

Mouhammadou Dabo

April 8, 2021

CS 1675: Intro to Machine Learning

Professor Milos Hauskrecht

Problem Assignment 9

**Problem 1.** K-means clustering

**Part a.**

**Sizes of groups:** 66, 98, 36

**Distance:** 76.4997



**Plot 1:** k-means clustering – 3 clusters

**Part b.**

**Sizes of groups:** 49, 36, 63, 52

**Distance:** 84.2883

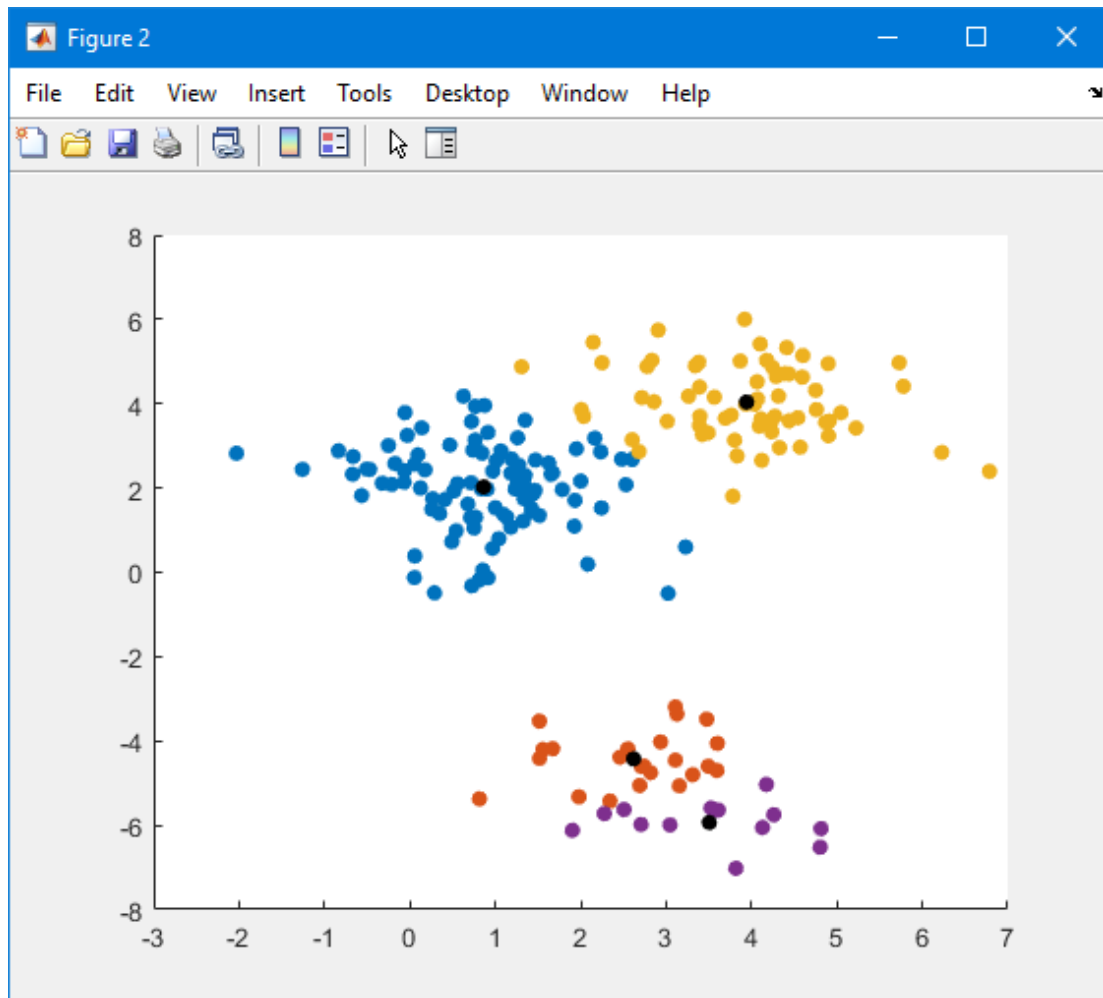


**Plot 2:** k-means clustering – 4 clusters

**Part c.**

**Sizes of groups:** 98, 23, 66, 13

**Distance:** 96.7493



**Plot 3:** k-means clustering – 4 clusters (rerun)

**Part d.**

Since the k-means clustering algorithm converges to centers minimizing the sum of squared center-point distances, the run with the smallest distance would be the best one.

A math expression that would help me compare these different clusterings and pick the best is:

$$bestCluster = \min (distances)$$

where

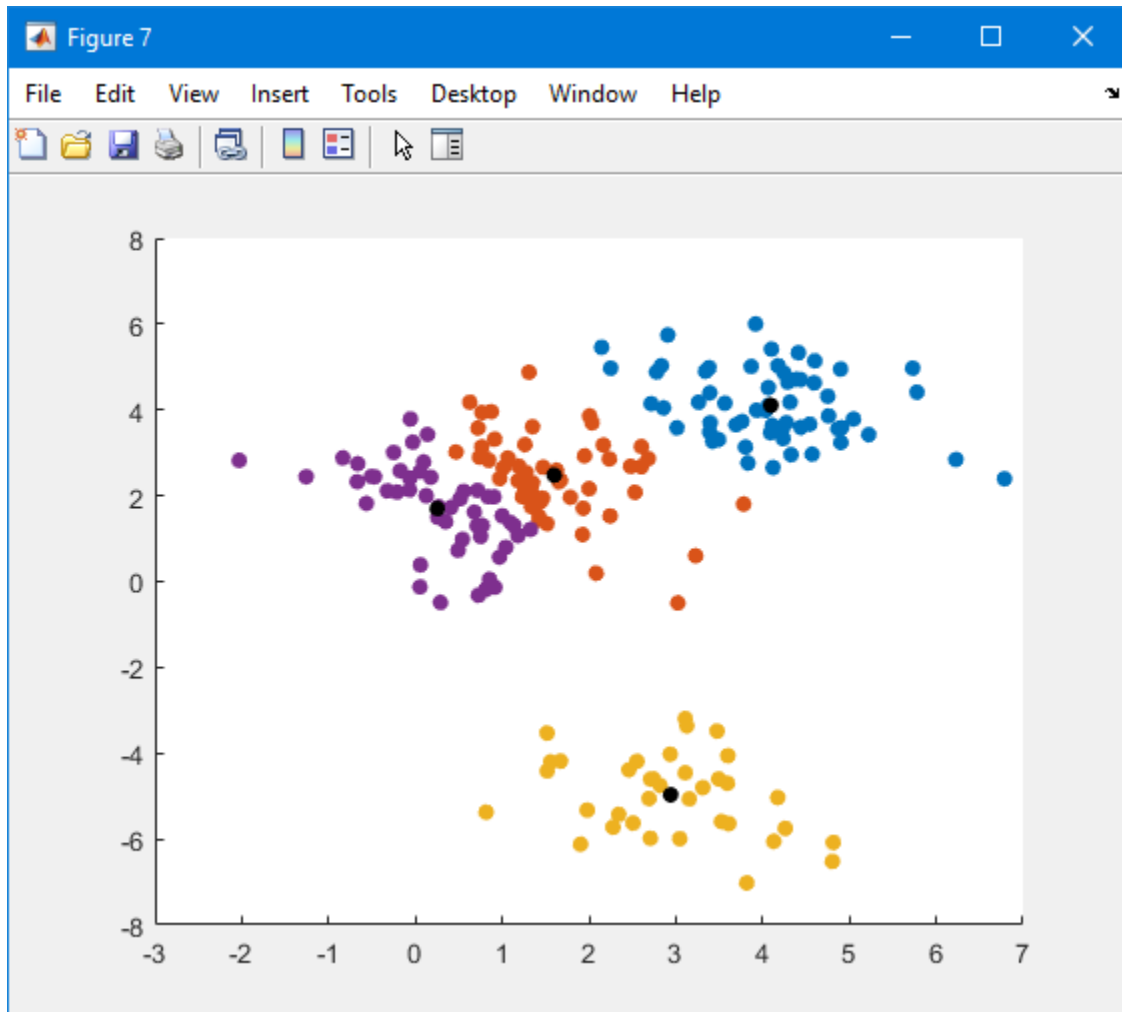
$$distance = \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - u_i\|^2$$

**Part e.**

<b>Run #</b>	<b>Cluster Sizes</b>				<b>Distance</b>
<b>1</b>	65	36	39	60	85.2187
<b>2</b>	98	10	66	26	96.5767
<b>3</b>	36	63	51	50	84.2718
<b>4</b>	39	36	65	60	85.2187
<b>5</b>	52	63	49	36	84.2883
<b>6</b>	18	66	18	98	96.0650
<b>7</b>	60	54	36	50	81.6453
<b>8</b>	96	36	40	28	83.3128
<b>9</b>	98	10	66	26	96.5767
<b>10</b>	36	63	69	32	82.2638
<b>11</b>	96	36	28	40	83.3128
<b>12</b>	23	98	66	13	96.7493
<b>13</b>	61	36	40	63	83.1815
<b>14</b>	51	36	52	61	81.7123
<b>15</b>	98	23	66	13	96.7493
<b>16</b>	63	36	49	52	84.2883
<b>17</b>	52	36	61	51	81.7123
<b>18</b>	69	63	32	36	82.2638
<b>19</b>	63	38	63	36	83.1438
<b>20</b>	38	36	98	28	83.1076
<b>21</b>	61	36	42	61	82.2850
<b>22</b>	60	36	54	50	81.6453
<b>23</b>	46	89	36	29	82.7662
<b>24</b>	31	36	92	41	82.7317
<b>25</b>	66	10	98	26	96.5767
<b>26</b>	66	10	26	98	96.5767
<b>27</b>	32	36	69	63	82.2638
<b>28</b>	36	89	46	29	82.7662
<b>29</b>	98	26	66	10	96.5767
<b>30</b>	38	36	29	97	83.2772

**Table 1:** K-means with k =4

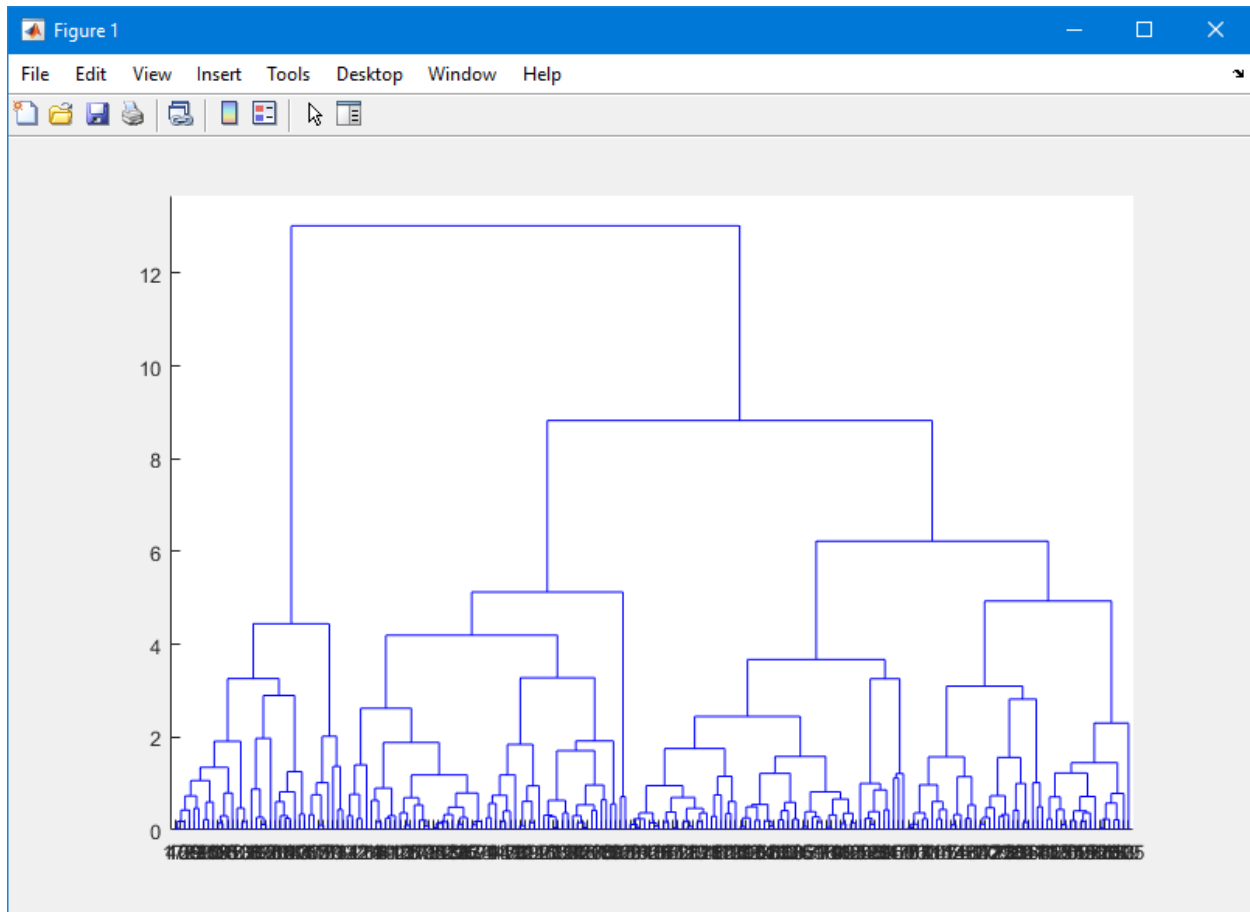
- By using the formula from Part d, the clustering that is the best is:
  - Run #7 with a cluster size of [60 54 36 50] and a distance of 81.6453



**Plot 4:** Run #7 – k-means clustering

## Problem 2. Hierarchical clustering

### Part a.



**Graph 1:** Dendrogram of Full Cluster

**Part b.**



**Plot 5:** Hierarchical Clustering Scatter Plot

- Comparing these results to the ones of Problem 2. Part e., the two clusters are different. One of the groupings match between the two (the purple from plot 5 and the yellow from plot 4).

**Problem 3.** Feature/Input ranking

**Part a.**

<b>Top</b>	<b>Dimension</b>	<b>Fisher score</b>
<b>1</b>	48	0.3192
<b>2</b>	25	0.2140
<b>3</b>	21	0.1910
<b>4</b>	70	0.1892
<b>5</b>	65	0.1693
<b>6</b>	40	0.1673
<b>7</b>	29	0.1650
<b>8</b>	19	0.1402
<b>9</b>	57	0.1255
<b>10</b>	20	0.1212
<b>11</b>	24	0.0995
<b>12</b>	30	0.0950
<b>13</b>	12	0.0858
<b>14</b>	47	0.0846
<b>15</b>	61	0.0607
<b>16</b>	10	0.0579
<b>17</b>	34	0.0527
<b>18</b>	27	0.0462
<b>19</b>	39	0.0461
<b>20</b>	41	0.0422

**Table 2:** Top 20 Fisher scores



**Part b.**

<b>Top</b>	<b>Dimension</b>	<b>AUROC score</b>
<b>1</b>	25	0.7340
<b>2</b>	29	0.6837
<b>3</b>	11	0.6695
<b>4</b>	47	0.6661
<b>5</b>	19	0.6315
<b>6</b>	34	0.6174
<b>7</b>	32	0.6021
<b>8</b>	30	0.6021
<b>9</b>	9	0.6000
<b>10</b>	56	0.5971
<b>11</b>	27	0.5953
<b>12</b>	60	0.5929
<b>13</b>	51	0.5881
<b>14</b>	26	0.5874
<b>15</b>	53	0.5845
<b>16</b>	7	0.5797
<b>17</b>	10	0.5709
<b>18</b>	61	0.5686
<b>19</b>	43	0.5567
<b>20</b>	44	0.5422

**Table 3:** Top 20 AUROC scores

- Comparing the results from part a to the ones from part b, the ordered lists are different. The two lists share only a small number of similar dimensions, but even so they are not ranked the same. Generally you would not find these two to be the same since the Fisher score is used to solve maximum likelihood, while AUROC is a performance metric for discrimination.