

به نام خدا

```
1 import re
2 import nltk
3 import pandas as pd
4
5
6 def tokenize(text_array):
7     for i in range(len(text_array)):
8         if not pd.isna(text_array.loc[i]):
9             text_array.loc[i] = text_array.loc[i].replace('\n', ' ')
10            text_array.loc[i] = re.sub(r'^https?:\/\/\.[^\s]*$', '', text_array.loc[i], flags=re.MULTILINE)
11            text_array.loc[i] = re.sub(r'[0-9]+', " ", text_array.loc[i])
12            text_array.loc[i] = re.sub(r'^a-zA-Z0-9', " ", text_array.loc[i].lower())
13            text_array.loc[i] = nltk.word_tokenize(text_array.loc[i])
14
15     return text_array
16
17 data_csv = pd.read_csv('dataset/train.csv')
18
19 pure_tokens = tokenize(data_csv['Comment'])
20
21 pure_tokens_dataframe = pd.DataFrame(pure_tokens)
22 pure_tokens_dataframe.to_csv('./one_pure_tokens.csv', index=False)
23
```

در این تابع برای ستون کامنت لینک ها، \n ها و کلماتی که با یک عدد شروع می شوند را جدا کردیم و سپس جملات را به کلمات توکنایز کردیم.

```
1 import nltk
2 import pandas as pd
3
4
5 def stopword_remover(text_array):
6     for i in range(len(text_array)):
7         if not pd.isna(text_array.loc[i]):
8             text_array.loc[i] = (text_array.loc[i].replace("'", "", 1).replace('"', "", 1).split(", ")).split(" ")
9     for i in range(len(text_array)):
10        if type(text_array.loc[i]) != float:
11            text_array.loc[i] = [word for word in text_array.loc[i] if word not in nltk.corpus.stopwords.words('english')]
12    return text_array
13
14
15 pure_tokens_csv = pd.read_csv('step1/one_pure_tokens.csv')
16
17 tokens_without_stopwords = stopword_remover(pure_tokens_csv['Comment'])
18
19 tokens_without_stopwords_dataframe = pd.DataFrame(tokens_without_stopwords)
20 tokens_without_stopwords_dataframe.to_csv('./two_tokens_without_stopwords.csv', index=False)
21
```

در این تابع stop word ها را پاکسازی کردیم.

```

1 import nltk
2 import pandas as pd
3
4
5 def stemming(text_array):
6     for i in range(len(text_array)):
7         if not pd.isna(text_array.loc[i]):
8             text_array.loc[i] = (text_array.loc[i]).replace("'", "", 1).replace('"', "", 1).split("'", "")
9     for i in range(len(text_array)):
10         if type(text_array.loc[i]) != float:
11             text_array.loc[i] = [nltk.PorterStemmer().stem(w) for w in text_array.loc[i]]
12     return text_array
13
14
15 tokens_without_stopwords_csv = pd.read_csv('step2/two_tokens_without_stopwords.csv')
16
17 tokens_with_stemming = stemming(tokens_without_stopwords_csv['Comment'])
18
19 tokens_with_stemming_dataframe = pd.DataFrame(tokens_with_stemming)
20 tokens_with_stemming_dataframe.to_csv('./three_tokens_with_stemming.csv', index=False)
21

```

در این تابع که استفاده نشده است ولی تست و بررسی شده است بر روی داده‌ها
stemming انجام داده‌ایم ولی با توجه به نتیجه نچندان خوب از
lemmetization استفاده کردیم.

```

1 import pandas as pd
2 from nltk.stem import WordNetLemmatizer
3
4
5
6 def lemmatizing(text_array):
7     for i in range(len(text_array)):
8         if not pd.isna(text_array.loc[i]):
9             text_array.loc[i] = (text_array.loc[i]).replace("'", "", 1).replace('"', "", 1).split("'", "")
10     for i in range(len(text_array)):
11         if type(text_array.loc[i]) != float:
12             text_array.loc[i] = [WordNetLemmatizer().lemmatize(w) for w in text_array.loc[i]]
13     return text_array
14
15
16 tokens_without_stopwords_csv = pd.read_csv('step2/two_tokens_without_stopwords.csv')
17
18 tokens_with_lemmatizing = lemmatizing(tokens_without_stopwords_csv['Comment'])
19
20 tokens_with_lemmatizing_dataframe = pd.DataFrame(tokens_with_lemmatizing)
21 tokens_with_lemmatizing_dataframe.to_csv('./three_tokens_with_lemmatizing.csv', index=False)
22

```

در این تابع از **lemmetizing** برای پیدا کردن ریشه کلمات استفاده کردیم.

```

1 import gensim
2 import pandas as pd
3
4
5 def word2Vec(text_array):
6     for i in range(len(text_array)):
7         if not pd.isna(text_array.loc[i]):
8             text_array.loc[i] = (text_array.loc[i]).replace("'", "", 1).replace('"', "", 1).replace(",", " ", " ")
9     text_array = text_array.apply(gensim.utils.simple_preprocess)
10    # print(text_array)
11    model = gensim.models.Word2Vec(window=3, min_count=3, workers=4)
12    model.build_vocab(text_array, progress_per=10000)
13    model.train(text_array, total_examples=model.corpus_count, epochs=70)
14    model.save('word2vec.model')
15
16    print(model.wv.get_vector('good'))
17
18    # for i in range(len(text_array)):
19    #     if type(text_array.loc[i]) != float:
20    #         text_array.loc[i]
21    #         print(str((i/13741)*100)+"%\r")
22    return text_array
23
24
25 tokens_with_lemmatizing_csv = pd.read_csv('step3/three_tokens_with_lemmatizing.csv')
26
27 extracted_key_tokens = word2Vec(tokens_with_lemmatizing_csv['Comment'])
28
29 extracted_key_tokens_dataframe = pd.DataFrame(extracted_key_tokens)
30 extracted_key_tokens_dataframe.to_csv('./four_extracted_key_tokens.csv', index=False)
31

```

در این تابع از word2vec آموزش داده شده مخصوص داده‌های train خودمان برای برداری کردن کلمات و سپس با گرفتن میانگین وزن دار کلمات هر داکيومنت برای برداری کردن داکيومنت ها استفاده کردیم. (این تابع به علت ضعیف‌تر بودن در برداری کردن داکيومنت ها کنار گذاشته شده و از تابع بعدی برای ادامه الگوریتم استفاده شده است)

```

1 import gensim
2 import pandas as pd
3 from nltk.tokenize import word_tokenize
4
5
6 def doc2Vec(text_array):
7     for i in range(len(text_array)):
8         if not pd.isna(text_array.loc[i]):
9             text_array.loc[i] = (text_array.loc[i]).replace("'", "", 1).replace('"', "", 1).replace(",", " ", " ")
10
11    text_array = [gensim.models.doc2vec.TaggedDocument(words=word_tokenize(_d.lower()), tags=[str(i)]) for i, _d in enumerate(text_array)]
12    documents_tags_dataframe = pd.DataFrame(text_array)
13    documents_tags_dataframe.to_csv('./documents_tags.csv', index=False)
14
15    model = gensim.models.Doc2Vec(dm=1, vector_size=65, hs=1, min_count=2, sample = 12000, window=3, alpha=0.025, min_alpha=0.00025)
16    model.build_vocab(text_array)
17    model.train(text_array, total_examples=model.corpus_count, epochs=70)
18    model.save('doc2vec.model')
19
20    return text_array
21
22
23 tokens_with_lemmatizing_csv = pd.read_csv('step3/three_tokens_with_lemmatizing.csv')
24
25 extracted_key_tokens = doc2Vec(tokens_with_lemmatizing_csv['Comment'])
26
27 extracted_key_tokens_dataframe = pd.DataFrame(extracted_key_tokens)
28 extracted_key_tokens_dataframe.to_csv('./four_point_one_extracted_key_tokens.csv', index=False)
29

```

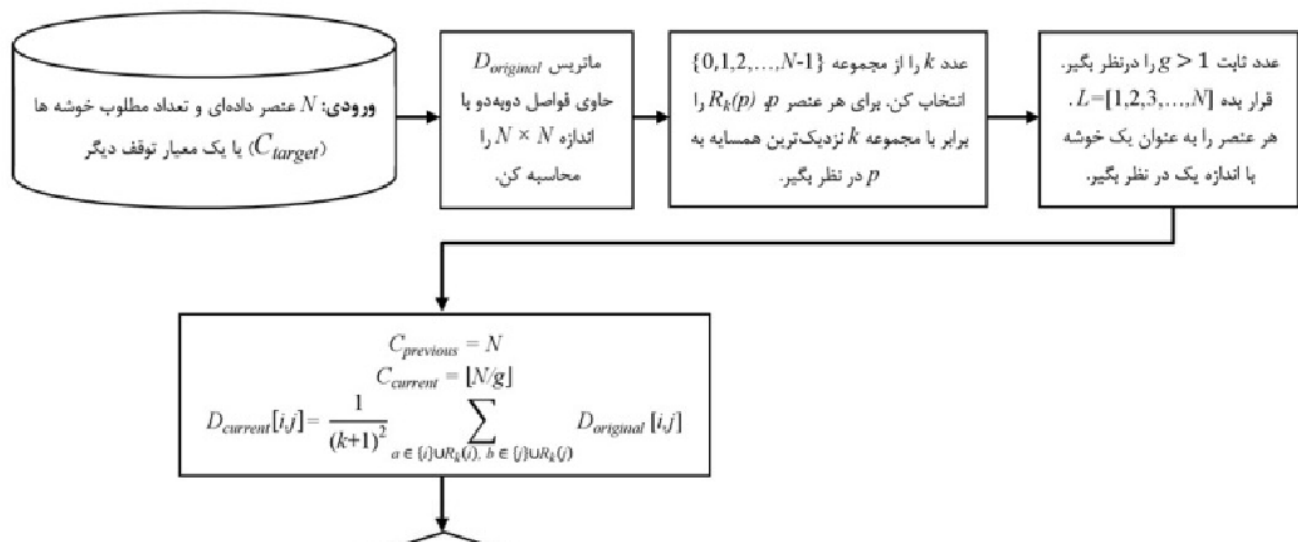
در این تابع با توجه به داکيومنت های موجود آن ها را با الگوریتم ، Doc2vec train کردیم و بردار هایی به ازای هر مقاله ساخته شد.

```
8 def save_function(original_cosine_distance, current_distance, indexes, document_cluster_labels, n_cluster_previous, n_cluster_current):
9     np.save('./original_cosine_distance', original_cosine_distance)
10    np.save('./current_distance', current_distance)
11    np.save('./indexes', indexes)
12    np.save('./first_current_distance_of_vectors', current_distance)
13
```

در این تابع مقادیر مرحله اول کلاسترینگ به علت زمان بر بودن و حجم زیاد ذخیره شدند تا در مرحله بعد این زمان از زمان اجرای ادامه کلاسترینگ کسر گردد.

```
14 def key_identification_clustering_one(document_vectors, n_clusters):
15
16     original_cosine_distance = np.zeros(shape=(8695, 8695))
17
18     for i in range(len(document_vectors.distances(document_vectors[0]))):
19         #calc cosine distances
20         original_cosine_distance[i] = document_vectors.distances(document_vectors[i])
21
22     original_cosine_distance = np.where(original_cosine_distance < 0, 0, original_cosine_distance)
23
24     #calc nearest neighbors
25     nearest_neighbors = NearestNeighbors(n_neighbors=k, metric="precomputed")
26     nearest_neighbors.fit(original_cosine_distance)
27
28     #Fetch kneighbors
29     distances, indexes = nearest_neighbors.kneighbors()
30     document_cluster_labels = [x for x in range(1, len(original_cosine_distance)+1)]
31     n_cluster_previous = len(Counter(document_cluster_labels).keys())
32     n_cluster_current = math.floor(n_cluster_previous/g)
33     current_distance = np.zeros(shape=(n_cluster_previous, n_cluster_previous))
34
35     for i in range(n_cluster_previous):
36         for j in range(n_cluster_previous):
37             current_distance_sigma = 0
38             for l in range(k):
39                 for p in range(k):
40                     current_distance_sigma += original_cosine_distance[indexes[i][l]][indexes[j][p]]
41
42             current_distance[i][j] = current_distance_sigma/((k+1)**2)
43             print("% ", ((i)/(n_cluster_previous))*100)
44         print(current_distance)
45
46     save_function(original_cosine_distance, current_distance, indexes, document_cluster_labels, n_cluster_previous, n_cluster_current)
47     return document_vectors
48
49 k = 5
50 g = 40
51 n_cluster_target = 3
52
53 model = gensim.models.doc2vec.Doc2Vec.load("step4/doc2vec.model")
54
55 key_identification_clustering_one(model.dv, n_cluster_target)
```

در این تابع که قسمت اول کلاسترینگ را محاسبه می کند یعنی این قسمت از فلوچارت را :



با استفاده از تابع فاصله کسینوسی فواصل بین داده‌ها و کلاسترها محاسبه شده است.

```

14 def load_function():
15     original_cosine_distance = np.load('step5_1/original_cosine_distance.npy')
16     knn_indexes = np.load('step5_1/indexes.npy')
17     document_cluster_labels = [x for x in range(len(original_cosine_distance))]
18     n_cluster_previous = len(Counter(document_cluster_labels).keys())
19     n_cluster_current = math.floor(n_cluster_previous/g)
20     current_distance = np.load('step5_1/first_current_distance_of_vectors.npy')
21     return original_cosine_distance, current_distance, knn_indexes, document_cluster_labels, n_cluster_previous, n_cluster_current

```

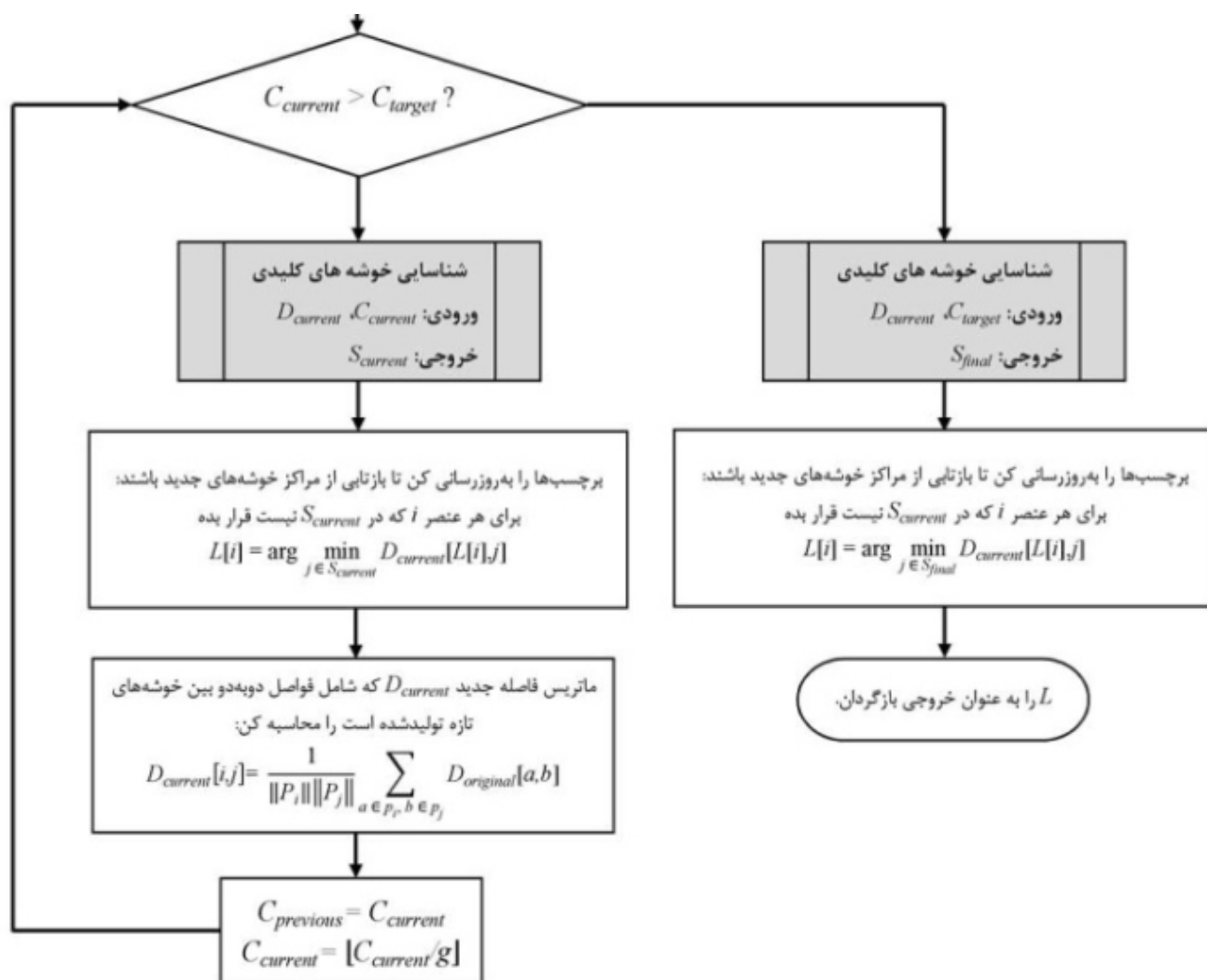
در این تابع متغیرهای ذخیره شده در قسمت اول الگوریتم load می‌شود.

```

23 def key_item_selection(current_distance_matrix, n_clusters):
24
25     mean_distance_matrix = np.mean(current_distance_matrix, axis=1)
26
27     #s in article
28     first_selected_key = np.argmin(mean_distance_matrix)
29     selected_keys = []
30     selected_keys.append(first_selected_key)
31
32     #k in article
33     unselected_items_for_key = [x for x in range(0, len(current_distance_matrix)) if x != first_selected_key]
34     temp_min_key_j = -1
35     temp_max_key_i = -1
36     n = len(selected_keys)
37
38     while n < n_clusters:
39         temp_max_distance = 0
40         for i in unselected_items_for_key:
41             temp_min_distance = 1000000000
42             for j in selected_keys:
43                 if current_distance_matrix[i][j] < temp_min_distance:
44                     temp_min_distance = current_distance_matrix[i][j]
45                     temp_min_key_j = j
46             if current_distance_matrix[i][temp_min_key_j] > temp_max_distance:
47                 temp_max_distance = current_distance_matrix[i][temp_min_key_j]
48                 temp_max_key_i = i
49             selected_keys.append(temp_max_key_i)
50             unselected_items_for_key.remove(temp_max_key_i)
51             n = n + 1
52             # print("key_item_selection : %", ((n/n_clusters))*100)
53     print(len(selected_keys))
54     print(selected_keys)
55
56     return selected_keys

```

در این تابع item های کلیدی طبق الگوریتم مقاله انتخاب می شوند برای شناخته شدن به عنوان مرکز کلاسترها



در این تابع که ادامه فایل قبل است این قسمت از فلوچارت در شکل زیر اجرا می شود

```

def key_identification_clustering_two(document_vectors, n_cluster_target):

    original_cosine_distance, current_distance, knn_indexes, document_cluster_labels, n_cluster_previous, n_cluster_current = load_function()

    while n_cluster_current > n_cluster_target:

        selected_keys = key_item_selection(current_distance, n_cluster_current)

        # update elements labels

        for i in range(len(document_cluster_labels)):
            min_distance = MAXINT
            for j in selected_keys:
                if current_distance[document_cluster_labels[i]][j] < min_distance:
                    min_distance = current_distance[document_cluster_labels[i]][j]
                    new_cluster_label = j
            document_cluster_labels[i] = new_cluster_label

        # update current distance matrix
        # p list in article
        counter = 0
        cluster_members_with_neighbors = {}
        for i in selected_keys:
            temp = list(locate(document_cluster_labels, lambda a: a == i))
            for l in temp:
                for j in range(k):
                    temp.append(knn_indexes[l][j])
                temp = list(set(temp))
            cluster_members_with_neighbors[counter] = temp
            counter += 1

        current_new_distance = np.zeros(shape=(n_cluster_current, n_cluster_current))
        for i in range(n_cluster_current):
            for j in range(n_cluster_current):
                current_distance_sigma = 0

                for l in cluster_members_with_neighbors[i]:
                    for p in cluster_members_with_neighbors[j]:
                        current_distance_sigma += original_cosine_distance[l][p]

                current_new_distance[i][j] = current_distance_sigma / (len(cluster_members_with_neighbors[i]) * len(cluster_members_with_neighbors[j]))

            print("update current distance matrix %", ((i)/(n_cluster_current))*100)

        # update n_cluster_previous & n_cluster_current
        n_cluster_previous = n_cluster_current
        n_cluster_current = math.floor(n_cluster_current/g)
        print("^^^")
        current_distance = current_new_distance

        c = 0
        cluster_label_dict_current = {}
        for i in selected_keys:
            cluster_label_dict_current[i] = c
            c += 1

        for i in range(len(document_cluster_labels)):
            document_cluster_labels[i] = cluster_label_dict_current[document_cluster_labels[i]]

        selected_keys = key_item_selection(current_distance, n_cluster_target)

        # update final elements labels
        for i in range(len(document_cluster_labels)):
            min_distance = MAXINT
            for j in selected_keys:
                if current_distance[document_cluster_labels[i]][j] < min_distance:
                    min_distance = current_distance[document_cluster_labels[i]][j]
                    new_cluster_label = j
            document_cluster_labels[i] = new_cluster_label

        c = 0
        cluster_label_dict_current = {}
        for i in selected_keys:
            cluster_label_dict_current[i] = c
            c += 1

        for i in range(len(document_cluster_labels)):
            document_cluster_labels[i] = cluster_label_dict_current[document_cluster_labels[i]]

    return document_cluster_labels

```


در این تابع که طبق فلوجارت و توضیحات مقاله نوشته شده است الگوریتم مقاله اجرا می شود و در نهایت لیبل های پیش بینی شده برای داده های **train** مشخص می گردد.

در کد صفحه ی بعد دستورات مربوط به اجرای الگوریتم و گرفتن اطلاعاتی راجع به صحت و درصد درستی نتایج داده شده نوشته شده است و همچنین بازنمایی داکيومنت ها در اسکتر پلات برای درک بیشتر نمایش داده شده است. و همچنین این الگوریتم همراه با الگوریتم خوشه بندی **k-means** اجرا شده تا تفاوت دقت این دو الگوریتم بررسی شود.

```

k = 3
g = 40
n_cluster_target = 3

model = gensim.models.doc2vec.Doc2Vec.load("step4/doc2vec.model")

labels = pd.read_csv('dataset/train.csv')

key_identification_clusters = key_identification_clustering_two(model.dv, n_cluster_target)

kmeans = KMeans(n_clusters=n_cluster_target)

li = np.zeros(shape=(8695, 65))
x = np.zeros(shape=(8695, 1))
y = np.zeros(shape=(8695, 1))

for i in range(len(model.dv)):
    li[i] = np.array(model.dv[i])
    x[i] = li[i][0]
    y[i] = li[i][1]

kmeans.fit(li)

#kmeans labels
# for i in range(len(kmeans.labels_)):
#     print(kmeans.labels_[i], labels["Topic"].iloc[i])

# Our labels
for i in range(len(key_identification_clusters)):
    print(key_identification_clusters[i], labels["Topic"].iloc[i])

Biology_0 = 0
Biology_1 = 0
Biology_2 = 0

Chemistry_0 = 0
Chemistry_1 = 0
Chemistry_2 = 0

Physics_0 = 0
Physics_1 = 0
Physics_2 = 0

for i in range(len(key_identification_clusters)):
    if key_identification_clusters[i] == 0:

        if labels["Topic"].iloc[i] == "Biology":
            Biology_0 += 1
        if labels["Topic"].iloc[i] == "Chemistry":
            Chemistry_0 += 1
        if labels["Topic"].iloc[i] == "Physics":
            Physics_0 += 1

    if key_identification_clusters[i] == 1:

        if labels["Topic"].iloc[i] == "Biology":
            Biology_1 += 1
        if labels["Topic"].iloc[i] == "Chemistry":
            Chemistry_1 += 1
        if labels["Topic"].iloc[i] == "Physics":
            Physics_1 += 1

    if key_identification_clusters[i] == 2:

        if labels["Topic"].iloc[i] == "Biology":
            Biology_2 += 1
        if labels["Topic"].iloc[i] == "Chemistry":
            Chemistry_2 += 1
        if labels["Topic"].iloc[i] == "Physics":
            Physics_2 += 1

#kmeans colors
kmeans_colors = kmeans.labels_

#key_identification_cluster colors
key_identification_cluster_colors = np.asarray(key_identification_clusters)

#main colors
main_colors = np.zeros(shape=(8695, 1))

for i in range(len(labels["Topic"])):
    if labels["Topic"].iloc[i] == "Biology":
        main_colors[i] = 0
    if labels["Topic"].iloc[i] == "Chemistry":
        main_colors[i] = 1
    if labels["Topic"].iloc[i] == "Physics":
        main_colors[i] = 2

print("Kmeans Accuracy:", (100*metrics.accuracy_score(main_colors, kmeans_colors)))
print("Key Identification Accuracy:", (100*metrics.accuracy_score(main_colors,
key_identification_cluster_colors)))

print("Biology_0", Biology_0)
print("Biology_1", Biology_1)
print("Biology_2", Biology_2)
print("Chemistry_0", Chemistry_0)
print("Chemistry_1", Chemistry_1)
print("Chemistry_2", Chemistry_2)
print("Physics_0", Physics_0)
print("Physics_1", Physics_1)
print("Physics_2", Physics_2)

plt.scatter(x, y, c=kmeans_colors)
plt.show()
plt.scatter(x, y, c=key_identification_cluster_colors)
plt.show()
plt.scatter(x, y, c=main_colors)
plt.show()

```