

**AI61003 LINEAR ALGEBRA FOR AI AND ML**  
**ASSIGNMENT 01 - PROBLEM 10**

**CASE 1      k = 4**

The document-cluster association produced by the clustering algorithm is as follows.

S.No.	Article Title	Cluster ID
01	<i>Linear Algebra</i>	3
02	<i>Data Science</i>	1
03	<i>Artificial Intelligence</i>	1
04	<i>European Central Bank</i>	4
05	<i>Financial Technology</i>	4
06	<i>International Monetary Fund</i>	4
07	<i>Basketball</i>	2
08	<i>Swimming</i>	2
09	<i>Cricket</i>	2

Cluster ID	Assignments
01	{ <i>Data Science, Artificial Intelligence</i> }
02	{ <i>Basketball, Swimming, Cricket</i> }
03	{ <i>Linear Algebra</i> }
04	{ <i>European Central Bank, Financial Technology, International Monetary Fund</i> }

**CASE 2      k = 8**

The document-cluster association produced by the clustering algorithm is as follows.

S.No.	Article Title	Cluster ID
01	<i>Linear Algebra</i>	7
02	<i>Data Science</i>	1
03	<i>Artificial Intelligence</i>	2
04	<i>European Central Bank</i>	6
05	<i>Financial Technology</i>	3
06	<i>International Monetary Fund</i>	4
07	<i>Basketball</i>	5
08	<i>Swimming</i>	8
09	<i>Cricket</i>	5

Cluster ID	Assignments
01	{ <i>Data Science</i> }

02	{ <i>Artificial Intelligence</i> }
03	{ <i>Financial Technology</i> }
04	{ <i>International Monetary Fund</i> }
05	{ <i>Basketball, Cricket</i> }
06	{ <i>European Central Bank</i> }
07	{ <i>Linear Algebra</i> }
08	{ <i>Swimming</i> }

### CASE 3      **k = 12**

9 objects cannot be clustered into 12 groups.

#### CHOOSING THE BEST CLUSTER SIZE

Let us first try to *manually* divide the nine articles into some compact clusters. Clearly, we can see three broad domains to which the articles belong.

- **Academia** : *Linear Algebra, Data Science, Artificial Intelligence*
- **Finance** : *Financial Technology, European Central Bank, International Monetary Fund*
- **Sports** : *Basketball, Swimming, Cricket*

Since we have to choose the best  $k$  value out of 4 and 8, we must select the one that more prominently expresses the underlying relationship among the articles in the form of the above three categories.

Choosing  $k = 8$  value to cluster 9 objects is a very bad idea because like this every cluster (except one) would get a single object and hence it would not be possible to deduce any meaningful relationship among the objects. This behaviour is prominent when the *k-means algorithm* was run on the feature vectors of the nine articles for  $k = 8$ . As we can see the resulting classifications, the only information we can deduce is that articles on *Basketball* and *Cricket* are similar. Besides that, one can draw no other conclusion about the articles from the produced clustering.

On the other hand, choosing  $k = 4$  gives excellent results. As we can see the resulting classifications, all articles on **sports** and **finance** are grouped under the same cluster. More importantly, the clusters corresponding to these two groups do not contain any other irrelevant articles (no contamination). This means that these two domains were perfectly identified by the algorithm. Now as far as the remaining three objects in the domain of **academia** are concerned, two of them fit inside the same cluster, that is better than all three of them being in different clusters. So it is clear that choosing  $k = 4$  gives a lot more information about the semantic relatedness among the articles that are clustered.

Therefore for this data,  $k = 4$  is the better choice.

```
In [1]: 1 import wikipedia
        2 import numpy as np
        3 from random import shuffle
        4 from sklearn.feature_extraction.text import TfidfVectorizer
        5 import sklearn.cluster as Clustering_Algorithm
```

```
In [2]: 1 def Get_Distance ( a , b ) :
        2     return np.linalg.norm(b-a)
        3
        4 def J_clust ( train , kmeans ) :
        5     s = 0
        6     for img in train :
        7         cluster_id = kmeans.predict(img.astype(np.float).reshape(1,-1))[0]
        8         s += Get_Distance(img, kmeans.cluster_centers_[cluster_id])
        9     return s / len(train)
```

```
In [3]: 1 articles = ['Linear Algebra', 'Data Science', 'Artificial intelligence', 'European Central Bank',
        2           'Financial technology', 'International Monetary Fund', 'Basketball', 'Swimming', 'Cricket']
```

```
In [4]: 1 wiki_content = []
        2 for article in articles :
        3     print('Loading content from article :', article)
        4     wiki_content.append(wikipedia.page(article).content)
```

```
Loading content from article : Linear Algebra
Loading content from article : Data Science
Loading content from article : Artificial intelligence
Loading content from article : European Central Bank
Loading content from article : Financial technology
Loading content from article : International Monetary Fund
Loading content from article : Basketball
Loading content from article : Swimming
Loading content from article : Cricket
```

```
In [5]: 1 vectorizer = TfidfVectorizer(stop_words={'english'})
        2 X = vectorizer.fit_transform(wiki_content).toarray()
```

```
In [6]: 1 def Do_Clustering ( train , k ) :
        2     kmeans = Clustering_Algorithm.KMeans(n_clusters=k, init='random', n_init=1)
        3     kmeans.fit(train)
        4     classes = list(kmeans.labels_)
        5     for i, cl in enumerate(classes) :
        6         print(articles[i], cl+1)
```

In [7]:

1

Do\_Clustering(X, 4)

Linear Algebra 3

Data Science 1

Artificial intelligence 1

European Central Bank 4

Financial technology 4

International Monetary Fund 4

Basketball 2

Swimming 2

Cricket 2

In [8]:

1

Do\_Clustering(X, 8)

Linear Algebra 7

Data Science 1

Artificial intelligence 2

European Central Bank 6

Financial technology 3

International Monetary Fund 4

Basketball 5

Swimming 8

Cricket 5