

In []:

```
1 import requests
2 from bs4 import BeautifulSoup
3 import lxml
4 import csv
5 import re
```

In []:

```
1 def remove_tags(input_string):
2     clean_string = re.sub(r'<.*?>', '', input_string)
3     return clean_string
4 def remove_extra_spaces(input_string):
5     return ' '.join(input_string.split())
```

In []:

```
1 file = open('Content4.csv', mode='a', encoding='utf-8', newline='')
2 writer = csv.writer(file)
3 writer.writerow(['Article_name', 'Year', 'Content'])
4
5 url='https://www.jmir.org/'
6 for j in range(2004,2024):
7     print(url+str(j))
8     links = url+str(j)
9     response = requests.get(url=links, stream=True)
10    for i in range(20):
11        # имя статьи
12        soup4 = BeautifulSoup(response.text, 'lxml')
13        soup4 = soup4.find_all('p', class_='h4 full-width-card-info-title')
14        z = (soup4)[i]
15        z = z.text # Здесь имя твое...
16        #=====
17        # Берем HXML так-же из списка
18        soup3 = BeautifulSoup(response.text, 'lxml')
19        soup3 = soup3.find_all('div', class_='full-width-card-info-group-buttons')
20        name = soup3[i]
21        a = str(name).rsplit('data-v-3802195d="" href=')[2]
22        a = a.split('target',1)[0].strip().replace('"','')
23        # Делаем запрос
24        response2 = requests.get(a)
25        soup2 = BeautifulSoup(response2.text, 'lxml')
26        b = str(soup2)
27        b = remove_tags(b)
28        b = remove_extra_spaces(b)
29        writer.writerow([z,j,b])
30        print(f'Скаченная статья {j} года № {i}')
31 file.close()
```

In [1]:

```

1 import pandas as pd
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from sklearn.cluster import KMeans
4 import matplotlib.pyplot as plt
5 import seaborn as sns

```

In [2]:

```

1 df = pd.read_csv('Content4.csv')
2 df.head()

```

Out[2]:

	Article_name	Year	Content
0	Using Claims Data to Examine Patients Using Pr...	2004	JMIR J Med Internet Res Journal of Medical Int...
1	Online Consumer Surveys as a Methodology for A...	2004	JMIR J Med Internet Res Journal of Medical Int...
2	A Multimedia Interactive Education System for ...	2004	JMIR J Med Internet Res Journal of Medical Int...
3	DietPal: A Web-Based Dietary Menu-Generating a...	2004	JMIR J Med Internet Res Journal of Medical Int...
4	Sex Differences in Youth-Reported Depressive S...	2004	JMIR J Med Internet Res Journal of Medical Int...

In [80]:

```
1 # List(df['Article_name'])
```

In [83]:

```

1 # Загрузка и предобработка данных
2
3 # corpus = df["Content"].tolist()
4 corpus = df[["Content", 'Article_name']]

```

In [84]:

```

1
2 # Вычисление TF-IDF векторов
3 vectorizer = TfidfVectorizer(stop_words="english")
4 X = vectorizer.fit_transform(corpus)

```

In [105]:

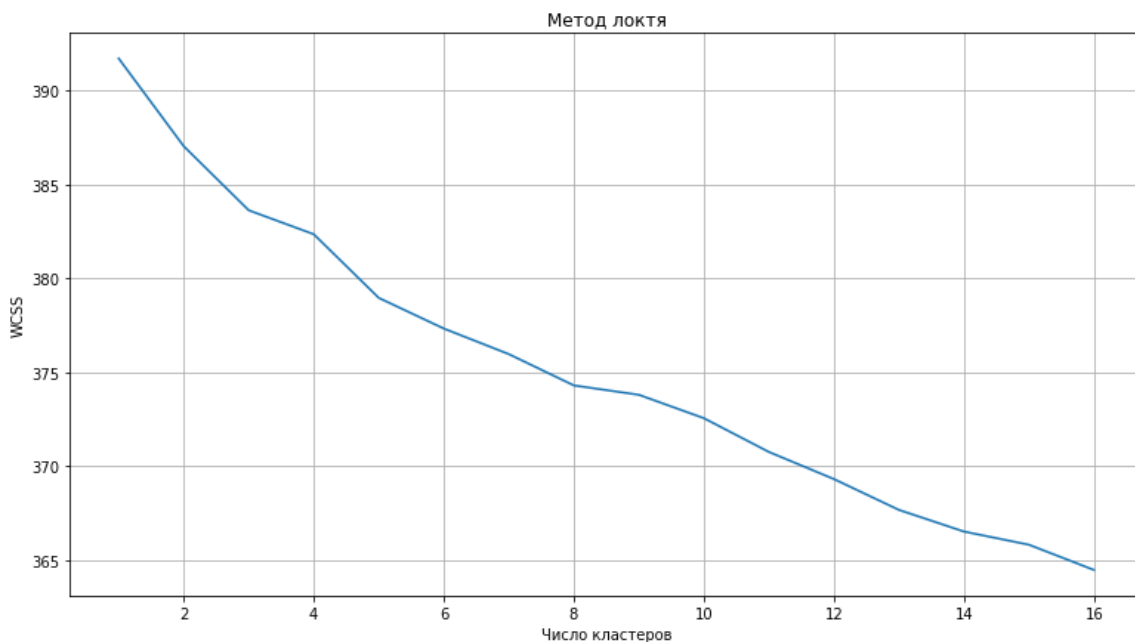
```
1 corpus = df['Article_name'].str.lower()
2 corpus.head()
```

Out[105]:

```
0    using claims data to examine patients using pr...
1    online consumer surveys as a methodology for a...
2    a multimedia interactive education system for ...
3    dietpal: a web-based dietary menu-generating a...
4    sex differences in youth-reported depressive s...
Name: Article_name, dtype: object
```

In [123]:

```
1 corpus = df['Article_name'].str.lower()
2 # Вычисление TF-IDF векторов
3 vectorizer = TfidfVectorizer(stop_words= m, analyzer='word', ngram_range=(1,1))
4 X = vectorizer.fit_transform(corpus)
5 # Задаем максимальное количество кластеров
6 max_clusters = 16
7 # Создаем список для значений WCSS
8 wcss = []
9
10 # Вычисляем WCSS для каждого числа кластеров от 1 до max_clusters
11 for i in range(1, max_clusters + 1):
12     kmeans = KMeans(n_clusters=i, init='k-means++', random_state=0)
13     kmeans.fit(X)
14     wcss.append(kmeans.inertia_)
15
16 # Строим график зависимости WCSS от числа кластеров
17 plt.figure(figsize=(13, 7))
18 plt.plot(range(1, max_clusters + 1), wcss)
19
20 plt.title('Метод локтя')
21 plt.xlabel('Число кластеров')
22 plt.ylabel('WCSS')
23 plt.grid(True)
24 plt.show()
```



In [21]:

```

1 # Обучение KMeans модели
2 k = 8 # количество кластеров
3 kmeans = KMeans(n_clusters=k, random_state=42).fit(X)
4 # Вывод результатов кластеризации
5 labels = kmeans.labels_
6 df["cluster"] = labels
7 df.head(7)

```

Out[21]:

	Article_name	Year	Content	слова	cluster
0	Using Claims Data to Examine Patients Using Pr...	2004	JMIR J Med Internet Res Journal of Medical Int...	[using, claims, data, to, examine, patients, u...	1
1	Online Consumer Surveys as a Methodology for A...	2004	JMIR J Med Internet Res Journal of Medical Int...	[online, consumer, surveys, as, a, methodology...	0
2	A Multimedia Interactive Education System for ...	2004	JMIR J Med Internet Res Journal of Medical Int...	[a, multimedia, interactive, education, system...	4
3	DietPal: A Web-Based Dietary Menu-Generating a...	2004	JMIR J Med Internet Res Journal of Medical Int...	[dietpal:, a, web-based, dietary, menu-generat...	7
4	Sex Differences in Youth-Reported Depressive S...	2004	JMIR J Med Internet Res Journal of Medical Int...	[sex, differences, in, youth-reported, depress...	4
5	Can Clinical Trials Requiring Frequent Partici...	2004	JMIR J Med Internet Res Journal of Medical Int...	[can, clinical, trials, requiring, frequent, p...	5
6	Online Pediatric Information Seeking Among Mot...	2004	JMIR J Med Internet Res Journal of Medical Int...	[online, pediatric, information, seeking, amon...	4

In [30]:

```
1 df['Article_name'][df['cluster'] == 7].head(10).values.tolist()
```

Out[30]:

```
['DietPal: A Web-Based Dietary Menu-Generating and Management System',  
'Providing a Web-based Online Medical Record with Electronic Communicatio  
n Capabilities to Patients With Congestive Heart Failure: Randomized Tria  
l',  
'Usage and Longitudinal Effectiveness of a Web-Based Self-Help Cognitive  
Behavioral Therapy Program for Panic Disorder',  
'Defining Participant Exposure Measures in Web-Based Health Behavior Chan  
ge Programs',  
'Rates and Determinants of Repeated Participation in a Web-Based Behavior  
Change Program for Healthy Body Weight and Healthy Lifestyle',  
'Effectiveness of a Web-Based Self-Help Intervention for Symptoms of Depr  
ession, Anxiety, and Stress: Randomized Controlled Trial',  
'Randomized Controlled Trial of an Internet-Based Versus Face-to-Face Dys  
pnea Self-Management Program for Patients With Chronic Obstructive Pulmona  
ry Disease: Pilot Study',  
'Therapist-Assisted, Internet-Based Treatment for Panic Disorder: Can Gen  
eral Practitioners Achieve Comparable Patient Outcomes to Psychologists?',  
'Integrating an eHealth Program for Pregnant Women in Midwifery Care: A F  
easibility Study Among Midwives and Program Users',  
'Comparison of Trial Participants and Open Access Users of a Web-Based Ph  
ysical Activity Intervention Regarding Adherence, Attrition, and Repeated  
Participation']
```