

Selected 3 Project

there is two phases into this project

phase 1

Project Requirements :

- 1- Make scrapping on the website and collect data
 - 2- clean the scrapped data
 - 3- auto correct the scrapped data
 - 4- Save your data in csv file "Data sheet"
- with the columns like " Author - title - publishing year - link for pdf - etc. "
- 5- cluster this data for topic modeling

The website for scrapping : <https://www.arab-books.com>

importing libraries

```
In [1]: from bs4 import BeautifulSoup
import requests
import csv
from itertools import zip_longest
import re
import pandas as pd
import string
import nltk
from nltk.tokenize import word_tokenize , sent_tokenize
import numpy as np
```

1- Make scrapping on the website: www.arab-books.com and collect data

Request to get all books from the website

there are **200 pages** on the website Each page have on average **30 book**

```
In [2]: main_src = []
for i in range(1,201):
    main_src.append(requests.get("https://www.arab-books.com/page/"+str(i)).text)
```

```
In [3]: main_soup = []
for src in main_src:
    main_soup.append(BeautifulSoup(src, 'lxml'))
```

```
In [4]: all_books_src = []
for soup in main_soup :
    all_books_src = all_books_src + soup.find_all('li',class_='post-item tie-standard')
```

```
In [5]: all_books_title=[]
for soup in main_soup :
    all_books_title = all_books_title + soup.find_all("div",{"class":"excerpt-book"})
```

```
In [6]: all_books_title[0:10]
```

```
In [7]: books_titles=[]
        for title in all_books_title :
            books_titles.append(title.text)
        len(books_titles)
```

Out[7]: 5984

```
In [8]: url_pattern = "http(^[^"]+)"
        URL_RegEx = re.compile(url_pattern)
        Books_urls=[]
        for Object in all_books_title :
            myResult = URL_RegEx.search(str(Object))
            Books_urls.append(myResult.group())
        Books_urls[:10]
```

```
Out[8]: ['https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d8%a7d9%84%d8%aa%d8%b6%d8%ad%d9%8a%d8%a9-%d8%b9%d
9%86%d8%af-%d8%a7d9%84%d8%ad%d9%8a%d9%88%d8%a7d9%86-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d9%85%d8%b9%d8%ac%d8%b2%d8%a9-%d8%a7d9%84%d8%b0%d
8%b1%d8%a9-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d9%84%d8%a7-%d8%aa%d8%aa%d8%ac%d8%a7d9%87d9%84-p
df/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d8%a7d9%84%d8%ad%d9%8a%d8%a7d8%a9-%d9%81d9%8a-%
d8%b3%d8%a8d9%8a%d9%84-%d8%a7d9%84d9%84d9%87-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d8%a7d9%84%d8%b9%d8%b8d9%85%d8%a9-%d9%81d9%8a-%
d9%83d9%84-%d9%85d9%83d8%a7d9%86-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d9%85d8%b9%d8%ac%d8%b2%d8%a9-%d8%a7d9%84%d8%ac%d
9%87d8%a7d8%b2-%d8%a7d9%84d9%85d9%86d8%a7d8%b9d9%8a-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d8%a7d9%84d8%b1d9%88d9%85d9%86d8%b3d9%8a%d
8%a9-%d8%b3d9%84d8%a7d8%ad-%d8%a8d9%8a%d8%af-%d8%a7d9%84d8%b4d9%8a%d8%b7d8%a7d9%86-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d8%b1d8%ad%d9%84d8%a9-%d9%81d9%8a-%d8%a7d9%84%
d9%83d9%88d9%86-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d8%a7d9%84d8%aa%d8%b5d9%85d9%8a%d9%85-%d9%81d
9%8a-%d8%a7d9%84d8%b7d8%a8d9%8a%d8%b9d8%a9-pdf/',
        'https://www.arab-books.com/books/%d9%83%d8%aa%d8%a7%d8%a8-%d9%85d8%b9%d8%ac%d8%b2%d8%a9-%d8%ae%d9%84d9%82-%
d8%a7d9%84d8%a5d9%86d8%b3d8%a7d9%86-pdf/']
```

```
In [9]: len(Books_urls)
```

Out[9]: 5984

We have now the url and title of **5984 books**

Request to get more information about each book

```
In [10]: i=0
        scoup_src = []
        for book in Books_urls :
            scoup_src.append(requests.get(book,allow_redirects=False).text)
            i+=1
            if(i%200==0):
                print(len(scoup_src))
```

```
200
400
600
800
1000
1200
1400
1600
1800
2000
2200
2400
2600
2800
3000
3200
3400
3600
3800
```

4000
4200
4400
4600
4800
5000
5200
5400
5600
5800

```
In [11]: scoop_soup = []  
         for page in scoop_src:  
             scoop_soup.append(BeautifulSoup(page, 'lxml'))  
         len(scoop_soup)
```

Out[11]: 5984

```
In [12]: scoop_books_info = []  
         scoop_books_summary=[]  
         for soop_page in scoop_soup :  
             scoop_books_info = scoop_books_info + soop_page.find_all('div',class_="book-info")  
             scoop_books_summary = scoop_books_summary + soop_page.find_all('div',class_="entry-content entry clearfix")
```

```
In [13]: scoop_books_info[0]
```

```
In [14]: scoop_books_summary[0]
```

```
In [15]: books_info=[]  
         for book_info in scoop_books_info:  
             new_book_info=[]  
             for info in book_info.find_all("li"):  
                 new_book_info.append(info.text.split(":"))  
             books_info.append(new_book_info)
```

```
In [16]: books_info[0]
```

```
In [17]: src_download_url=[]  
         for book_url in scoop_soup:  
             src_download_url = src_download_url + book_url.find_all('div',class_="down-link-bottom")
```

```
In [17]: src_download_url[:5]
```

```
In [18]: Books_download_urls=[]  
         for src_url in src_download_url :  
             myResult = URL_RegEx.search(str(src_url))  
             Books_download_urls.append(myResult.group())  
         Books_download_urls[:10]
```

```
In [19]: books_titles[0]
```

```
In [20]: scoop_books_summary[0]
```

```
In [21]: books_summary=[]  
         for text in scoop_books_summary:  
             books_summary.append(text.find_all('p'))  
         books_summary
```

```
In [22]: books_summary[0]
```

```
In [23]: Books_download_urls[0]
```

```
In [26]: books_Authors = []  
         books_Class = []
```

```
books_Language=[]
books_Pages =[]
books_publisher=[]
book_Size =[]
book_format =[]
for Book_Info in books_info :
    books_Authors.append(Book_Info[0][1])
    books_Class.append(Book_Info[1][1])
    books_Language.append(Book_Info[2][1])
    books_Pages.append(Book_Info[3][1])
    books_publisher.append(Book_Info[4][1])
    book_Size.append(Book_Info[5][1])
    book_format.append(Book_Info[6][1])
```

In [24]: `len(book_Size)`

Now we have all information about our 5984 books

Lets save the data in a csv

In [28]: `data_list=[books_titles,books_Authors,books_summary,books_Class,books_Language,books_Pages,books_publisher,book_format,book_Size]
exported=zip_longest(*data_list)`

In [29]: `with open("Scrapped_Data.csv","w",encoding="utf-8") as Scrapped_Data:
 wr=csv.writer(Scrapped_Data)
 wr.writerow(["BookTitle","Author","BookSummary","Class","Language","Pages","publisher","BookSize","format","DownloadURL"])
 wr.writerows(exported)`

2-Clean Scrapped Data

In [25]: `#read Scrapped_Dataset to clean it
data=pd.read_csv("Scrapped_Data.csv")
data.head(2)`

Out[25]:

	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
0	كتاب التوضيحية عند الحيوان PDF	الكاتب هارون يحيى الكاتب هارون ...	الروايات <p>والكتب العربية تعتبر من الروابط...	تحميل الكتب الأدب العربي	عربي	160 صفحة	عدنان أوكتار	5 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	كتاب معجزة الذرة PDF	الكاتب هارون يحيى الكاتب هارون ...	الروايات <p>والكتب العربية تعتبر من الروابط...	تحميل الكتب تطوير الذات	عربي	136 صفحة	عدنان أوكتار	3.41 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...

In [26]: `data.shape`

Out[26]: (5984, 11)

In [27]: `data.tail(2)`

Out[27]:

	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
--	-----------	--------	-------------	-------	----------	-------	-----------	----------	--------	-------------

	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
5982	القران الكريم باللغة الانجليزية - Quran wit...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN books.com/books/%
5983	Le Saint Coran en français	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN books.com/bc

In [28]:

```
#Remove duplicated and null rows
data=data.dropna()
data=data.drop_duplicates()
data.describe()
```

Out[28]:

	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
count	5898	5898	5898	5898	5898	5898	5898	5898	5898	
unique	5792	600	5370	290	21	791	784	1143	11	
top	كتاب العالم نهاية العالم PDF	كتب أحمد خالد توفيق	[<p><span style="font-size: 24px; font-weight:...	تحميل كتب اسلامية	عربي	1 صفحة		1 ميغابايت	pdf	https://www.books.com/books/%d9%83%d8%
freq	3	298	69	1230	3837	133	686	640	4048	

In [29]:

```
data.head(2)
```

Out[29]:

	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
0	كتاب التضحية عند الحيوان PDF	الكاتب هارون يحيى الكاتب هارون ...	[<p>الروايات والكتب العربية تعتبر من الروابط...	تحميل كتب الأدب العربي	عربي	160 صفحة	عدنان أوكتار	5 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	كتاب معجزة الذرة PDF	الكاتب هارون يحيى الكاتب هارون ...	[<p>الروايات والكتب العربية تعتبر من الروابط...	تحميل كتب تطوير الذات	عربي	136 صفحة	عدنان أوكتار	3.41 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...

In [30]:

```
data[["BookTitle","Author","BookSummary"]].head()
```

Out[30]:

	BookTitle	Author	BookSummary
0	كتاب التضحية عند الحيوان PDF	... كتب الكاتب هارون الكاتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط ب<p>
1	كتاب معجزة الذرة PDF	... كتب الكاتب هارون الكاتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط ب<p>
2	كتاب لا تتجاهل PDF	... كتب الكاتب هارون الكاتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط ب<p>
3	كتاب الحياة في سبيل الله PDF	... كتب الكاتب هارون الكاتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط ب<p>
4	كتاب العظمة في كل مكان PDF	... كتب الكاتب هارون الكاتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط ب<p>

```
In [31]: #cleaning our data from any thing except arabic letters
pattern="[^ء-ي-ا- ]"
data["BookTitle"].replace(pattern, ' ', regex=True,inplace=True)
data["Author"].replace(pattern, ' ', regex=True,inplace=True)
data["Class"].replace(pattern, ' ', regex=True,inplace=True)
data["BookSummary"].replace(pattern, ' ', regex=True,inplace=True)
```

```
In [32]: data["BookTitle"][0]
```

```
Out[32]: 'كتاب التضحية عند الحيوان'
```

```
In [33]: #handle the start and end spaces
data["BookTitle"]=data["BookTitle"].apply(lambda title:title.strip())
data["BookTitle"][0]
```

```
Out[33]: 'كتاب التضحية عند الحيوان'
```

```
In [34]: data["Author"][0]
```

```
Out[34]: 'كتب الكاتب هارون يحيى كتب الكاتب هارون يحيى'
```

```
In [35]: data["Class"][0]
```

```
Out[35]: 'تحميل كتب الأدب العربي تحميل كتب متنوعة'
```

```
In [36]: data["BookSummary"][0]
```

```
Out[36]: 'الروايات والكتب العربية تعتبر من الروابط بيننا وبين تاريخنا كعرب او بيننا وبين اللغة العربية وبين تاريخنا واحداث هذا التاريخ الروايات العربية الي الان ومع تقدم العصور والتكنولوجيا لازالت قائمة وحاضرة كتاب التضحية عند الحيوان للكاتب هارون يحيى تعتبر صفحات التضحية والتعاون والرحمة والمحبة والاحترام في البيئة التي يعيش فيها ولعل كل انسان يسعى إلى أن يحيا وسط مجتمع تنتشر فيه هذه الفضائل وفي هذا الكتاب وقع تناول هذه الفضائل بالبحث هنا وتجمع فيه جميع هذه الفضائل هو الحيوان فثمة حيوانات تلقى نفسها إلى الهلاك والخطر في مواجهة أعدائها من أجل حماية حياة صغارها وثمة طيور تعرض نفسها للخطر أيضا من أجل حماية فراخها وفي مناطق القطب المتجمد تمكث طيور البطريق لأشهر دون أكل ودون حركة حتى توفر الحرارة والدفع لفراخها ونجد كذلك حيوانات تتحمل مسؤولية رعاية من انقطعت به السبل من صغار الحيوانات وبقي يتيمًا بلا أم ونجد من أمثال هذه النماذج الكثير الكثير كتاب التضحية عند لحيوان عدنان أوكتار المعروف باسم هارون يحيى هو كاتب وباحث تركي مسلم ولد في أنقرة عام ١٩٨٠ وعاش فيها حتى عام ٢٠١٠ م عندما انتقل إلى إسطنبول حيث التحق بكلية الفنون الجميلة في جامعة المعمار سنان وخلال سنواته الجامعية قام ببحوث مفصلة في الفلسفة المادية والايديولوجية السائدة التي تحيط به قام بإنشاء مؤسسة البحث العلمي في تركيا ركز كتاباته على تفنيد وتكذيب نظريات التطور والارتقاء والنشوء وبيان تناقضها حسب رأيه كما يركز في كتاباته أيضا على م موضوعات الماسونية والصهيونية والإلحاد له أكثر من مئة كتاب حول قيم وأخلاقيات القرآن وحول مواضيع إيمانية عديدة ومختلفة وترجمت إلى العديد من اللغات ال عالمية واسمه القلمي هارون يحيى ومن أشهر كتبه كتاب أطلس الخلق الذي يقع في صفحة ويتحدث عن رفض نظرية النشوء والارتقاء لداروين ومن أشهر م ولفاته أطلس الخلق نهاية الدارونية معجزة الذرة الإعجاز في خلق النباتات خلق الكون المفاهيم الأساسية في القرآن هل فكرت في الحقيقة الكوارث التي جلبتها الدارونية للإنسانية معجزة النمل معجزة النحل معجزة الجهاز المناعي استعمل عدنان أوكتار الاسم المستعار جاويد بالجن في بعض كتبه ولكن القسم الأكبر منها نشره بالاسم المستعار هارون يحيى وهذا الاسم المستعار يتكون من اسمي نبيين في إشارة إلى ذكرى النبي هارون والنبي يحيى اللذين حاربوا أفكار الإلحاد كتاب التضحية عند الحيوان ترون بأنفسكم نماذج من السلوك عند بعض الحيوانات يؤثر الاستغراب والدهشة ويبعث على التأمل ذلك أنه سلوك ينطق بالحكمة والعقل ويبين عن وعي خارق وبعد نظر ع ميق وعند قراءة هذه النماذج سوف يقفز إلى أذهانكم السؤال التالي كيف يمكن لحيوان أن يعرف كل هذا كيف يمكنه أن يفكر في كل هذا وسوف تجدون الإجابة على هذه التساؤلات في الكتاب الذي بين أيديكم إن الذي خلقها جميعها هو الله الرحمن الرحيم وهو الذي يلهمها طريقها فتسير فيه طائفة فكل ما في السما وات وما في الأرض طوع أمره سبحانه وتعالى يمكنك أيضا قراءة وتحميل روايات عربي من خلال مكتبتكم ' المكتبة العربية مثل
```

```
In [37]: #remove all this unimportant spaces
space_patern="[ ]+"
data["BookSummary"].replace(space_patern, ' ', regex=True,inplace=True)
```

```
In [38]: data["BookSummary"][0]
```

```
Out[38]: 'الروايات والكتب العربية تعتبر من الروابط بيننا وبين تاريخنا كعرب او بيننا وبين اللغة العربية وبين تاريخنا واحداث هذا التاريخ الروايات العربية الي الان ومع تقدم الع صور والتكنولوجيا لازالت قائمة وحاضرة كتاب التضحية عند الحيوان للكاتب هارون يحيى تعتبر صفحات التضحية والتعاون والرحمة والمحبة والتكافل من الأخلاق الحسنة الت ي يسعى كل مجتمع لنشرها بين أفرادها وهذه الميزات تسبغ على الإنسان الذي يتحلى بها المحبة والاحترام في البيئة التي يعيش فيها ولعل كل انسان يسعى إلى أن يحيا وس ط مجتمع تنتشر فيه هذه الفضائل وفي هذا الكتاب وقع تناول هذه الفضائل والفضائل بالدرس بيد أن الذي يتصف بالتضحية والإيثار هنا والذي يقوم بأعمال التعاون والتكافل ويفض شفقة ومحبة ورحمة ليس كائنًا بشريًا إن الذي نتناوله بالبحث هنا وتجمع فيه جميع هذه الفضائل هو الحيوان فثمة حيوانات تلقى نفسها إلى الهلاك والخطر في مواجه ة أعدائها من أجل حماية حياة صغارها وثمة طيور تعرض نفسها للخطر أيضا من أجل حماية فراخها وفي مناطق القطب المتجمد تمكث طيور البطريق لأشهر دون أكل ودو ن حركة حتى توفر الحرارة والدفع لفراخها ونجد كذلك حيوانات تتحمل مسؤولية رعاية من انقطعت به السبل من صغار الحيوانات وبقي يتيمًا بلا أم ونجد من أمثال هذه ال نماذج الكثير الكثير كتاب التضحية عند الحيوان عدنان أوكتار المعروف باسم هارون يحيى هو كاتب وباحث تركي مسلم ولد في أنقرة عام ١٩٨٠ وعاش فيها حتى عام ٢٠١٠ م عندما انتقل إلى إسطنبول حيث التحق بكلية الفنون الجميلة في جامعة المعمار سنان وخلال سنواته الجامعية قام ببحوث مفصلة في الفلسفة المادية والايديولوجية السائدة التي تحيط به قام
```

بإنشاء مؤسسة البحث العلمي في تركيا تركز كتاباته على تنفيذ وتكذيب نظريات التطور والارتقاء والنشوء وبيان تناقضها حسب رأيه كما يركز في كتاباته أيضا على موضوعات الماسونية والصهيونية والإلحاد له أكثر من مئة كتاب حول قيم وأخلاقيات القرآن وحول مواضيع إيمانية عديدة ومختلفة وترجمت إلى العديد من اللغات العالمية واسمه القلم ي هارون يحيى ومن أشهر كتبه كتاب أطلس الخلق الذي يقع في صفحة ويتحدث عن رفض نظرية النشوء والارتقاء لداروين ومن أشهر مؤلفاته أطلس الخلق نهاية الدارونية معجزة الذرة الإعجاز في خلق النباتات خلق الكون المفاهيم الأساسية في القرآن هل فكرت في الحقيقة الكوارث التي جلبتها الداروينية للإنسانية معجزة النمل معجزة النحل مع جزة الجهاز المناعي استعمل عدنان أوكطار الاسم المستعار جاويد بالجَن في بعض كتبه ولكن القسم الأكبر منها نشره بالاسم المستعار هارون يحيى وهذا الاسم المستعار يتكرر من اسمي نبيين في إشارة إلى ذكرى النبي هارون والنبي يحيى اللذين حاربوا أفكار الإلحاد كتاب التضحية عند الحيوان ترون بأنفسكم نماذج من السلوك عند بعض الحيوانات يؤثر الاستغراب والدهشة وبيع على التأمل ذلك أنه سلوك ينطق بالحكمة والعقل ويبين عن وعي خارق وبعد نظر عميق وعند قراءة هذه النماذج سوف يقفز إلى أذهانكم السؤال التالي كيف يمكن لحيوان أن يعرف كل هذا كيف يمكنه أن يفكر في كل هذا وسوف تجدون الإجابة على هذه التساؤلات في الكتاب الذي بين أيديكم إن الذي خلفها جميعها هو الله الرحمن الرحيم وهو الذي يلهمها طريقها فتسير فيه طائعة فكل ما في السماوات وما في الأرض طوع أمره سبحانه وتعالى يمكنك أيضا قراءة وتحميل ر

' وايات عربي من خلال مكتبكم المكتبة العربية مثل

```
In [39]: #remove redundant words unusefal
data["BookTitle"].replace(r"كتاب", "", regex=True, inplace=True)
data["Author"].replace(r"الكاتب", "", regex=True, inplace=True)
data["Class"].replace(r"كتب", "", regex=True, inplace=True)
```

```
In [40]: data[["BookTitle", "Author", "BookSummary", "Class"]].head()
```

	BookTitle	Author	BookSummary	Class
0	التضحية عند الحيوان	كتب هارون يحيى كتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	تحميل الأدب العربي تحميل متنوعة
1	معجزة الذرة	كتب هارون يحيى كتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	تحميل تطوير الذات تحميل متنوعة
2	لا تتجاهل	كتب هارون يحيى كتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	تحميل الأدب العربي
3	الحياة في سبيل الله	كتب هارون يحيى كتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	تحميل اسلامية
4	العظمة في كل مكان	كتب هارون يحيى كتب هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	تحميل الأدب العربي

```
In [41]: #cleaning from redundant words unusefal
data["Class"]=data["Class"].apply(lambda text:text.replace("1.", "", "تحميل"))
data["Author"]=data["Author"].apply(lambda text:text.replace("1.", "", "كتب"))
#split Classes into multiclass & split Authers too
data["Class"]=data["Class"].apply(lambda row:row.split("تحميل"))
data["Author"]=data["Author"].apply(lambda row:row.split("كتب"))
```

```
In [42]: data[["BookTitle", "Author", "BookSummary", "Class"]].head()
```

	BookTitle	Author	BookSummary	Class
0	التضحية عند الحيوان	[هارون يحيى , هارون يحيى]	...الروايات والكتب العربية تعتبر من الروابط بين	[الأدب العربي , متنوعة]
1	معجزة الذرة	[هارون يحيى , هارون يحيى]	...الروايات والكتب العربية تعتبر من الروابط بين	[تطوير الذات , متنوعة]
2	لا تتجاهل	[هارون يحيى , هارون يحيى]	...الروايات والكتب العربية تعتبر من الروابط بين	[الأدب العربي]
3	الحياة في سبيل الله	[هارون يحيى , هارون يحيى]	...الروايات والكتب العربية تعتبر من الروابط بين	[اسلامية]
4	العظمة في كل مكان	[هارون يحيى , هارون يحيى]	...الروايات والكتب العربية تعتبر من الروابط بين	[الأدب العربي]

```
In [43]: data["Author"][0]
```

```
Out[43]: [ ' هارون يحيى ' , ' هارون يحيى ' ]
```

```
In [44]: #handle the start and end spaces
data["Author"]=data["Author"].apply(lambda Authorlist:[Author.strip() for Author in Authorlist])
data["Author"][0]
```

```
Out[44]: [ 'هارون يحيى' , 'هارون يحيى' ]
```

```
In [45]: data["Class"][0]
```

```
Out[45]: [ ' الأدب العربي ' , ' متنوعة ' ]
```

```
In [46]: #handle the start and end spaces
```

```
data["Class"] = data["Class"].apply(lambda classlist: [myclass.strip() for myclass in classlist])
data["Class"][0]
```

['الأدب العربي' , 'متنوعة']

```
#delete all duplicated values
data["Author"] = data["Author"].apply(lambda Authorlist: list(set(Authorlist)))
data["Class"] = data["Class"].apply(lambda classlist: (list(set(classlist))))
```

```
data.head(3)
```

BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL	
0	التضحية عند الحيوان	هارون[يحيى	الروايات والكتب العربية تعتبر من...الروابط بينن	[الأدب] العربي، متنوعة	عربي	160 صفحة	عدنان أوكتار	5 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	معجزة الذرة	هارون[يحيى	الروايات والكتب العربية تعتبر من...الروابط بينن	[متنوعة] تطوير الذات	عربي	136 صفحة	عدنان أوكتار	3.41 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...
2	لا تتجاهل	هارون[يحيى	الروايات والكتب العربية تعتبر من...الروابط بينن	[الأدب] العربي	عربي	100 صفحة	عدنان أوكتار	3 ميغابايت	pdf	https://www.arab-books.com/books/%d9%83%d8%aa%...

```
data["format"].unique()
```

```
array([' pdf', ' ', ' excel', ' word', ' PDF', ' 15.11', ' Pdf', ' حبيب',  
      ' P', ' PDF ZIP', ' PDF'], dtype=object)
```

```
data[data["format"]== ' 15.11']
```

BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	Download	
3936	أحكام الأضحية والزكاة	الشيخ محمد بن صالح [العثيمين]	كتاب زاد الداعية إلى الله للشيخ ... العلامة محمد	[العقيدة]	العربية	34 صفحة	دار الثريا	15.11 ميغابايت	15.11	https://www.arabbooks.com/books/%d9%83%d8%aa

```
#cleaning
PDF=[" pdf"," PDF"," Pdf"," P"," PDF ZIP"," PDf"," 15.11 "," ","حبيب"]
data["format"].replace(PDF,"PDF",inplace=True)
```

```
data["format"] = unique()
```

```
data[format].unique()
```

```
array(['PDF', ' excel', ' word'], dtype=object)
```

```
data["experiments"] = unique()
```

```
data["Language"].unique()

array(['عربي', 'إنجليزي', 'العربية', 'العربية', 'العربية', 'الإنجليزية',
      'الألمانية', 'عربية - المانية', 'العربية - الألمانية', 'العربية',
      'اللغة العربية', 'العربية', 'العربية', 'اللغة العربية', 'العربية'],
      dtype=object)
```

```
Arabic=[' العربية ', ' اللغة العربية.', ' العربية', ' العربية', ' اللغة العربيه', ' اللغة العربية', ' اللة العربية', ' الألمانية - العربية - الألمانية',  
English=[' الإنجليزية ', ' إنجليزي', ' English', ' الانجليزية']  
German=' الألمانية'
```

```
#cleaning
data["Language"].replace(Arabic,"اللغة العربية",inplace=True)
data["Language"].replace(English,"اللغة الانجليزية",inplace=True)
```



```
data["Language"].replace(German,"اللغة الألمانية",inplace=True)
data["Language"].unique()
```

```
Out[55]: array(['اللغة العربية', 'اللغة الانجليزية', 'اللغة الألمانية'],
      dtype=object)
```

```
In [56]: data[["BookTitle","Author","BookSummary","Class"]].head()
```

```
Out[56]:
```

	BookTitle	Author	BookSummary	Class
0	التضحية عند الحيوان	هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	[الأدب العربي, متنوعة]
1	معجزة الذرة	هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	[متنوعة, تطوير الذات]
2	لا تتجاهل	هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	[الأدب العربي]
3	الحياة في سبيل الله	هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	[اسلامية]
4	العظمة في كل مكان	هارون يحيى	...الروايات والكتب العربية تعتبر من الروابط بين	[الأدب العربي]

```
In [63]: Authors=[]
for Author in data["Author"]:
    if len(Author)==1:
        Authors.append(Author)
Authors[10:15]
```

```
Out[63]: [['محمود سالم'],
          ['محمود سالم'],
          ['محمود سالم'],
          ['هارون يحيى'],
          ['هارون يحيى']]
```

```
In [64]: Authors=[]
for Author in data["Author"]:
    if len(Author)>1:
        Authors.append(Author)
Authors[10:15]
```

```
Out[64]: [['', 'الدكتور سليم حسن'],
          ['', 'الدكتور سليم حسن'],
          ['', 'الدكتور سليم حسن'],
          ['', 'الدكتور سليم حسن'],
          ['', 'الدكتور سليم حسن']]
```

```
In [65]: #cleaning the column from fake values
data["Author"]= data["Author"].apply(lambda x:[item for item in x if item not in ['']])
```

```
In [66]: Authors=[]
for Author in data["Author"]:
    if len(Author)>1:
        Authors.append(Author)
Authors[10:15]
```

```
Out[66]: [['كوثر محمود عبد الرسول', 'محمد رياض'],
          ['ابن جزي الكلبي', 'ابن بطوطة'],
          ['جايا جرانت', 'أندرو جرانت'],
          ['تشنسر إلتون', 'أدريان جوستيك'],
          ['بيتر سي كايرو', 'ديفيد إل دوليتش', 'ستيفن إتش راينسميث']]
```

```
In [69]: classs=[]
for bookclass in data["Class"]:
    if len(bookclass)==1:
        classs.append(bookclass)
classs[10:15]
```

```
Out[69]: [['الأدب العربي'], ['اسلامية'], ['الفلسفة والمنطق'], ['اسلامية'], ['اسلامية']]
```

```
In [74]: classs=[]
for bookclass in data["Class"]:
    if len(bookclass)>1:
        classs.append(bookclass)
classs[150:155]
```

```
Out[74]: [['الأدب العربي', 'روايات عربية', 'الروايات العالمية المترجمة'],
['الأدب العربي', 'روايات عربية', 'الروايات العالمية المترجمة'],
['الأدب العربي', 'روايات عربية', 'الروايات العالمية المترجمة'],
['الأدب العربي', 'روايات عربية', 'الروايات العالمية المترجمة'],
['الأدب العربي', 'روايات عربية', 'الروايات العالمية المترجمة']]
```

```
In [75]: #cleaning the column from fake values
data["Class"] = data["Class"].apply(lambda x:[item for item in x if item not in ['']])
```

```
In [76]: data.head(3)
```

```
Out[76]:
```

	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
0	التضحية عند الحيوان	هارون [يحيى]	الروايات والكتب العربية تعتبر من...الروابط بين	الأدب [العربي، متنوعة]	اللغة العربية	160 صفحة	عدنان أوكتار	5 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	معجزة الذرة	هارون [يحيى]	الروايات والكتب العربية تعتبر من...الروابط بين	متنوعة [تطوير الذات]	اللغة العربية	136 صفحة	عدنان أوكتار	3.41 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
2	لا تتجاهل	هارون [يحيى]	الروايات والكتب العربية تعتبر من...الروابط بين	الأدب [العربي]	اللغة العربية	100 صفحة	عدنان أوكتار	3 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...



now our dataset is clean

lets save it in CSV

```
In [77]: data.to_csv("CleanScrappedDataset.csv")
data.to_json("CleanScrappedDataset.json")
```

3-auto correct the scrapped data

3.1-Auto contextual_correct

```
In [78]: from ar_corrector.corrector import Corrector
corr = Corrector()
```

```
In [79]: dataset=pd.read_json("CleanScrappedDataset.json")
```

```
In [80]: #Lenth befor context correction
len(dataset["BookSummary"][0])
```

```
Out[80]: 2831
```

```
In [81]: #Len after context correction
len(dataset["BookSummary"].iloc[0:1].apply(lambda row:corr.contextual_correct(row))[0])
```

```
Out[81]: 2824
```

we notice the diferrent between two lengths

lets apply to all data

```
In [82]: #context correction this take too time
dataset["BookSummary"] = dataset["BookSummary"].apply(lambda row:corr.contextual_correct(row))
```

Tokenization to make Auto correction on words

```
In [84]: dataset['BookSummary'] = dataset.apply(lambda row: nltk.word_tokenize(row['BookSummary']), axis=1)
```

```
In [86]: dataset["BookSummary"].iloc[0][:10]
```

Out[86]:

'الروايات'
'والكتب'
'العربية'
'تعتبر'
'من'
'الروابط'
'بيننا'
'وبين'
'تاريخنا'
'كعرب'

3.2-Auto word correct the scrapped data

```
In [88]: dataset.head(3)
```

```
Out[88]:
```

BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
-----------	--------	-------------	-------	----------	-------	-----------	----------	--------	-------------

0	التصحية عند الحيوان	هارون[يحيى	الروايات، والكتب، العربية، تعتبر، من، ...الروابط	الأدب، العربي، [متنوعة]	اللغة العربية	160 صفحة	عدنان أوكتار	5 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	معجزة الذرة	هارون[يحيى	الروايات، والكتب، العربية، تعتبر، من، ...الروابط	متنوعة، تطوير [الذات]	اللغة العربية	136 صفحة	عدنان أوكتار	3.41 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
2	لا تتجاهل	هارون[يحيى	الروايات، والكتب، العربية، تعتبر، من، ...الروابط	الأدب، العربي، [الأدب]	اللغة العربية	100 صفحة	عدنان أوكتار	3 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...

```
In [89]: check=lambda word : type(corr.spell_correct(word,1)) is type([])
         get =lambda word : corr.spell_correct(word,1)[0][0]
```

```
In [90]: #this take alot of time to correct all words
dataset["BookSummary"] = dataset["BookSummary"].apply(lambda row: [get(word) if check(word) else word for word in row["BookSummary"]])
```

Remove Aabic Stop words

```
In [92]: import arabicstopwords.arabicstopwords as stp
stp.classed_stopwords_list()
```

Out[92]: dict_keys(['حم' و 'غداة' و 'جنوب' و 'ذواتا' و 'حي' و 'لازلنا' و 'زمان' و 'عوض' و 'بنا' و 'أجمع' و 'أو' و 'إيه' و 'لن' و 'ها' و 'ش' و 'غرب' و 'أولئك
لازلتم' و 'هج' و 'ههب' و 'عسى' و 'ل' و 'أنفا' و 'لازلتن' و 'هن' و 'هم' و 'هل' و 'هاؤم' و 'ارتد' و 'هي' و 'هو' و 'لازالنا' و 'تلكما' و 'وقت' و 'أولئك
' و 'نحو' و 'حسب' و 'نحن' و 'لن' و 'الذين' و 'أب' و 'أخ' و 'قبل' و 'بدون' و 'مادامتا' و 'بنس' و 'د' و 'ذا' و 'العمر' و 'هاهنا' و 'كما' و 'الستم' و 'ل
ستن' و 'حتى' و 'لدى' و 'ذه' و 'ذي' و 'ذو' و 'أي' و 'أو' و 'أف' و 'أن' و 'أم' و 'نعما' و 'هيت' و 'هيا' و 'مايرج' و 'حينما' و 'هلا' و 'إنما' و 'ج
ل' و 'كخ' و 'ت' و 'عما' و 'يكن' و 'بكم' و 'مازلنا' و 'ع' و 'رويدك' و 'دون' و 'أولئك' و 'كي' و 'هؤلاء' و 'لها' و 'ي' و 'مكانكما' و 'كم' و 'كل' و 'م
اقتتن' و 'ماانفكتكما' و 'لعل' و 'إن' و 'لازال' و 'متي' و 'مايرحوا' و 'خلال' و 'مازلتن' و 'مازلتم' و 'تتا' و 'راح' و 'هلم' و 'لولا' و 'مايرحمتا' و 'ماقتنت
ماقتنتم' و 'تان' و 'مرة' و 'أصلا' و 'وشكان' و 'كانما' و 'إياهما' و 'أسفل' و 'ص' و 'أمامك' و 'م' و 'مادمتما' و 'إياهم' و 'مافانكوا' و 'استحال' و 'إليك' و
'أمسى' و 'انبرى' و 'ماقتنتا' و 'انتم' و 'لازالوا' و 'انتن' و 'ليستا' و 'ماانفكتن' و 'أوشك' و 'طفق' و 'الزالت' و 'صه' و 'لازالا' و 'هناك' و 'ليس' و 'هاذي
ن' و 'مازالوا' و 'الإلا' و 'سوف' و 'ب' و 'أمين' و 'آء' و 'هاتين' و 'ليت' و 'و' و 'أكثر' و 'إلى' و 'أثناء' و 'ابن' و 'بعحث' و 'فلان' و 'ليل' و 'تين' و 'أ
صبح' و 'لكما' و 'كأين' و 'قد' و 'ماانفكت' و 'قط' و 'لحظة' و 'شمال' و 'ثم' و 'أجل' و 'عليك' و 'دونك' و 'ذلكما' و 'حما' و 'أنت' و 'التي' و 'أوه' و 'أن
ا' و 'مثل' و 'هذين' و 'س' و 'لازلتما' و 'ك' و 'أنى' و 'حمو' و 'حمي' و 'هما' و 'جبر' و 'أبدا' و 'للذان' و 'عند' و 'ظل' و 'سرعان' و 'خلف' و 'تحو
ل' و 'هاته' و 'هاتي' و 'حوالى' و 'جل' و 'تينك' و 'أولاء' و 'آنذاك' و 'بعد' و 'حين' و 'ذبت' و 'أيان' و 'قلما' و 'مازالا' و 'خ' و 'إزاء' و 'بعض' و 'فيم
ا' و 'لنا' و 'مكانن' و 'تلكم' و 'ذين' و 'اللائي' و 'تلقاء' و 'خلا' و 'حيث' و 'ماقتنوا' و 'صار' و 'أيما' و 'يمين' و 'رب' و 'ماذا' و 'لست' و 'ماقتن' و 'ذا
ن' و 'إياه' و 'ذاك' و 'ة' و 'ن' و 'حذار' و 'ماانفكتا' و 'ع' و 'مادمت' و 'مادمتن' و 'مايرحنا' و 'ى' و 'ذات' و 'واها' و 'أما' و 'جدا' و 'أمد' و 'اللواتي' و
'أمس' و 'لسن' و 'معاذ' و 'الذين' و 'حيثما' و 'تحت' و 'غير' و 'حاشا' و 'ماقتى' و 'مايرحتن' و 'مايرحتم' و 'أمام' و 'ج' و 'أضحى' و 'التي' و 'هاتان' و
هكذا' و 'مايرحنا' و 'كلتا' و 'بضع' و 'اللقين' و 'سواء' و 'السمتا' و 'إليك' و 'إليك' و 'غدا' و 'حاي' و 'كلما' و 'تبدل' و 'مكانكم' و 'كذلك' و 'شتا
و' و 'الآن' و 'أولالك' و 'إياي' و 'ممن' و 'إياك' و 'كيت' و 'بعدهما' و 'أخلوقي' و 'حرى' و 'مادمتا' و 'ط' و 'ه' و 'عاد' و 'كيف' و 'إياكم' و 'إياكن' و 'سب

جان' و 'هذان' و 'أقبل' و 'الألاء' و 'مما' و 'ليسوا' و 'لما' و 'ضحة' و 'مابرحا' و 'أها' و 'ماداموا' و 'آ' و 'كاد' و 'عل' و 'عن' و 'مافتنتما' و 'ز' و 'لازل
ت' و 'بيد' و 'ق' و 'ظ' و 'ماداما' و 'كان' و 'كلاهما' و 'بين' و 'هيهات' و 'إذا' و 'إياكما' و 'الذي' و 'بل' و 'نخ' و 'طاق' و 'أبو' و 'أبي' و 'هذا' و 'لاس
يما' و 'هنا' و 'ساعما' و 'ذنيك' و 'أه' و 'ساعة' و 'ح' و 'ابتدأ' و 'صباح' و 'ضمن' و 'لكيلا' و 'بمن' و 'كيفما' و 'سوى' و 'إياها' و 'تاتك' و 'أقل' و 'إليك
ا' و 'لنلا' و 'كذا' و 'أض' و 'بما' و 'مانفكا' و 'قيم' و 'إياهن' و 'فا' و 'هذي' و 'أبا' و 'هذه' و 'تجاه' و 'شطر' و 'حينننذ' و 'أيضا' و 'شرق' و 'رجع' و
'عامه' و 'إذن' و 'غروب' و 'عندما' و 'وراءك' و 'إذ' و 'ألا' و 'لكي' و 'مابرحت' و 'لكن' و 'لكم' و 'يومنذ' و 'لكنما' و 'قيلما' و 'بي' و 'ذ' و 'مانفكنا' و
'بك' و 'بن' و 'به' و 'اللذان' و 'مانفككنم' و 'تاره' و 'بس' و 'مابرحن' و 'فقط' و 'أنشأ' و 'مانفككن' و 'بخ' و 'متلما' و 'شبه' و 'بها' و 'مازلتما' و 'فوق' و
'شهر' و 'مانفكت' و 'إي' و 'اللاتي' و 'مازال' و 'إن' و 'كان' و 'كليهما' و 'مافتنا' و 'كليكما' و 'شرع' و 'ليسا' و 'مكانك' و 'ليست' و 'وي' و 'أعلى' و 'لس
نا' و 'ثمة' و 'لات' و 'حسيما' و 'ث' و 'طالما' و 'نوا' و 'مع' و 'مانفك' و 'نفس' و 'مذ' و 'يكما' و 'ذوي' و 'نوو' و 'مه' و 'من' و 'حبذا' و 'مادام' و
'عدا' و 'بينما' و 'وراء' و 'بات' و 'عدس' و 'وا' و 'كلا' و 'على' و 'علق' و 'مازلن' و 'حار' و 'وراءكن' و 'مازلتا' و 'وراءكم' و 'مادمت' و 'أخي' و 'أخ
و' و 'قام' و 'ئ' و 'إنما' و 'أين' و 'بطان' و 'الآلي' و 'ض' و 'ذلكم' و 'ذلكن' و 'أية' و 'نعم' و 'أيا' و 'هاك' و 'ذواتي' و 'مازلت' و 'أننذ' و 'أخذ' و 'ذل
ك' و 'عندنذ' و 'أخا' و 'ويكان' و 'إيانا' و 'كرب' و 'لي' و 'لو' و 'لن' و 'له' و 'لم' و 'لك' و 'بلا' و 'مادمن' و 'ذاتك' و 'إما' و 'هنالك' و 'مادامت' و
'ء' و 'ريث' و 'طق' و 'جميع' و 'حول' و 'ر' و 'مهما' و 'تي' و 'كأي' و 'لوما' و 'ته' و 'ف' و 'فو' و 'في' و 'أ' و 'أول' و 'مافتننا' و 'مساء' و 'بماذ
[['ا' و 'عين' و 'لا' و 'جنب' و 'تلك' و 'بله' و 'منذ' و 'أنتما' و 'وراءكما' و 'لازلن' و 'ما' و 'بلي



In [93]:

```
#remove stop words
dataset['BookSummary'] = dataset['BookSummary'].apply(lambda row: [word for word in row if not stop.is_stop(word)]
```

In [94]:

```
dataset[['BookTitle', 'Author', 'Class', "BookSummary"]].head()
```

Out[94]:	BookTitle	Author	Class	BookSummary
0	التضحية عند الحيوان	هارون [يحيى]	[الأدب العربي, متنوعة]	...الروايات, والكتب, العربية, تعتبر, الروابط, تا]
1	معجزة الذرة	هارون [يحيى]	[متنوعة, تطوير الذات]	...الروايات, والكتب, العربية, تعتبر, الروابط, تا]
2	لا تتجاهل	هارون [يحيى]	[الأدب العربي]	...الروايات, والكتب, العربية, تعتبر, الروابط, تا]
3	الحياة في سبيل الله	هارون [يحيى]	[اسلامية]	...الروايات, والكتب, العربية, تعتبر, الروابط, تا]
4	العظمة في كل مكان	هارون [يحيى]	[الأدب العربي]	...الروايات, والكتب, العربية, تعتبر, الروابط, تا]

Stemming

In [95]:

```
from nltk.stem.isri import ISRIStemmer
st = ISRIStemmer()
```

In [96]:

```
dataset['BookSummary']=dataset['BookSummary'].apply(lambda x: [st.stem(item) for item in x])
```

In [97]:

```
dataset['BookSummary'].iloc[0][:10]
```

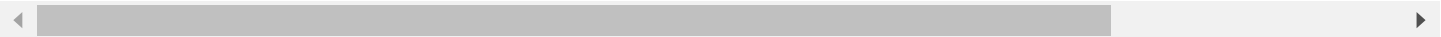
Out[97]:

```
['روي' و 'كتب' و 'عرب' و 'عبر' و 'ربط' و 'ارخ' و 'كعرب' و 'او' و 'لغة' و 'عرب']
```

In [99]:

```
dataset.head(3)
```

Out[99]:	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
0	التضحية عند الحيوان	هارون [يحيى]	روي, كتب, عرب, عبر, ربط, ارخ, ... ,كعرب, او, لغة	الأدب العربي, [متنوعة]	اللغة العربية	160 صفحة	عدنان أوكتار	5 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	معجزة الذرة	هارون [يحيى]	روي, كتب, عرب, عبر, ربط, ارخ, ... ,كعرب, او, لغة	متنوعة, [تطوير الذات]	اللغة العربية	136 صفحة	عدنان أوكتار	3.41 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
2	لا تتجاهل	هارون [يحيى]	روي, كتب, عرب, عبر, ربط, ارخ, ... ,كعرب, او, لغة	الأدب [العربي]	اللغة العربية	100 صفحة	عدنان أوكتار	3 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...



4-Save your data in csv file "Data sheet"

with the following columns " Author - title - publishing year - link for pdf - etc. "

```
In [100... dataset.to_csv("finalCleanedData.csv",header=True)
```

Save your data in json file to save our lists as it is

```
In [101... dataset.to_json("finalCleanedData.csv")
```

5-Cluster this data for topic modeling

```
In [102... mydata=pd.read_json("finalCleanedData.csv")
```

```
In [103... dataset.head()
```

Out[103...	BookTitle	Author	BookSummary	Class	Language	Pages	publisher	BookSize	format	DownloadURL
0	التضحية عند الحيوان	هارون [يحيى	روي, كتب, عرب, [عبر, ربط, ارخ, ... ,كعرب, او, لغة	الأدب العربي, [متنوعة	اللغة العربية	160 صفحة	عدنان أوكتار	5 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
1	معجزة الذرة	هارون [يحيى	روي, كتب, عرب, [عبر, ربط, ارخ, ... ,كعرب, او, لغة	متنوعة, [تطوير الذات	اللغة العربية	136 صفحة	عدنان أوكتار	3.41 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
2	لا تتجاهل	هارون [يحيى	روي, كتب, عرب, [عبر, ربط, ارخ, ... ,كعرب, او, لغة	الأدب العربي	اللغة العربية	100 صفحة	عدنان أوكتار	3 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
3	الحياة في سبيل الله	هارون [يحيى	روي, كتب, عرب, [عبر, ربط, ارخ, ... ,كعرب, او, لغة	[اسلامية]	اللغة العربية	100 صفحة	عدنان أوكتار	2.11 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...
4	العظمة في كل مكان	هارون [يحيى	روي, كتب, عرب, [عبر, ربط, ارخ, ... ,كعرب, او, لغة	الأدب العربي	اللغة العربية	140 صفحة	عدنان أوكتار	5.3 ميغابايت	PDF	https://www.arab-books.com/books/%d9%83%d8%aa%...

```
In [104... #clustering by BookSummary base
mydata=mydata[["BookTitle", "Author", "BookSummary"]]
mydata.head()
```

Out[104...	BookTitle	Author	BookSummary
0	التضحية عند الحيوان	هارون [يحيى	... ,روي, كتب, عرب, عبر, ربط, ارخ, كعرب, او, لغة
1	معجزة الذرة	هارون [يحيى	... ,روي, كتب, عرب, عبر, ربط, ارخ, كعرب, او, لغة
2	لا تتجاهل	هارون [يحيى	... ,روي, كتب, عرب, عبر, ربط, ارخ, كعرب, او, لغة
3	الحياة في سبيل الله	هارون [يحيى	... ,روي, كتب, عرب, عبر, ربط, ارخ, كعرب, او, لغة
4	العظمة في كل مكان	هارون [يحيى	... ,روي, كتب, عرب, عبر, ربط, ارخ, كعرب, او, لغة

```
In [105... mydata["BookSummary"][0][:10]
```

```
Out[105... ['روي', 'كتب', 'عرب', 'عبر', 'ربط', 'ارخ', 'كعرب', 'او', 'لغة', 'عرب']
```

```
In [106...] mydata["BookSummary"]=mydata["BookSummary"].apply(lambda row: ' '.join(row))
```

```
In [107...] mydata["BookSummary"][0]
```

```
Out[107...] 'روي كتب عرب عبر ربط ارج كعرب او لغة عرب ارج حدث ارج روي عرب الي الن قدم عصر كتولوج قثم حضر كتب ضحي حيو كتب هار بجى عبر صفح ضح'
ي تعا رحم حبة كفل خلق حسن سعى جمع نشر فرد ميز سينغ انس حلى حبة حرم بيئ يعيش انس سعى يحا وسط جمع نشر فضل كتب وقع نول خصل فضل درس تص
ف ضحي يثر يقم أعمال تعا كفل يفض شفق محب رحم كائ بشر نول بحث جمع فضل حيو حيو لقي هلك خطر وجه عدئ حمي حبة صغر طير عرض خطر حمي فرخ
نطق قطب تجمد مكث طير طروق أشهر حرك وفر حرر دفء فرخ نجد حيو حمل سؤل رعي قطع سبل صغر حيو وبق يتي نجد مثل اذج كثر كثر كتب ضحي حيو عدن
كطر عرف بسم هار يحى كتب بحث ترك سلم ولد نقر عام وعش عام نقل لإسطنبول تحق بكل فنن جمل جمع عمر سنن سنا جمع بحث فصل لسف ادة ايديولوجية سند ت
حط إنشاء وُسس بحث علم ترك ركز كتب فند كذب نظر تطر رقة نشء وبى نقض راه ركز كتب وضع اسن صهو لحد منة كتب قيم أخلاق قرأ مواضيع يمن عدد خلف
رجم عدد لغت علم وسم قلم هار يحى شهر كتب كتب طلס خلق يقع صفح حدث رفض نظر نشء رقة درو شهر ولف طلס خلق نهى درن عجز لذر عجز خلق نيت
خلق كون مفاهيم سسي قرأ فكر حقق كرت جلب داروينية سان عجز نمل عجز نحل عجز جهاز نعي عمل عدن كطر اسم عار جاويد يلج كتب قسم كير نشر اسم عار هار
يحى اسم عار يثك اسم نبي شرة نكرى نبي هار نبي يحى الل حرب فكر لحد كتب ضحي حيو ترن أنفس اذج سلك حيو يثر غرب دهش بعث أمل سلك نطق حكم عقل
ويب وعي خرق نظر عمق قرء اذج قفز ذهن سؤل تلي يمكن لحو عرف يملك فكر تجد جبة سؤل كتب ايد خلق الل رحم رحم لهم طرق تسر طلع سمو ارض طوع امر
وبى على يملك قرء حمل روي عرب كبت كتب عرب
```

```
In [108...] mydata.head()
```

	BookTitle	Author	BookSummary
0	التضحية عند الحيوان	هارون[يحيى]	...روي كتب عرب عبر ربط ارج كعرب او لغة عرب ارج حد
1	معجزة الذرة	هارون[يحيى]	...روي كتب عرب عبر ربط ارج كعرب او لغة عرب ارج حد
2	لا تتجاهل	هارون[يحيى]	...روي كتب عرب عبر ربط ارج كعرب او لغة عرب ارج حد
3	الحياة في سبيل الله	هارون[يحيى]	...روي كتب عرب عبر ربط ارج كعرب او لغة عرب ارج حد
4	العظمة في كل مكان	هارون[يحيى]	...روي كتب عرب عبر ربط ارج كعرب او لغة عرب ارج حد

kmean clustering using CountVectorizer

```
In [109...] from sklearn.feature_extraction.text import CountVectorizer
count_vect = CountVectorizer()
bow = count_vect.fit_transform(mydata['BookSummary'].values)
bow
```

```
Out[109...] <5898x15872 sparse matrix of type '<class 'numpy.int64'>'
with 663630 stored elements in Compressed Sparse Row format>
```

```
In [110...] terms = count_vect.get_feature_names()
```

```
In [113...] from sklearn.cluster import KMeans
model = KMeans(n_clusters = 20,init='k-means++',random_state=99)
model.fit(bow)
```

```
Out[113...] KMeans(n_clusters=20, random_state=99)
```

```
In [114...] labels = model.labels_
cluster_center=model.cluster_centers_
```

```
In [115...] cluster_center
```

```
Out[115...] array([[0.          , 0.          , 0.          , ..., 0.          , 0.          ,
0.          ],
[0.          , 0.          , 0.          , ..., 0.          , 0.          ,
0.          ],
[0.          , 0.          , 0.          , ..., 0.          , 0.          ,
0.          ],
...,
```

```
[0.      , 0.      , 0.      , ..., 0.      , 0.      ,
 0.      ],
[0.      , 0.      , 0.      , ..., 0.      , 0.      ,
 0.      ],
[0.00355872, 0.00355872, 0.00355872, ..., 0.      , 0.      ,
 0.      ]])
```

In [116... terms[1:10]

Out[116... ['ءات', 'ءاتئ', 'ءاذ', 'ءام', 'ءامئ', 'ءة', 'ءله', 'ءمر', 'ءمن']

In [117... `from sklearn import metrics`
`silhouette_score = metrics.silhouette_score(bow, labels, metric='euclidean')`

In [118... *# which tells us that clusters are far away from each other*
`silhouette_score`

Out[118... 0.18003436494177047

In [119... *# Giving Labels/assigning a cluster to each point/text*
`mydata['Bow Class Labal'] = model.labels_ # the last column you can see the Label numebbers`
`mydata.head(2)`

Out[119...

	BookTitle	Author	BookSummary	Bow Class Labal
0	التضحية عند الحيوان	هارون يحيى	...روي كتب عرب عبر ربط ارخ كعرب او لغة عرب ارخ حد	9
1	معجزة الذرة	هارون يحيى	...روي كتب عرب عبر ربط ارخ كعرب او لغة عرب ارخ حد	9

In [120... `mydata.groupby(['Bow Class Labal'])['BookSummary'].count()`

Out[120... Bow Class Labal

0	64
1	430
2	17
3	33
4	16
5	151
6	84
7	2321
8	96
9	961
10	287
11	62
12	223
13	194
14	110
15	197
16	116
17	28
18	227
19	281

Name: BookSummary, dtype: int64

In [124... *#Refrence credit - to find the top 10 features of cluster centriod*
`print("Top terms per cluster:")`
`order_centroids = model.cluster_centers_.argsort()[:, :-1]`
`terms = count_vect.get_feature_names()`
`for i in range(20):`
 `print("Cluster %d:" % i, end='')`
 `for ind in order_centroids[i, :10]:`
 `print(' %s-' % terms[ind], end='')`
 `print()`

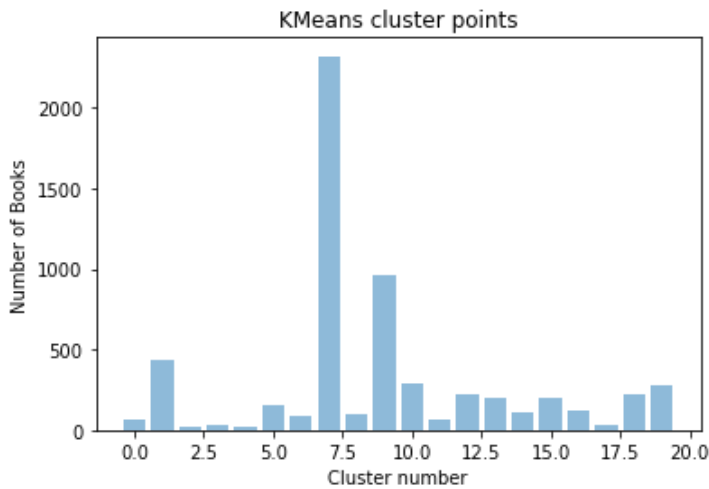
Top terms per cluster:
Cluster 0: ادب - ولد - قال - عرب - حفظ - الي - لسف - نصر - انس - كتب
Cluster 1: روي - كتب - عرب - ادب - شهر - حمد - حفظ - بدى - خلد - وفق
Cluster 2: ون - ين - ال - لل - ان - يم - اب - لر - ور - يه

Cluster 3: -كتب- قطب- سيد- سلم- فكر- رجل- ادب- جمع- شهد- اخو
 Cluster 4: -لل- ين- ون- يم- ان- ام- ال- ول- اب- يل
 Cluster 5: -سلم- كتب- فكر- حمد- عمر- دكتور- علم- انه- جمع- حدث
 Cluster 6: -علم- كتب- نفس- درس- عام- بشر- رود- عقل- عرف- جمع
 Cluster 7: -كتب- علم- قرء- حمل- حمد- جمع- سلم- عام- يمك- عمل
 Cluster 8: -كتب- عرب- عام- الل- ارخ- حمد- روي- رسل- قصص- سلم
 Cluster 9: -كتب- عرب- ارخ- علم- عام- جمع- سلم- روي- حمد- درس
 Cluster 10: -كتب- سلم- حدث- جمع- حمد- علم- الل- كثر- دمشق- عيد
 Cluster 11: -حكم- كتب- سرح- وفق- حيث- شهر- ادب- دبي- عمل- كانت
 Cluster 12: -يكي- عدد- ومك- طفل- لئر- ديز- بيطط- كتب- تحر- جال
 Cluster 13: -كتب- دين- علم- جلل- سيط- علماء- حفظ- طلب- فسر- كثر
 Cluster 14: -كتب- شيخ- قرضاو- سلم- وسف- علم- حصل- تمك- سنة- عصر
 Cluster 15: -رمض- يوم- شهر- سلم- مسك- بدأ- فطر- وقت- وفق- رؤء
 Cluster 16: -كتب- غزل- حمد- علم- ولد- ه- طوس- دين- فقه- صوف
 Cluster 17: -ون- ين- لل- ال- ان- يم- اب- نت- لس- يه
 Cluster 18: -كتب- حمد- سلم- علم- لغز- لسل- نظر- عدد- ختخ- نقل
 Cluster 19: -الل- سلم- سعد- علم- كتب- ولف- عيد- دكتور- دين- قال

```

In [132... # visually how points or reviews are distributed across 10 clusters
import matplotlib.pyplot as plt
plt.bar([x for x in range(20)], mydata.groupby(['Bow Class Labal'])['BookSummary'].count(), alpha = 0.5)
plt.title('KMeans cluster points')
plt.xlabel("Cluster number")
plt.ylabel("Number of Books")
plt.show()

```



kmean clustering using TF-IDF

```

In [128... from sklearn.feature_extraction.text import TfidfVectorizer
tfidf_vect = TfidfVectorizer()
tfidf = tfidf_vect.fit_transform(mydata['BookSummary'].values)
tfidf.shape

```

Out[128... (5898, 15872)

```

In [129... from sklearn.cluster import KMeans
model_tf = KMeans(n_clusters = 20, random_state=99)
model_tf.fit(tfidf)

```

Out[129... KMeans(n_clusters=20, random_state=99)

```

In [130... labels_tf = model_tf.labels_
cluster_center_tf=model_tf.cluster_centers_

```

```

In [131... cluster_center_tf

```

```

Out[131... array([[0., 0., 0., ..., 0., 0., 0.],
        [0., 0., 0., ..., 0., 0., 0.],
        [0., 0., 0., ..., 0., 0., 0.],
        ...,

```



```
[0., 0., 0., ..., 0., 0., 0.],
[0., 0., 0., ..., 0., 0., 0.],
[0., 0., 0., ..., 0., 0., 0.]])
```

```
In [133... # to understand what kind of words generated as columns by BOW
terms1 = tfidf_vect.get_feature_names()
```

```
In [134... terms1[1:10]
```

```
Out[134... ['اات', 'ااتد', 'ااذ', 'ام', 'امد', 'اة', 'له', 'مر', 'من']
```

```
In [135... from sklearn import metrics
silhouette_score_tf = metrics.silhouette_score(tfidf, labels_tf, metric='euclidean')
```

```
In [136... silhouette_score_tf
```

```
Out[136... 0.14474513124824395
```

```
In [137... mydata1 = mydata
mydata1['TF-IDF class label'] = model_tf.labels_
mydata1.head(5)
```

```
Out[137...
      BookTitle      Author      BookSummary  Bow Class Labal  TF-IDF class label
0  التضحية عند الحيوان  [هارون
  يحيى]  ...روي كتب عبر ربط ارخ كعرب او لغة عرب ارخ حد      9      3
1  معجزة الذرة  [هارون
  يحيى]  ...روي كتب عبر ربط ارخ كعرب او لغة عرب ارخ حد      9      3
2  لا تتجاهل  [هارون
  يحيى]  ...روي كتب عبر ربط ارخ كعرب او لغة عرب ارخ حد      9      3
3  الحياة في سبيل الله  [هارون
  يحيى]  ...روي كتب عبر ربط ارخ كعرب او لغة عرب ارخ حد      9      3
4  العظمة في كل مكان  [هارون
  يحيى]  ...روي كتب عبر ربط ارخ كعرب او لغة عرب ارخ حد      9      3
```

```
In [138... # How many points belong to each cluster ->

mydata1.groupby(['TF-IDF class label'])['BookSummary'].count()
```

```
Out[138... TF-IDF class label
0      119
1      369
2      116
3     2332
4        79
5      293
6      224
7        62
8      288
9        96
10     194
11     131
12        98
13     100
14     141
15     147
16        36
17     627
18     248
19     198
Name: BookSummary, dtype: int64
```

```
In [139... #Refrence credit - to find the top 10 features of cluster centriod
print("Top terms per cluster:")
order_centroids = model_tf.cluster_centers_.argsort()[:, :-1]
```

```

for i in range(20):
    print("Cluster %d:" % i, end='')
    for ind in order_centroids[i, :10]:
        print(' %s-' % terms1[ind], end='')
    print()

```

Top terms per cluster:

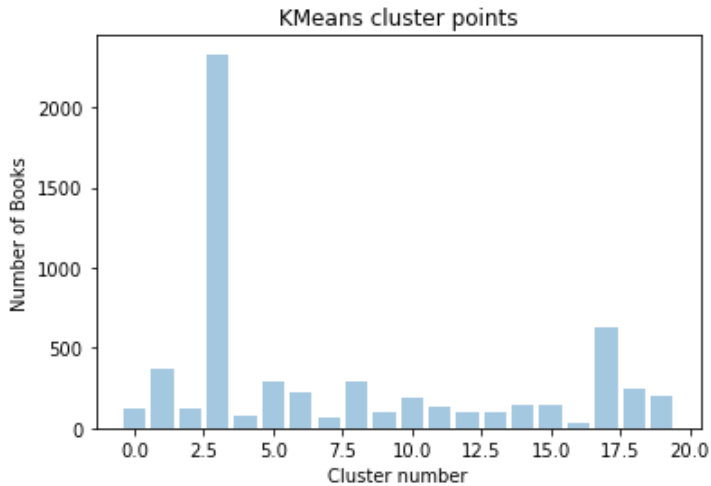
Cluster 0: دمشق- يمي- فتى- سنة- بلغ- عصر- تقى- شيخ- حنبل- جمع
Cluster 1: روي- خلد- بدي- ادب- وفق- كتب- عرب- لسل- نجب- الك
Cluster 2: غزل- طوس- صوف- هـ- حمد- اسماعيلي- شفع- علم- فقه- مات
Cluster 3: كتب- علم- سلم- حمد- روي- الل- جمع- ارخ- قرأ- شيخ
Cluster 4: كتب- فقي- راهيم- حسن- ندي- كيت- كند- ندق- لخص- عرب
Cluster 5: دكتور- نبل- سلم- فكر- فرق- لطب- كتب- عمر- عام- انه
Cluster 6: يكي- ديز- ومك- لثر- بطط- تحر- طفل- كة- رسم- جال
Cluster 7: سرح- حكم- وفق- درا- حيت- دبي- بصف- كتب- ثوب- قلد
Cluster 8: لغز- ختخ- حمد- سحل- برا- لسل- كتب- سلم- حرب- اسكندرية
Cluster 9: الب- حدث- نبي- ضعف- صحح- حاديث- كتب- امم- نصر- حمد
Cluster 10: جلل- سيط- دين- كتب- علم- فسر- علماء- طلب- حفظ- اعت
Cluster 11: ين- ون- حنبل- لل- اب- ذهب- جهد- يم- وب- هـ
Cluster 12: فى- وسع- سلم- أحكام- نبي- عجز- ليل- صلى- فقه- وقع
Cluster 13: صلب- كتب- عقد- حرب- ارخ- زكر- عيس- عضو- وهى- حدث
Cluster 14: يمك- حمل- قرء- وكل- قلب- ثعب- ياب- كتب- عمل- عجز
Cluster 15: نفس- علم- كتب- انس- درس- نصر- رود- لسف- سلك- فرع
Cluster 16: نجي- شرح- يدو- دكتور- ورق- تجد- تيب- فيسيولوجي- قنت- بشت
Cluster 17: كتب- عرب- ارخ- روي- درس- عام- سلم- جمع- فكر- حمد
Cluster 18: فسر- سور- كثر- دمشق- كتب- الل- كشك- سلم- هـ- بصر
Cluster 19: رمض- مسك- يوم- فطر- شهر- رؤة- ايو- توقع- بدأ- وقت

In [140.. *# visually how points or reviews are distributed across 10 clusters*

```

plt.bar([x for x in range(20)], mydata1.groupby(['TF-IDF class label'])['BookSummary'].count(), alpha = 0.4)
plt.title('KMeans cluster points')
plt.xlabel("Cluster number")
plt.ylabel("Number of Books")
plt.show()

```



In []:

In []: