

# Lab 1 Clustering Group 10

Liuxi Mei

Xiaochen Liu

2025-02-03

## SimpleKmeans

Apply “SimpleKMeans” to your data. In Weka euclidian distance is implemented in SimpleKmeans. You can set the number of clusters and seed of a random algorithm for generating initial cluster centers. Experiment with the algorithm as follows:

- 1) Choose a set of attributes for clustering and give a motivation. (Hint: always ignore attribute “name”. Why does the name attribute need to be ignored?)

The set is chosen is the attributes excluding attributes ‘Names’, ‘Protein’, and ‘Fat’. The reason for excluding ‘Protein’ and ‘Iron’ data is that they both have quite similar mean values across clusters, see picture below , which means that the differences among cluster are insignificant and contribute quite less to clustering results.

‘Name’ is also ignored as it is a string attribute. For clustering, Euclidean distance can not be calculated for strings and Weka will also report a mistake.

```
Number of iterations: 3
Within cluster sum of squared errors: 3.2290308976556164
Missing values globally replaced with mean/mode
```

  

```
Cluster centroids:
```

Attribute	Full Data	Cluster#			
		0	1	2	3
	(27)	(8)	(11)	(7)	(1)
Energy	207.4074	341.875	170.4545	115.7143	180
Protein	19	18.75	22.1818	13.8571	22
Fat	13.4815	28.875	8.1818	4.8571	9
Calcium	43.963	8.75	19.1818	77	367
Iron	2.3815	2.4375	2.3455	2.3571	2.5

Figure 1: Clustering result with all numeric attributes

- 2) Experiment with at least two different numbers of clusters, e.g. 2 and 5, but with the same seed value 10.

Number of iterations: 3				
Within cluster sum of squared errors: 1.2005472547937988				
Missing values globally replaced with mean/mode				
Cluster centroids:				
Attribute	Full Data (27)	Cluster#		
		0 (8)	1 (12)	2 (7)
Energy	207.4074 +/-101.2078	341.875 +/-46.3633	174.5833 +/-35.7045	110 +/-46.9929
Fat	13.4815 +/-11.257	28.875 +/-5.4625	8.75 +/-4.3719	4 +/-3.6968
Calcium	43.963 +/-78.0343	8.75 +/-0.7071	11.8333 +/-5.5895	139.2857 +/-109.5928

Figure 2: Calculated results with 3 clusters, seed set to 10

Number of iterations: 4						
Within cluster sum of squared errors: 0.46682463309541183						
Missing values globally replaced with mean/mode						
Cluster centroids:						
Attribute	Full Data (27)	Cluster#				
		0 (6)	1 (11)	2 (6)	3 (1)	4 (3)
Energy	207.4074 +/-101.2078	361.6667 +/-31.728	168.1818 +/-29.349	98.3333 +/-38.8158	180 +/-0	270 +/-27.8388
Fat	13.4815 +/-11.257	31 +/-4.1952	8 +/-3.6878	3.1667 +/-3.2506	9 +/-0	20.6667 +/-4.0415
Calcium	43.963 +/-78.0343	8.6667 +/-0.8165	12.0909 +/-5.7871	101.3333 +/-48.0985	367 +/-0	9 +/-0

Figure 3: Calculated results with 5 clusters, seed set to 10

- 3) Then try with a different seed value, i.e. different initial cluster centers. Compare the results with the previous results. Explain what the seed value controls.

With seed = 20, following results are given. Both the numbers of instances in clusters and the instances belong to clusters have changed. As a result, the variance and mean value for each cluster are different from the previous results.

When initializing this algorithm, N (corresponding to cluster number) initial ‘means’ are generated randomly within the data domain, which is controlled by the seed number. Different seeds give different initial ‘means’, resulting different results in this case.

```

Number of iterations: 4
Within cluster sum of squared errors: 1.172520467950207
Missing values globally replaced with mean/mode
Cluster centroids:

```

Attribute	Full Data (27)	Cluster#		
		0 (3)	1 (15)	2 (9)
Energy	207.4074 +/-101.2078	151.6667 +/-30.1386	144.3333 +/-49.5648	331.1111 +/-54.0704
Fat	13.4815 +/-11.257	7.6667 +/-2.3094	6.2 +/-4.3948	27.5556 +/-6.4636
Calcium	43.963 +/-78.0343	227.6667 +/-120.6703	28.3333 +/-30.647	8.7778 +/-0.6667

Figure 4: Calculated results with 3 clusters, seed set to 20

- 4) Do you think the clusters are “good” clusters? (Are all of its members “similar” to each other? Are members from different clusters dissimilar?)

Those are not ‘good’ clusters, at least not appropriate for the application of SimpleKmeans method. If they are obviously similar to each other, the clusters should include the same instances despite different seed values.

By applying ‘standardize’ filter to raw data and with seed set to 20, the regenerated results are shown as the figures below. When cluster number is chosen to be 3, cluster 0 and 1 different calcium values but very similar Energy and Fat values. And cluster 2 seems to be a more independent. For 5 clusters, cluster 0 and 4 are similar in Energy and Fat values but unsimilar calcium values. The similarity in calcium values is also noticed between cluster 0 and 1.

- 5) What does each cluster represent? Choose one of the results. Make up labels (words or phrases in English) which characterize each cluster.

Take figure 6 as an example, cluster 0 is the one with moderate energy, moderate fat and high calcium level. Cluster 1 is the one with moderate energy, moderate fat and low calcium level. Cluster 3 is the one with high energy, high fat and moderate energy level.

Number of iterations: 4  
 Within cluster sum of squared errors: 0.5497295323397319  
 Missing values globally replaced with mean/mode

Cluster centroids:

Attribute	Full Data (27)	Cluster#				
		0 (3)	1 (5)	2 (4)	3 (7)	4 (8)
Energy	207.4074 +/-101.2078	151.6667 +/-30.1386	86 +/-29.0259	228.75 +/-31.4576	352.8571 +/-37.1772	166.25 +/-21.9984
Fat	13.4815 +/-11.257	7.6667 +/-2.3094	1.6 +/-0.8944	16 +/-3.1623	30.1429 +/-4.4508	7.25 +/-2.5495
Calcium	43.963 +/-78.0343	227.6667 +/-120.6703	60 +/-36.4417	7.5 +/-1.9149	8.7143 +/-0.7559	14.125 +/-5.4625

Figure 5: Calculated results with 5 clusters, seed set to 20

Number of iterations: 4  
 Within cluster sum of squared errors: 1.1725204679502073  
 Missing values globally replaced with mean/mode

Cluster centroids:

Attribute	Full Data (27)	Cluster#		
		0 (3)	1 (15)	2 (9)
Energy	0 +/-1	-0.5508 +/-0.2978	-0.6232 +/-0.4897	1.2223 +/-0.5343
Fat	0 +/-1	-0.5165 +/-0.2052	-0.6468 +/-0.3904	1.2502 +/-0.5742
Calcium	0 +/-1	2.3541 +/-1.5464	-0.2003 +/-0.3927	-0.4509 +/-0.0085

Figure 6: Calculated results with 3 clusters, seed set to 20, with standardized data

```

Number of iterations: 4
Within cluster sum of squared errors: 0.5497295323397318
Missing values globally replaced with mean/mode

Cluster centroids:

```

		Cluster#				
Attribute	Full Data	0	1	2	3	4
	(27)	(3)	(5)	(4)	(7)	(8)
Energy	0	-0.5508	-1.1996	0.2109	1.4371	-0.4067
	+/-1	+/-0.2978	+/-0.2868	+/-0.3108	+/-0.3673	+/-0.2174
Fat	0	-0.5165	-1.0555	0.2237	1.4801	-0.5536
	+/-1	+/-0.2052	+/-0.0795	+/-0.2809	+/-0.3954	+/-0.2265
Calcium	0	2.3541	0.2055	-0.4673	-0.4517	-0.3824
	+/-1	+/-1.5464	+/-0.467	+/-0.0245	+/-0.0097	+/-0.07

Figure 7: Calculated results with 5 clusters, seed set to 20, with standardized data

## MakeDensityBasedClusters

Now with MakeDensityBasedClusters, SimpleKMeans is turned into a density-based cluster. You can set the minimum standard deviation for normal density calculation. Experiment with the algorithm as the follows:

- 1) Use the SimpleKMeans clusterer which gave the result you haven chosen in 5).
- 2) Experiment with at least two different standard deviations. Compare the results. (Hint: Increasing the standard deviation to higher values will make the differences in different runs more obvious and thus it will be easier to conclude what the parameter does)

With 0.001 and 0.5 selected for minStdDev, the same cluster centroids for 2 custers are initilized, see figure 8. However, the fitted cluster vary and cluster have 9 and 8 instances when the minStdDev is 0.001 and 0.5 respectively. The standard deviations for fitted normal distribution models are also different to fulfill the setting of minStdDev. As a result, log likelihood is -1.7 for the model with 0.001 minStdDev and 3.8 for the other.

Number of iterations: 2  
 Within cluster sum of squared errors: 1.9487869540091085  
 Missing values globally replaced with mean/mode

Cluster centroids:

Attribute	Full Data (27)	Cluster#	
		0 (8)	1 (19)
Energy	0	1.3286	-0.5594
Fat	0	1.3675	-0.5758
Calcium	0	-0.4513	0.19

Figure 8: Initial clusters with centroids with MakeDensityBasedClusters applied