

## 4 Task 4

### 4.1 a)

```
1 #scale
2 cows["age_scaler"]=(cows["age"]-cows["age"].min())/(cows["age"
   ].max()-cows["age"].min())
3 cows["milk_scaler"]=(cows["milk"]-cows["milk"].min())/(cows["
   milk"].max()-cows["milk"].min())
4 cows
5
6 ndata=np.array(cows[["age_scaler","milk_scaler"]])
7
8 # Compute numerical similarity
9 def get_numsims_M(ndata):
10     lens=len(ndata)
11
12     SM=np.zeros((lens,lens))
13
14     for i in range(lens):
15         j=i+1
16         while j <lens:
17             tep=(ndata[i][0]-ndata[j][0])**2+(ndata[i][1]-ndata
18                 [j][1])**2
19             tep=np.sqrt(tep)
20             SM[i][j]=SM[j][i]=tep
21             j=j+1
22         return SM.round(4)
23 re=get_numsims_M(ndata)
24 re
```

The euclidean distances between cows is shown in Fig.4

```
array([[0.      , 0.2857, 0.52   , 0.3627, 0.9371, 1.2289],
       [0.2857, 0.      , 0.52   , 0.5429, 1.0976, 1.4142],
       [0.52   , 0.52   , 0.      , 0.3308, 0.6615, 0.9923],
       [0.3627, 0.5429, 0.3308, 0.      , 0.5759, 0.8781],
       [0.9371, 1.0976, 0.6615, 0.5759, 0.      , 0.3308],
       [1.2289, 1.4142, 0.9923, 0.8781, 0.3308, 0.      ]])
```

Figure 4: Euclidean distance

### 4.2 b)

The Goodall distances between cows is shown in Fig.5

```
array([[0.      , 0.7037, 0.7037, 1.      , 1.      , 1.      ],
       [0.7037, 0.      , 1.      , 0.75   , 1.      , 0.4537],
       [0.7037, 1.      , 0.      , 0.75   , 0.75   , 0.7037],
       [1.      , 0.75   , 0.75   , 0.      , 0.4537, 0.75   ],
       [1.      , 1.      , 0.75   , 0.4537, 0.      , 1.      ],
       [1.      , 0.4537, 0.7037, 0.75   , 1.      , 0.      ]])
```

Figure 5: Goodall distances

```
1  #construct dict
2  dic={"Ayrshir": 1/2,"Holstein": 1/3,"Finnccattle": 1/6,
3       "lively":1/6,"kind":1/3,"calm":1/2,
4       "rock":1/3,"country":1/3,"classical":1/3}
5
6  #Get categorical features
7  ccols=["race","character","music"]
8  cdata=np.array(cows[ccols])
9  cdata
10
11 #Compute categorical similarity
12 def get_catsim_M(cdata,dic):
13     lens,f_num=cdata.shape
14
15     CM=np.zeros((lens,lens))
16
17     for i in range(lens):
18         j=i+1
19         while j <lens:
20             tep=0
21             for x in range(f_num):
22                 if cdata[i][x]==cdata[j][x]:
23                     tep=tep+1-dic[cdata[i][x]]**2
24             CM[i][j]=CM[j][i]=1-tep/f_num
25             j=j+1
26     return CM.round(4)
27
28 re=get_catsim_M(cdata,dic)
29
30 re
```

### 4.3 c)

Distance using both numerical and categorical features is shown in Fig.6

```
array([[0.    , 1.507, 1.735, 2.1   , 2.658, 2.941],
       [1.507, 0.    , 2.253, 1.838, 2.814, 2.167],
       [1.735, 2.253, 0.    , 1.632, 1.953, 2.194],
       [2.1   , 1.838, 1.632, 0.    , 1.352, 2.164],
       [2.658, 2.814, 1.953, 1.352, 0.    , 2.069],
       [2.941, 2.167, 2.194, 2.164, 2.069, 0.    ]])
```

Figure 6: Combined distance measure

```
1 numsim=get_numsims_M(ndata)
2 catsim=get_catsim_M(cdata,dic)
3
4 a=2/5
5
6 mixSim=a*numsim/np.std(numsim)+(1-a)*catsim/np.std(catsim)
7 mixSim=mixSim.round(3)
8 mixSim
```

## 4.4 d)

### 4.4.1 Histograms

```
1 #Show Histograms
2 def show_hist(data,bins):
3     x = []
4     for i in range(6):
5         x.extend(data[i][i+1:6])
6     # print(x)
7     plt.hist(x, bins=bins)
8     plt.gca().set(title='Frequency Histogram', ylabel='
        Frequency')
9     plt.show()
```

After test different numbers of bins in all three cases, I choose 8,6,8 in each plot as their number of bins to illustrate the clusters of cows in data. The figures are shown in Fig.7 I think the measurement combined numerical and categorical features can better cluster cows. Because more features are used for calculation, it means that we have considered more comprehensively, resulting better effect.

### 4.4.2 Graph-based clustering

First we can create a complete distance graph between cows. and we know the distances of cows from the task above. Then We should remove edges with longest distances(largest

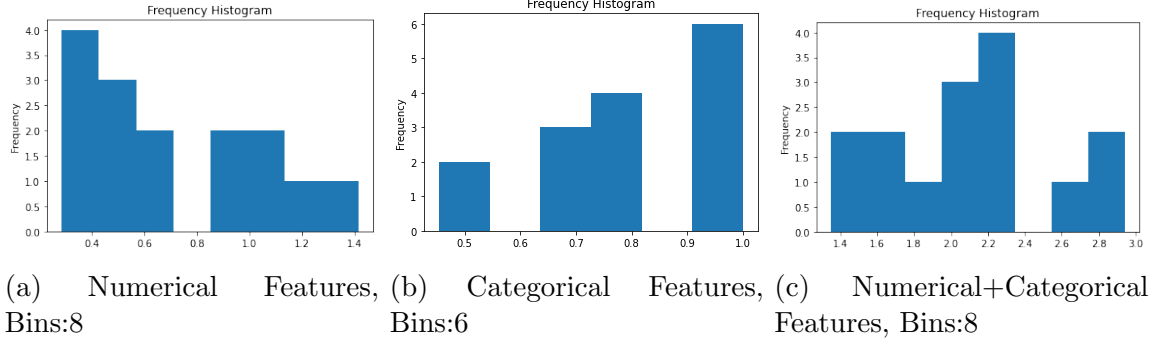


Figure 7: Histograms in 3 different cases

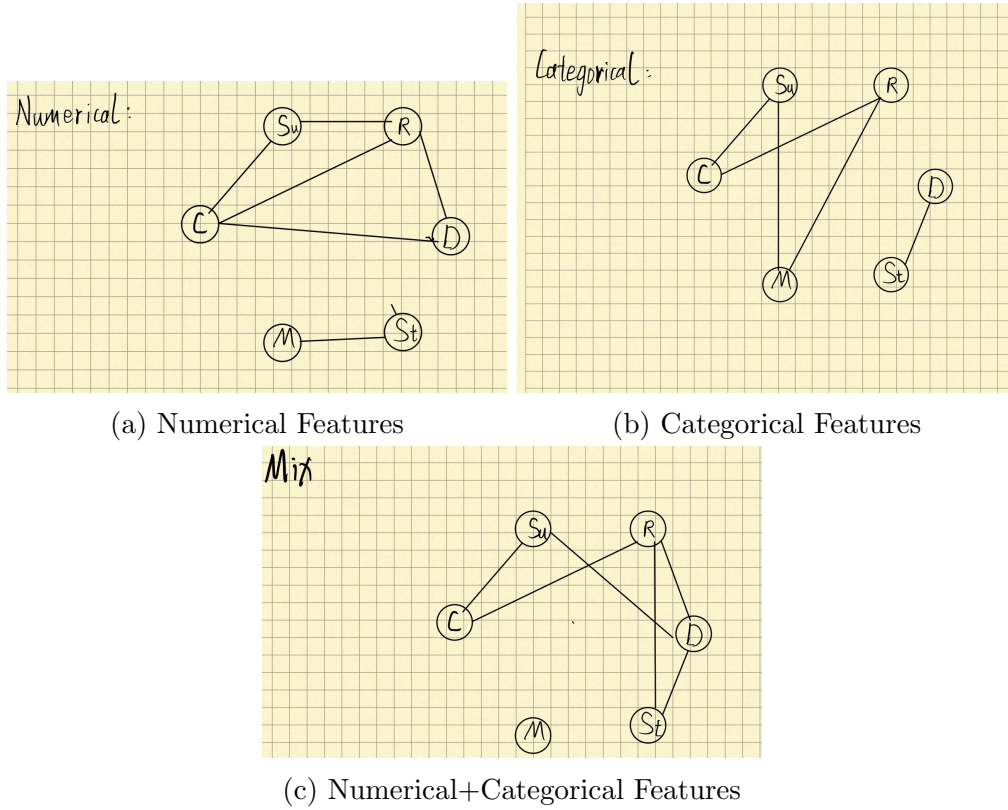


Figure 8: Graph clustering in 3 different cases

weights) until the graph is broken into two connected components. These components are clusters of cows. I think the measurement of mix features including numerical and categorical can better cluster cows from Fig.8. The results are not so radical to make the number of cows out of balance among the different clusters. It balances the application of two different types of data into the similarity, the more data types are used in the results, making the results more inclusive.