

Global complexity versus Elbow method to choose the hyper-parameter k of k -means

Ghribi Saif Eddine
Supervised by
Mr Jean-Louis Dessalles
Mr Pierre-Alexandre Murena

March 23, 2019



ATHENS NETWORK
ADVANCED TECHNOLOGY
HIGHER EDUCATION NETWORK

Contents

1	Introduction	3
2	Data used in the project	3
3	Elbow method	3
3.1	Overview	3
3.2	Example using Iris data set	3
4	Global Complexity	4
4.1	Definition	4
4.2	Implementation	4
5	Elbow method and Global Complexity comparison	5
5.1	The main idea	5
5.2	Implementation	5
5.2.1	Function to evaluate Global Complexity and the distortion	5
5.2.2	Function to plot the comparison for a given data name	5
5.3	Results for the four data sets	6
5.4	Commentary	8
6	Conclusion	8
7	References	9

1 Introduction

During the Athens week and thanks to the subject From Complexity to Intelligence (TPT37), I was introduced to the theory of complexity and its relation with machine learning. In the lab of Machine Learning of this subject the idea of choosing the optimal number of clusters with global complexity was well introduced and explained. Thus, since I already worked with k-means in several projects I have always used the elbow method to determine this hyper-parameter. In this short report I will explain and compare these two methods.

2 Data used in the project

In order to compare the two methods I have chosen to use 4 of the widely used data sets in machine learning for testing. These data sets will be all retrieved from SKlearn .

- Sklearn Iris : This data sets consists of 3 different types of irises.
- Sklearn Make Moons : Make two interleaving half circles.
- Sklearn Make blobs : Generate isotropic Gaussian blobs for clustering.
- Sklearn Wine data : The wine dataset is a classic and very easy multi-class classification dataset.

To store the data for further use I created a dictionary in a function called data to retrieve data by its name:

```
1 def data(index):
2     dictionary = {} # dictionary to store data
3     d = load_iris()
4     X = d['data']
5     dictionary["iris"] = X # add iris data (index "iris")
6     X, y = make_moons(noise=.05, random_state=1)
7     dictionary["moons"] = X
8     X, y = make_blobs(random_state=1)
9     dictionary["blobs"] = X
10    d = load_wine()
11    dictionary["wine"] = d['data']
12    return dictionary[index]
```

Listing 1: Function to retrieve data

3 Elbow method

3.1 Overview

The elbow method is used to determine the optimal number of clusters in k-means. if k increases, average distortion will decrease, each cluster will have fewer constituent instances, and the instances will be closer to their respective centroids.

The value of k at which improvement in distortion declines the most is called the elbow, at which we should stop dividing the data into further clusters.

The distortion function here is considered as the sum of squared distances of samples to their closest cluster center.

3.2 Example using Iris data set

In order to plot the elbow method I used a list to store the sum of squared distances of samples to their closest cluster center for a set of integers ,a set of the number of clusters.

```

1 # fix the figure size
2 figure(num=None, figsize=(8, 6), dpi=90, facecolor='w', edgecolor='k')
3 l_elbow, l_kolo = eval_complex(max_clusters, "iris")
4 plt.plot(range(1, max_clusters), l_elbow, label="elbow method")
5 plt.plot(6, l_elbow[5], "o", label="elbow point = 6")
6 plt.xlabel("Number of clusters")
7 plt.ylabel("Sum of squared distances of samples to their closest cluster center")
8 plt.title("Elbow method")
9 plt.legend()

```

Listing 2: Plot Elbow method

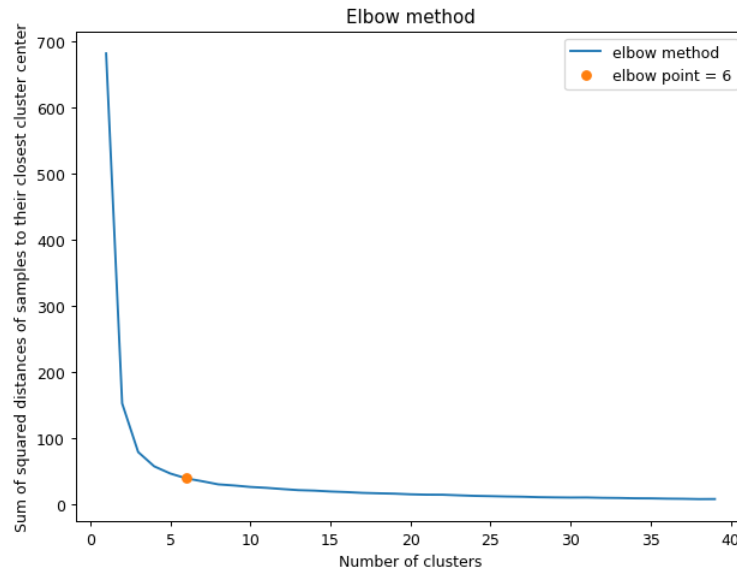


Figure 1: Elbow method

4 Global Complexity

4.1 Definition

The Global complexity was first introduced in the fourth lab. The definition of the global complexity is :

$$GlobalComplexity = Complexity(Data|Prototypes) + Complexity(Prototypes) \quad (1)$$

The prototypes here are the centers of clusters known as centroids. The function that determines this complexity is already implemented by Mr Jean-Louis Dessalles and Mr Pierre-Alexandre Murena as follows :

4.2 Implementation

```

1 def Cmatrix(M) :
2     return np.sum(np.log2(1 + abs_app(M)))
3
4 def gradCmatrix(M) :
5     return np.sign(M) / (1.0 + abs_app(M))
6
7 def C_prototypes(M) :
8     return Cmatrix(M) + np.log2(M.shape[0])
9

```

```

10 def C_data_knowing_prototypes(D, P):
11     NearestPrototype = associatePointToPrototype(D, P)
12     return Cmatrix(D - P[NearestPrototype, :])
13
14 def globalComplexity_clustering(X, P):
15     return C_data_knowing_prototypes(X, P) + C_prototypes(P)

```

Listing 3: Global Complexity implemented by Mr Jean-Louis Dessalles and Mr Pierre-Alexandre Murena

5 Elbow method and Global Complexity comparison

5.1 The main idea

To compare these two methods one shall plot the reaction of the Global Complexity when the hyper-parameter k varies and the distortion function which is the sum of squared distances of samples to their closest cluster center.

5.2 Implementation

5.2.1 Function to evaluate Global Complexity and the distortion

```

1 def eval_complex(max_clusters, data_index):
2     # list to store Sum of squared distances of samples to their closest cluster
3     # center.
4     l_elbow = []
5     # list to store the general complexity
6     l_kolo = []
7     for i in range(1, max_clusters):
8         # choose data
9         X = data(data_index)
10        # fit model (k-means)
11        kmeans = KMeans(n_clusters=i, random_state=1).fit(X)
12        # global complexity after K-Means using the function
13        # compareComplexitiesWithKMeans(Data,K) from the lab session
14        Cmatrix, global_complexity = compareComplexity(X, i)
15        # store the values in corresponding lists
16        l_elbow.append(kmeans.inertia_)
17        l_kolo.append(global_complexity)
18    return l_elbow, l_kolo

```

Listing 4: Function to evaluate Global Complexity and the distortion

5.2.2 Function to plot the comparison for a given data name

```

1 def plot(max_clusters, data_name):
2     # retrieve the two lists of elbow method and global complexity to plot
3     l_elbow, l_kolo = eval_complex(max_clusters, data_name)
4     # fix the figure size
5     fig, ax1 = plt.subplots(figsize=(12, 8))
6     # use subplot to show the two plots according to their scales.
7     # plot the global complexity
8     ax1.set_xlabel("Number of clusters K")
9     ax1.set_ylabel('Global Complexity', color='red')
10    ax1.plot(range(1, max_clusters), l_kolo, color="red")
11    # plot the minimum of global complexity
12    ax1.plot(np.argmin(l_kolo)+1, l_kolo[np.argmin(l_kolo)],
13            "o", label="min complexity="+str(np.argmin(l_kolo)+1))
14    ax1.legend()
15    # plot the elbow method which is Sum of squared distances of samples to their
16    # closest cluster center.

```

```

16 ax2 = ax1.twinx()
17 ax2.set_ylabel(' Elbow Method', color='green')
18 ax2.plot(range(1, max_clusters), l_elbow, color="green")
19 # plot the title of the plot accordinf the data used
20 fig.suptitle("Data:" + data_name, size=16)
21 fig.tight_layout()
22 fig.subplots_adjust(top=0.90)

```

Listing 5: Function to plot the comparison for a given data name

5.3 Results for the four data sets

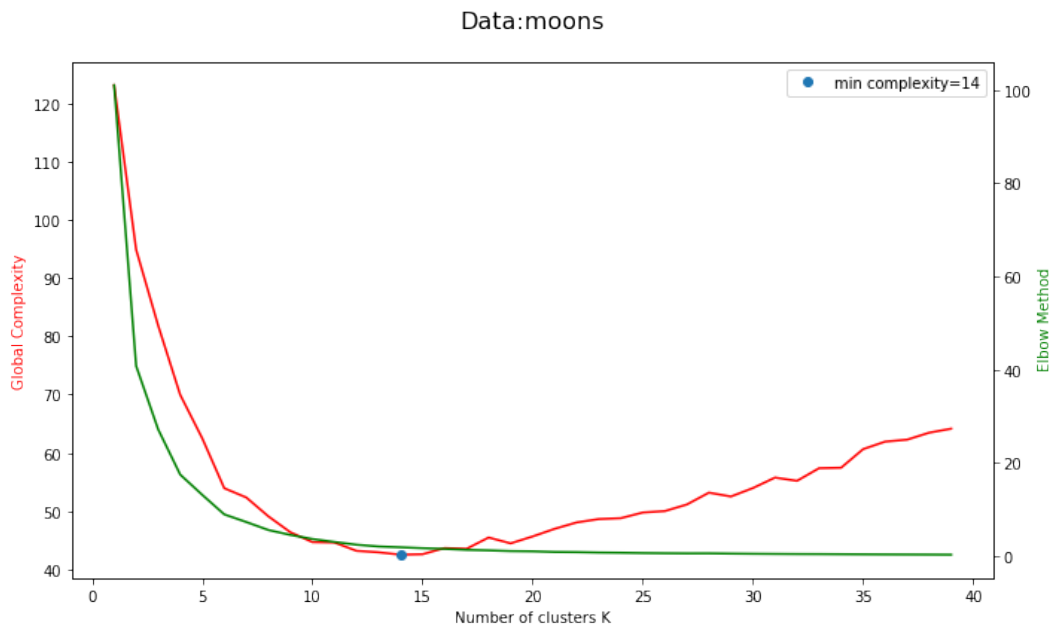


Figure 2: Elbow method vs Global Complexity (Moons data set)

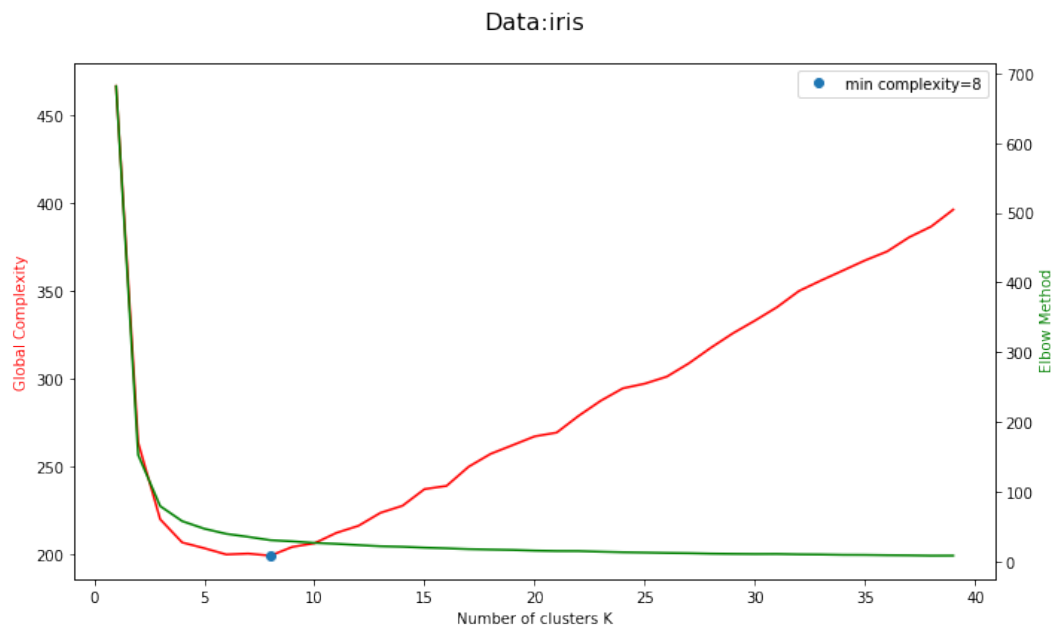


Figure 3: Elbow method vs Global Complexity (Iris data set)

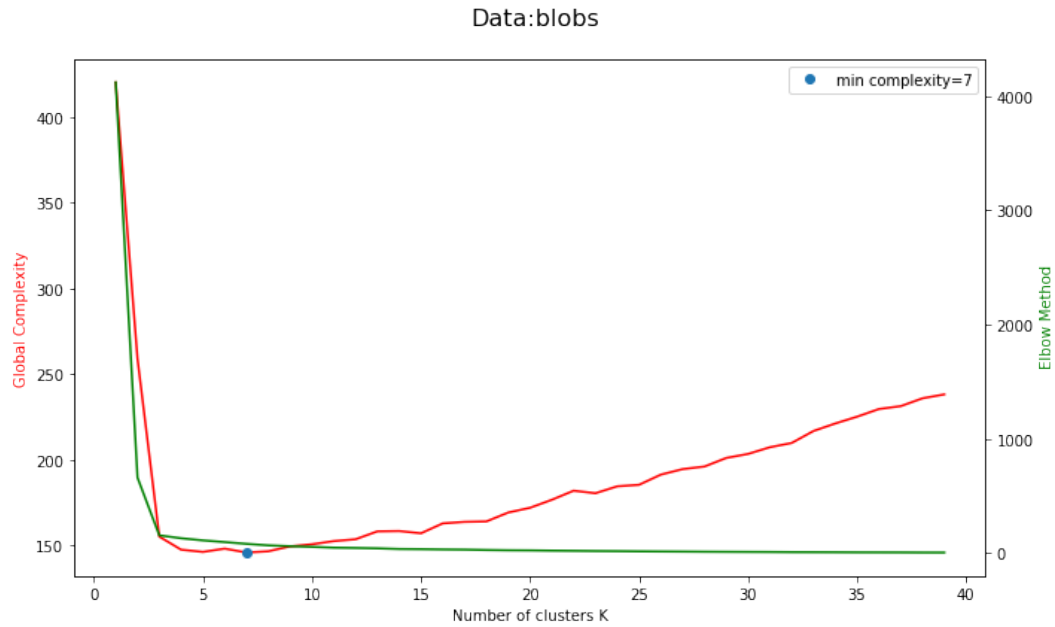


Figure 4: Elbow method vs Global Complexity (Blobs data set)

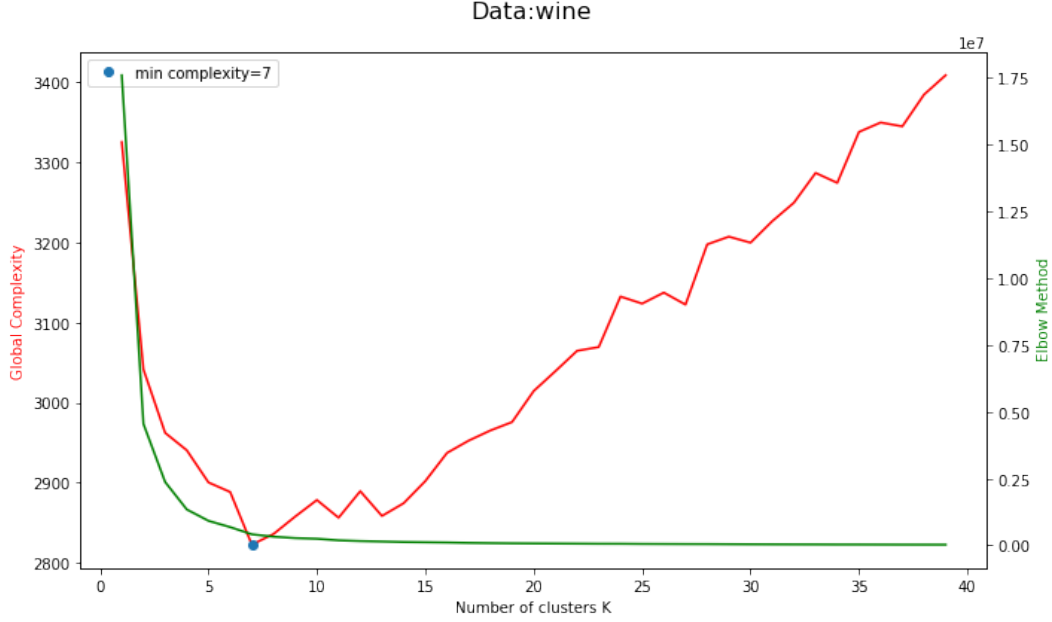


Figure 5: Elbow method vs Global Complexity (Wine data set)

5.4 Commentary

The results are certainly fascinating the Argmin of the Global Complexity is so close to what we choose manually with the elbow method. We notice also that for the four data sets the global complexity is always a strict convex function which will come very handy if we want to have a unique optimal point, a unique number of clusters.

6 Conclusion

Unsupervised learning is a very effective branch of machine learning. Nevertheless, the hyper-parameters in unsupervised learning like the number of clusters in k-means and the number of hidden layers in neural networks in supervised learning for example are one of the biggest challenges for machine learning experts. The complexity here surprisingly explained very well how the complexity decreases if we approach the optimal number of clusters. Thus, a more formal research about this topic would lead to brilliant results.

7 References

The implementation of global complexity used in this project is implemented by Mr Jean-Louis Dessalles and Mr Pierre-Alexandre Murena .You can contact them for further information.