffff" data-n-head="ssr" name="msapplicati
            on-navbutton-color"/><meta content="default" data-n-head="ssr" nam
         ▶ soup.find_all('h1')[0].text
In [7]:
```

Out[7]: 'List of companies in India'

```
  | for i in soup.find_all('h2'):

 In [8]:
                 print(i.text.strip())
             TCS
             Accenture
             Cognizant
             ICICI Bank
             HDFC Bank
             Wipro
             Infosys
             Capgemini
             Tech Mahindra
             Genpact
             HCLTech
             Amazon
             Axis Bank
             Concentrix Corpo...
             IBM
             Reliance jio
             Larsen & Toubro ...
             HDB Financial Se...
             Vodafone Idea
             Teleperformance
             Reliance Retail
             Kotak Mahindra B...
             Reliance Industr...
             Deloitte
             Bharti Airtel
             BYJU'S
             Tata Motors
             Flipkart
             WNS
             IndusInd Bank
 In [9]:
             company = soup.find_all('div', class_ = 'company-content-wrapper')
In [10]:
             company
   Out[10]: [<div class="company-content-wrapper"><div class="company-content">
             t"><div class="company-logo"><img alt="Tata Consultancy Services 1
             ogo" class="lazy" data-src="https://static.ambitionbox.com/assets/
             v2/images/rs:fit:200:200:false:false/bG9jYWw6Ly8vbG9nb3Mvb3JpZ2luY
             WxzL3Rjcy5qcGc.webp" height="100" onerror="this.onerror=null; this.
             src='/static/icons/company-placeholder.svg';" src="https://static.
             ambitionbox.com/static/icons/company-placeholder.svg" width="100"/
             ></div> <div class="company-info-wrapper"><div class="company-inf</pre>
             o"><div class="left"><a href="/overview/tcs-overview"><h2 class="c
             ompany-name bold-title-1" title="TCS">
             TCS
             h2></a> <div class="rating-wrapper"><p class="rating badge-large r
             ating-35"><i class="icon icon-star"></i>
             3.9
                                                                                </
             p> <a class="review-count sbold-Labels" href="https://www.ambition</pre>
             harr and harris are 14 and harris are all s
```

```
In [11]: ▶ len(company)
```

Out[11]: 30

```
TCS ** 3.9 (53.3k Reviews)

#6 Best Mega Company - 2021 +1 more

Public

Public

Stances

Mumbai,Maharashtra + 275 more

1 Lakh+ Employees (India)

BPO IT Services & Consulting Forbes Global 2000

Fortune India 500 Public Mumbai,Maharashtra
```

```
In [12]:
                                                  name = []
                                                  rating = []
                                                  reviews = []
                                                  type_of_comp = []
                                                  headquarter = []
                                                  how_old = []
                                                  no\_of\_emp = []
                                                  for i in company:
                                                                  name.append(i.find('h2').text.strip())
                                                                  rating.append(i.find('p', class_='rating').text.strip())
                                                                  reviews.append(i.find('a', class_='review-count').text.strip())
                                                                  type_of_comp.append(i.find_all('p', class_='infoEntity')[0].text.st
                                                                  headquarter.append(i.find_all('p', class_='infoEntity')[1].text.str
                                                                  how_old.append(i.find_all('p', class_='infoEntity')[2].text.strip()
                                                                  no_of_emp.append(i.find_all('p', class_='infoEntity')[3].text.strip
                                                  df = pd.DataFrame({'Company Name': name, 'Company Rating' : rating, 'C
                                                                                                                          'type_of_comp': type_of_comp, 'headquarter' : headqua
                                                                                                                         'no_of_emp': no_of_emp})
```

In [13]: ► df

Out[13]:		Company Name	Company Rating	Company Reviews	type_of_comp	headquarter
	0	TCS	3.9	(53.3k Reviews)	Public	Mumbai,Maharashtra + 275 more
	1	Accenture	4.1	(33.8k Reviews)	Public	Dublin + 133 more
	2	Cognizant	4.0	(31.1k Reviews)	Private	Teaneck. New Jersey. + 102 more
	3	ICICI Bank	4.0	(35k Reviews)	Public	Mumbai,Maharashtra + 1072 more
	4	HDFC Bank	4.0	(41.3k Reviews)	Public	Mumbai,Maharashtra + 1295 more
	5	Wipro	3.9	(30.4k Reviews)	Public	Bangalore/Bengaluru,Karnataka + 229 more
	6	Infosys	3.9	(29.6k Reviews)	Public	Bengaluru/Bangalore,Karnataka + 127 more
	7	Capgemini	3.9	(25.1k Reviews)	Public	Paris + 79 more
	8	Tech Mahindra	3.7	(22.6k Reviews)	Public	Pune,Maharashtra + 213 more
	9	Genpact	4.0	(20.3k Reviews)	Public	New York,New York + 69 more
1	10	HCLTech	3.8	(19.5k Reviews)	Public	Noida,Uttar Pradesh + 141 more
	11	Amazon	4.2	(19.8k Reviews)	Public	Seattle,Washington + 387 more
1	12	Axis Bank	3.9	(19.3k Reviews)	Public	Mumbai,Maharashtra + 1123 more
1	13	Concentrix Corpo	4.0	(15.6k Reviews)	Public	Fremont,California + 70 more
1	14	IBM	4.2	(15.8k Reviews)	Public	Armonk,New York + 121 more
1	15	Reliance jio	4.0	(15.4k Reviews)	Public	Navi Mumbai,Maharashtra + 1019 more
•	16	Larsen & Toubro	4.1	(25.1k Reviews)	Public	Mumbai,Maharashtra + 536 more
,	17	HDB Financial Se	4.0	(13.9k Reviews)	Private	Ahmedabad,Gujrat + 812 more

	Company Name	Company Rating	Company Reviews	type_of_comp	headquarte
18	Vodafone Idea	4.2	(13.5k Reviews)	Public	Gandhinagar,Gujrat + 578 more
19	Teleperformance	3.6	(15.1k Reviews)	Private	Paris + 122 mor
20	Reliance Retail	4.1	(17.8k Reviews)	Private	Navi Mumbai,Maharashtra 739 mor
21	Kotak Mahindra B	3.9	(14k Reviews)	Public	Mumbai,Maharashtra + 48 mor
22	Reliance Industr	4.1	(46.9k Reviews)	Public	Navi Mumbai,Maharashtra 519 mor
23	Deloitte	4.1	(10.8k Reviews)	Private	New York,New York + 132 mor
24	Bharti Airtel	4.1	(12.8k Reviews)	Public	Gurgaon/Gurugram,Haryana 540 mor
25	BYJU'S	3.4	(14k Reviews)	Private	Bangalore,Karnataka + 26 mor
26	Tata Motors	4.1	(12.3k Reviews)	Public	Pune,Maharashtra + 394 mor
27	Flipkart	4.2	(12.5k Reviews)	Public	Bangalore,Karnataka + 46 mor
28	WNS	3.7	(7.2k Reviews)	Private	Mumbai,Maharashtra + 2 mor
29	IndusInd Bank	3.8	(6.9k Reviews)	Public	Gurgaon/Gurugram,Haryana 578 mor

In [14]:

Out[14]: (30, 7)

Extract all details from all 333 pages

```
In [15]:  | final = pd.DataFrame()
             for j in range(1,334):
                 headers={'User-Agent':'Mozilla/5.0 (Windows NT 6.3; Win 64; x64) A
                 url = requests.get('https://www.ambitionbox.com/list-of-companies?p
                 webpage= requests.get('https://www.ambitionbox.com/list-of-companie
                 soup = BeautifulSoup(webpage, 'lxml')
                 company = soup.find_all('div', class_ = 'company-content-wrapper')
                 name = []
                 rating = []
                 reviews = []
                 type_of_comp = []
                 headquarter = []
                 how_old = []
                 no\_of\_emp = []
                 for i in company:
                     try:
                         name.append(i.find('h2').text.strip())
                         name.append(np.nan)
                     try:
                         rating.append(i.find('p', class_='rating').text.strip())
                     except:
                         rating.append(np.nan)
                     try:
                         reviews.append(i.find('a', class_='review-count').text.stri
                     except:
                         reviews.append(np.nan)
                     try:
                         type_of_comp.append(i.find_all('p', class_='infoEntity')[0]
                     except:
                         type_of_comp.append(np.nan)
                     try:
                         headquarter.append(i.find_all('p', class_='infoEntity')[1].
                     except:
                         headquarter.append(np.nan)
                     try:
                         how_old.append(i.find_all('p', class_='infoEntity')[2].text
                     except:
                         how_old.append(np.nan)
                         no_of_emp.append(i.find_all('p', class_='infoEntity')[3].te
                     except:
                         no of emp.append(np.nan)
                 df = pd.DataFrame({'Company Name': name, 'Company Rating' : rating,
                                'type_of_comp': type_of_comp, 'headquarter' : headqua
                                'no_of_emp': no_of_emp})
                 final = final.append(df, ignore index=True)
```

C:\Users\Lenovo\AppData\Local\Temp\ipykernel_19552\825628277.py:52: F
utureWarning: The frame.append method is deprecated and will be remov
ed from pandas in a future version. Use pandas.concat instead.
 final = final.append(df, ignore index=True)

▶ final.tail() In [16]: Out[16]: Company Company Company type_of_comp headquarter how_olc Name Rating **Reviews** St. (53 Gurgaon/Gurugram 11-50 Employees 9979 Andrews 4.5 NaN Reviews) + 2 more (India) Inst... Navi Sarla (53 24 years 9980 4.1 Private Mumbai, Maharashtra Advantech Reviews) olo + 5 more Gurgaon/Gurugram Pathways (53 9981 4.4 NaN NaN Schools Reviews) + 2 more (53 Noida, Uttar Pradesh 23 years Vega 9982 Private 4.1 Industries Reviews) + 9 more olc 501-11 Or Yehuda + 7 Tahal (53 9983 4.3 71 years old Employees Group Reviews) more (Global In [17]: final.shape (9984, 7)Out[17]: In []: M