Algorithm-

START

1. Goto https://www.lib.lsu.edu/
2. Select Databases tab and then 'w' from Browse Databases.
3. Click on the 'Web of Science' link.
   Note – Instead of above 3 steps, we can directly go to the link -
   https://sso.paws.lsu.edu/login?service=https%3A%2F%2Fwebauth.shib.lsu.edu%2Fidp%2FAuthn%2FExtCas%3Fconversation%3De1s1&entityId=https%3A%2F%2Flibezp.lib.lsu.edu%2Fshibboleth

4. Login with the username and password.
5. Select the advanced search tab and write the query.
   For example – if we want the institutions that have cited 'MIT' in the years from 1998 to 2018 in the field of Astronomy and Astrophysics, then search for 'MIT' after selecting 'OG' field tag, add the required institution (in our case its MIT) and then click ok. You will be redirected to the query page with the query 'OG=(Massachusetts Institute of Technology (MIT))'. Now sleect the topic from the WC= Web of Science Category field tag (in our case its Astronomy & Astrophysics) and write a complete query as -
   'OG=(Massachusetts Institute of Technology (MIT)) and WC=Astronomy & Astrophysics'
   Click on search.
   Note that time span is selected to be 'All years (1900 – 2018)'
6. Now click on the results in the search history.
7. After this refine the search by selecting the publication year one by one.
   This can be done by two ways-
    a) Select the publication year from the left menu and click refine
    OR
    b) Goto the link http://apps.webofknowledge.com.libezp.lib.lsu.edu/RAMore.do?product=WOS&search_mode=AdvancedSearch&SID=7BVklzaxLTCSPkDWooE&qid=1&ra_mode=more&ra_name=PublicationYear&colName=WOS&viewType=raMore
   and select years one by one and click refine.
8. Once you get the refined results, click on the 'Create Citation Report' option on the right side  menu.
9. Select 'Analyze' under 'Citing articles', then select the 'Organizations-Enhanced' option in the left menu. Select the 'Sort by **Record count**' as 'Show **500**'.
10. Copy the url and run the script scraper.py in the terminal with arguments as – url , name of institute you are scraping information for, year for which you are scraping information for.
    For example -
    python3 scraper.py "http://wcs.webofknowledge.com/RA/analyze.do?product=WOS&SID=7BVklzaxLTCSPkDWooE&field=OG_OrgEnhancedName_OrgEnhancedName_en&yearSort=false" MIT '2018'

    The above command will make a folder by the name 'MIT' in the current directory and make a file 2018.json with the citing institutes list and no of citations to MIT in that particular year.

11. The scraper.py is as follows -
    1) import os , sys
    2) from selenium  import webdriver
    3) from selenium.webdriver.common.keys import Keys
    4) import re

```
5)  from bs4 import BeautifulSoup
6)  from selenium.webdriver.common.by import By
7)  from selenium.webdriver.support.ui import WebDriverWait
8)  from selenium.webdriver.support import expected_conditions as EC
9)  import json
10) import errno
11) browser = webdriver.Firefox()
12) str1=[]
13) #yearList=['2018','2017','2016','2015,2014,2013,2012,2011,2010,2009,2008,2007
    ,2006,2005,2004,2003,2002,2001,2000,1999,1998,1997,1996,1995]
14) def start():
15)         name=sys.argv[2]
16)         make_directory(name)
17)         url=sys.argv[1]
18)         print(url)
19)         browser.get(url)
20)         elem =
    WebDriverWait(browser,10).until(EC.presence_of_element_located((By.XPATH,
    '/html/body/div[2]/div[6]/form/div/div[2]/div[2]/div/div[2]/div[4]/table')))
21)         source = browser.page_source
22)         scrape(source,name)
23)
24) def scrape(source,name):
25)         soup = BeautifulSoup(source,'html.parser')
26)
27)         file=open('./'+name+'/'+sys.argv[3]+'.json',"w")
28)         str1=[]
29)         str1.append("{")
30)         for tr in soup.findAll("tr",{"class":"RA-NEWRAresultsEvenRow"}):
31)                 tds = tr.find_all('td')
32)                 str1.append("\"")
33)                 str1.append(tds[1].text)
34)                 str1.append("\"")
35)                 str1.append(":")
36)                 str1.append("\"")
37)                 str1.append(tds[2].text.strip())
38)                 str1.append("\"")
39)                 str1.append(",")
40)         for tr in soup.findAll("tr",{"class":"RA-NEWRAresultsOddRow"}):
41)                 tds = tr.find_all('td')
42)                 str1.append("\"")
43)                 str1.append(tds[1].text)
44)                 str1.append("\"")
45)                 str1.append(":")
46)                 str1.append("\"")
47)                 str1.append(tds[2].text.strip())
48)                 str1.append("\"")
49)                 str1.append(",")
50)         var = ''.join(str1)
51)         file.write(var)
52)         file.close()
53)         correct_json_data(name)
```

```
54)
55) def correct_json_data(name):
56)         instituteFileRead=open('./'+name+'/'+sys.argv[3]+'.json',"r")
57)         finalstr=(instituteFileRead.read())
58)                finalstr = finalstr[:-1]
59)         finalstr+="}"
60)         instituteFileRead.close()
61)         instituteFileWrite=open('./'+name+'/'+sys.argv[3]+'.json',"w")
62)         var = ''.join(finalstr)
63)         instituteFileWrite.write(var)
64)         instituteFileWrite.close()
65)
66) def make_directory(name):
67)     path = "./"+name
68)     os.makedirs(path, exist_ok=True)
69)     print ("Path is created")
70)
71) start()
```

12. Now to get the data for the remaining years follow these steps recusively,
    a)Goto the link http://apps.webofknowledge.com.libezp.lib.lsu.edu/RAMore.do?product=WOS&search_mode=AdvancedSearch&SID=7BVklzaxLTCSPkDWooE&qid=1&ra_mode=more&ra_name=PublicationYear&colName=WOS&viewType=raMore
    b)Select the year and click refine.
    c)Once you get the refined results, click on the 'Create Citation Report' option on the right side  menu. Select 'Analyze' under 'Citing articles'.
    d)Now execute the same previous command with only difference in the 3rd argument, changing it to the year the results are refined in. Note – the same url can be used once the result is refined manually, under one session.

13. After executing step 12 for all the years we want data for, we will get files 2018.json, 2017.json,2016.json,.... in the MIT directory inside the current directory.

14. Run the script test.py in the same directory
    python3 test.py

    The above command will create a directory "CitingInstitutes" inside "MIT" in the current - directory, with InstituteName.json files which have the count of the no. of citations from the year 1998 to 2018 for the InstituteName Institute to MIT.

15. The test.py is as follows
    ```
    1)  import os , sys
    2)  import json
    3)  from pprint import pprint
    4)
    5)  def make_directory():
    6)      path = "./MIT/CitingInstitutes"
    7)      os.makedirs(path, exist_ok=True)
    8)      print ("Path is created")
    9)
    10) def makeitwork():
    11)         institutesNameList=['UNIVERSITY OF CALIFORNIA
            SYSTEM','MASSACHUSETTS INSTITUTE OF TECHNOLOGY
            MIT','UNITED STATES DEPARTMENT OF ENERGY DOE','ISTITUTO
            NAZIONALE DI FISICA NUCLEARE','HARVARD
    ```

UNIVERSITY','CHINESE ACADEMY OF SCIENCES','CONSEJO SUPERIOR DE INVESTIGACIONES CIENTIFICAS CSIC','UNIVERSITY OF CHICAGO','UNIVERSITE PARIS SACLAY','RUSSIAN ACADEMY OF SCIENCES','CEA','SORBONNE UNIVERSITE','UNIVERSITY OF CALIFORNIA BERKELEY','CNRS NATIONAL INSTITUTE OF NUCLEAR PARTICLE PHYSICS IN2P3','GODDARD SPACE FLIGHT CENTER','SAPIENZA UNIVERSITY ROME','UNIVERSITE SORBONNE PARIS CITE USPC COMUE','UNIVERSITY OF MICHIGAN','OHIO STATE UNIVERSITY','JOHNS HOPKINS UNIVERSITY','EUROPEAN ORGANIZATION FOR NUCLEAR RESEARCH CERN','ISTITUTO NAZIONALE GEOFISICA E VULCANOLOGIA INGV','UNIVERSITY OF TORONTO','UNIVERSITY OF OXFORD','DEUTSCHES ELEKTRONEN SYNCHROTRON DESY','STANFORD UNIVERSITY','UNIVERSITY OF WASHINGTON','UNIVERSITY OF WISCONSIN MADISON','NORTHWESTERN UNIVERSITY','INDIAN INSTITUTE OF TECHNOLOGY SYSTEM IIT SYSTEM','JOINT INSTITUTE FOR NUCLEAR RESEARCH RUSSIA','UNIVERSITY OF LONDON','UNIVERSITY OF BIRMINGHAM','CALIFORNIA INSTITUTE OF TECHNOLOGY','PENNSYLVANIA COMMONWEALTH SYSTEM OF HIGHER EDUCATION PCSHE']

```
12)      stringToWriteInFile=[]
13)      for i in range (1998,2019):
14)          make_directory()
15)          yearFile=open('./MIT/'+str(i)+'.json',"r")
16)          str1=(yearFile.read())
17)          with open('./MIT/'+str(i)+'.json') as f:
18)              data = json.load(f)
19)          for j in range (35):
20)              try:
21)                  if (i==1998):
22)                      stringToWriteInFile.append("{")
23)
         instituteFile=open('./MIT/CitingInstitutes/'+institutesNameList[j]
   +'.json',"a")
24)                  stringToWriteInFile.append("\""+str(i)
   +"\""+":"+"\""+data[institutesNameList[j]]+"\",")
25)                  var = ''.join(stringToWriteInFile)
26)                  instituteFile.write(var)
27)                  instituteFile.close()
28)              except:
29)                  print ("Exception handled")
30)              stringToWriteInFile=[]
31)
32)      """    In the code below I have written the code to automatically
   remove the last comma ',' in the json data and place a '}' in place
         of it using some of the string operations
33)      """
34)
35)      for i in range (35):    #35 because I have considered only 35
   institutes as of now to get the no of citations from them to MIT and the
                      #35 institutes are mentioned in the string
   institutesNameList
```

```
36)        instituteFileRead=open('./MIT/CitingInstitutes/'+institutesNameList[i]
   +'.json',"r")
37)            finalstr=(instituteFileRead.read())
38)            print (finalstr)
39)            finalstr = finalstr[:-1]
40)            finalstr+="}"
41)            instituteFileRead.close()
42)

   instituteFileWrite=open('./MIT/CitingInstitutes/'+institutesNameList[i]
   +'.json',"w")
43)            var = ''.join(finalstr)
44)            instituteFileWrite.write(var)
45)            instituteFileWrite.close()
46) makeitwork()
```

16. Note that only the data of the institutes mentioned in the list by the name 'institutesNameList' in the test.py code will be made. Though we have data for 500 institutes per year but here we are only mentioning 35 institute names in the list. Also the variable i of the second loop is hard coded to be 35 i.e. length of string.

17. Once test.py is run, we will get data in json files. Now it can be converted to csv files by the following command -
Example-
        json2csv -i STANFORD\ UNIVERSITY.json  -o STANFORD\ UNIVERSITY.csv

The csv files will be stored in the current directory.


The data is ready to be used for models.


END