## **Question 1**

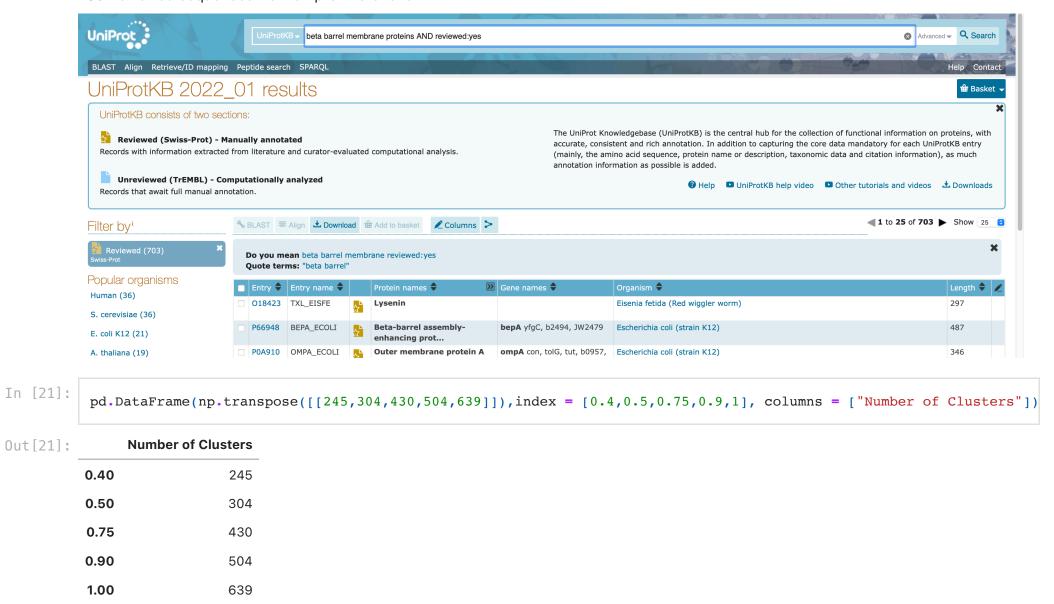
```
In [16]:
          import string
          import pandas as pd
          import numpy as np
          aas = set(string.ascii_uppercase) - set("BJOZXU")
          a = "AMENLNMDLLYMAAAVMMGLAAIGAAIGIGILGGKFLEGAARQPDLIPLLRTQFFIVMGLVDAIPMIAVG LGLYVMFAVA".replace(" ", '')
          b = "AADVSAAVGATGQSGMTYRLGLSWDWDKSWWQTSTGRLTGYWDAGYTYWEGGDEGAGKHSLSFAP VFVYEFAGDSIKPFIEAGIGVAAFSGTRVGDQNLGSSLNFEDR
          c = "MALLPAAPGAPARATPTRWPVGCFNRPWTKWSYDEALDGIKAAGYAWTGLLTASKPSLHHATATPEY LAALKQKSRHAA".replace(" ", '')
          def comp(a):
              return {x : a.count(x)/len(a) for x in set(aas)}
          def hamm(a,b):
              d1 = comp(a)
              d2 = comp(b)
              return sum([abs(d1[x]-d2[x]) for x in aas])
          def euclidean(a,b):
              d1 = comp(a)
              d2 = comp(b)
              return sum([(d1[x] - d2[x])**2 for x in aas])**0.5
          eu = [euclidean(x,y)  for x,y  in [(a,b),(b,c),(c,a)]]
          ha = [hamm(x,y) for x,y in [(a,b),(b,c),(c,a)]]
          df = pd.DataFrame([eu,ha], index = ['Euclidean', "Hamming"], columns = ['1,2','2,3','1,3'])
          df
Out[16]:
                       1,2
                                2,3
                                         1,3
```

## Question 2

703 Reviewed sequences from uniprot were taken

Euclidean 0.201062 0.201130 0.220868

Hamming 0.665728 0.726633 0.843354



Header of each file in CD-HIT Cluster is attached

```
Download all files
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                810aa, >sp|P0A940|BAMA_ECOLI... at 78.77%
     Browse clusters by size
Browse clusters by length
Distribution of clusters
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 810aa, >sp | P0A940 | BAMA_ECOLI ... at 78.77% 803aa, >sp | B5F878 | BAMA_SALA4 ... at 79.83% 810aa, >sp | B5F878 | BAMA_SALA4 ... at 79.83% 810aa, >sp | B5E076 | BAMA_ECOSE ... at 78.77% 810aa, >sp | B7MP17 | BAMA_ECOSE ... at 78.77% 810aa, >sp | B7MP17 | BAMA_ECOSI ... at 78.64% 810aa, >sp | B7MP17 | BAMA_ECOSI ... at 78.77% 809aa, >sp | P0A942 | BAMA_ECOSI ... at 78.77% 809aa, >sp | B5Y114 | BAMA_ECOSI ... at 78.78% 809aa, >sp | B5Y114 | BAMA_ELEP3 ... at 78.86% 809aa, >sp | G5T439 | BAMA_ELEP3 ... at 78.86% 809aa, >sp | G5T439 | BAMA_ELEP3 ... at 78.86% 810aa, >sp | B7M170 | BAMA_ECOSI ... at 78.77% 810aa, >sp | B7M170 | BAMA_ECOSS ... at 78.77% 810aa, >sp | B7M170 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C5CASI | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSA ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | B1LGY4 | BAMA_ECOSM ... at 78.77% 810aa, >sp | B1LGY4 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ECOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_ESCOSM ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SALTI ... at 79.10% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BAMA_SHIFE ... at 78.77% 810aa, >sp | C3EA78 | BA
Raw output
Download all files
                                                                                                                                                                                                                                                                                                                                                                                                                                                       >Cluster 0
Browse clusters by size
Browse clusters by length
Distribution of clusters
Raw output
Download all files
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 BY 0

810aa, >sp | P0A940 | BAMA_ECOLI ... at 78.778
803aa, >sp | B5F878 | BAMA_SALA4 ... at 79.838
810aa, >sp | B5F878 | BAMA_SALA4 ... at 79.838
810aa, >sp | B5F878 | BAMA_SALA4 ... at 79.018
810aa, >sp | B5E976 | BAMA_ECOSE ... at 78.778
810aa, >sp | B7M977 | BAMA_ECOS1 ... at 78.778
810aa, >sp | B7M977 | BAMA_ECO51 ... at 78.778
809aa, >sp | B5Y174 | BAMA_ECO45 ... at 78.778
809aa, >sp | B5Y174 | BAMA_ECO57 ... at 78.788
809aa, >sp | B5Y174 | BAMA_ECO57 ... at 78.788
814aa, >sp | G5F458 | BAMA_ECO57 ... at 78.888
810aa, >sp | B5T447 | BAMA_ECO55 ... at 78.778
810aa, >sp | B7M170 | BAMA_ECO55 ... at 78.778
810aa, >sp | B7M170 | BAMA_ECO55 ... at 78.778
810aa, >sp | B7M170 | BAMA_ECO55 ... at 78.778
809aa, >sp | C5CAT3 | BAMA_ECO55 ... at 78.778
809aa, >sp | C5CAT3 | BAMA_ECO55 ... at 78.778
810aa, >sp | B7M170 | BAMA_ECO55 ... at 78.778
810aa, >sp | S7M170 | BAMA_ECO55 ... at 78.778
810aa, >sp | A72H87 | BAMA_ECO55 ... at 78.778
810aa, >sp | C4SR78 | BAMA_ECO55 ... at 78.778
810aa, >sp | C4SR78 | BAMA_ECO55 ... at 78.778
810aa, >sp | S1M246 | BAMA_ECO55 ... at 78.778
810aa, >sp | S1M246 | BAMA_ECO55 ... at 78.778
810aa, >sp | S1M246 | BAMA_ECO55 ... at 78.778
810aa, >sp | S1M246 | BAMA_ECO55 ... at 78.778
810aa, >sp | S1M246 | BAMA_ECO55 ... at 78.778
810aa, >sp | S1M74 | BAMA_ECO55 ... at 78.778
810aa, >sp | S2M21 | BAMA_ECO55 ... at 79.554
805aa, >sp | C95R5 | BAMA_SHIFE ... at 78.778
810aa, >sp | Q2NR15 | BAMA_SHIFE ... at 78.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 78.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 78.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 78.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 78.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 78.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.778
810aa, >sp | S2M15 | BAMA_SHIFE ... at 8.
                                                                                                                                                                                                                                                                                                                                                                                                                                                       >Cluster 0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                810aa, >sp|P0A940|BAMA_ECOLI... at 78.77%
  Browse clusters by size
  Distribution of clusters
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             795aa, >sp | C5BTR5 | BAMA_EDWT9... at 78.49%
805aa, >sp | A4W6S2 | BAMA_ENT38... at 80.12%
810aa, >sp | A5V06 | BAMA_SALNS... at 79.38%
803aa, >sp | Q5PD65 | BAMA_SALPA... at 79.38%
  Raw output
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          810aa, >sp | B54878 | BAMA_SALA4... at 93.52% | B5F878 | BAMA_SALA4... at 93.52% | B5F878 | BAMA_SALA4... at 95.93% | 810aa, >sp | B5E076 | BAMA_SCO5E... at 100.00% | 810aa, >sp | B7M57 | BAMA_SCO5E... at 100.00% | 810aa, >sp | B7M57 | BAMA_ECO5T... at 100.00% | 810aa, >sp | B7M57 | BAMA_ECO5T... at 100.00% | 809aa, >sp | B5Y1194 | BAMA_ECO5T... at 100.00% | 809aa, >sp | B5Y1194 | BAMA_KLEF7... at 91.57% | 805aa, >sp | Q57731 | BAMA_SALCH... at 93.54% | 810aa, >sp | B7LEN8 | BAMA_ECO5T... at 93.54% | 810aa, >sp | B7LEN8 | BAMA_ECO5T... at 100.00% | 8008aa, >sp | B7LEN8 | BAMA_ECO5T... at 100.00% | 8008aa, >sp | B7LEN8 | BAMA_ECO5T... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8055a... at 100.00% | 8008aa, >sp | B7LEN8 | 8056aa, >sp | 8056aa, >sp
  Browse clusters by size
Browse clusters by length
Distribution of clusters
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              805aa, >sp | 057731 | BAMA_SALCH... at 93.54% 810aa, >sp | BTMGN8 | BAMA_ECO55... at 100.00% 810aa, >sp | B7MIY0 | BAMA_ECO8A... at 100.00% 810aa, >sp | A7ZHR7 | BAMA_ECO8A... at 100.00% 810aa, >sp | C4ZRR9 | BAMA_ECO8M... at 100.00% 810aa, >sp | B1XD46 | BAMA_ECOBM... at 100.00% 810aa, >sp | B1XD46 | BAMA_ECOBM... at 100.00% 810aa, >sp | B1LGX9 | BAMA_ECOBM... at 98.89% 810aa, >sp | C4ZRP4 | BAMA_ECOBM... at 98.89% 810aa, >sp | S7LW74 | BAMA_ECOTT... at 100.00% 802aa, >sp | B7LW74 | BAMA_ECOTT... at 93.66% 803aa, >sp | Q8Z9A3 | BAMA_SALPC... at 93.76% 803aa, >sp | Q8Z9A3 | BAMA_SALPC... at 93.66% 810aa, >sp | Q8Z9A3 | BAMA_SALPC... at 100.00% 810aa, >sp | P0A943 | BAMA_SHIPS... at 100.00% 810aa, >sp | P0A943 | BAMA_SHIPS... at 100.00% 810aa, >sp | B7UIM2 | BAMA_ECOZT... at 100.00% 810aa, >sp | B7UIM2 | BAMA_SLDS... at 95.31% 803aa, >sp | S2DS6 | BAMA_SALPS... at 95.31% 810aa, >sp | B7UIM2 | BAMA_SLDS... at 95.31% 810aa, >sp | B7UIM2 | BAMA_SLDS... at 95.31% 810aa, >sp | B7UIM2 | BAMA_SLDS... at 95.31% 810aa, >sp | A7ZWC3 | BAMA_SHIBS... at 100.00% 810aa, >sp | A7ZWC3 | BAMA_SHIBS... at 100.00% 810aa, >sp | A7ZWC3 | BAMA_SCOLS... at 100.00% 810aa, >sp | A7ZWC3 | BAMA_ECOMS... at 99.88% 810aa, >sp | P0A941 | BAMA_ECOLS... at 100.00% 810aa, >sp | B1GG4 | BAMA_ECOLS... at 100.00% 810aa, >sp | B6HEF1 | BAMA_ECOLS... at 100.00% 810aa, >sp | B5F1GG4 | BAMA_ECOLS... at 100.00% 804aa, >sp | B5F1GG4 | BAMA_ECOLS... at 93.78% 804aa, >sp | B5F1GG4 | BAMA_SALDC... at 93.78% 804aa, >sp | B5F1GG4 | BAMA_SALDC.... at 93.78% 804aa, >sp | B5F1GG4 | BAMA_SALDC.... 
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   810aa, >sp B7LGN8 BAMA_ECO55... at 100.009
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                810aa, >sp B5RHG2 BAMA_SALG2... at 95.31%
```

Raw output

Since Pisces server is down, these questions could not be done. Instead the comaparision is done between uniport and CD-HIT at cutoffs in question 5

## Question 5

1.00

639

640