

siamese+cat	52.90	3.33	47.54	17.22	83.55
skunk	12.50	2.50	49.85	18.19	8.89
mole	16.67	0.00	55.31	29.79	13.75
hippopotamus	36.59	23.96	39.36	24.34	5.64
leopard	8.75	20.60	52.14	2.31	0.00
mongoose	37.16	11.25	33.84	22.08	6.25
spider+monkey	55.00	39.06	9.38	31.67	12.78
humpback+whale	60.31	51.03	19.06	20.34	0.00
elephant	22.81	49.87	6.25	12.29	6.88
gorilla	56.19	54.50	7.50	44.17	18.93
ox	14.68	16.79	32.39	13.12	53.33
fox	79.26	18.75	48.18	27.88	4.38
sheep	10.62	73.69	0.00	7.50	30.68
seal	44.49	49.84	6.25	18.18	47.03
chimpanzee	84.36	74.51	8.75	41.93	36.62
hamster	18.37	16.98	35.00	36.60	71.44
squirrel	18.48	9.17	22.58	21.56	14.76
rhinoceros	10.59	22.85	25.17	3.75	2.50
rabbit	11.47	38.16	18.89	38.89	60.52
bat	32.30	62.25	19.97	34.91	5.56
giraffe	16.95	40.89	9.57	8.69	3.33
wolf	61.48	46.81	41.97	18.40	5.00
chihuahua	29.84	1.25	44.91	9.80	73.55
rat	24.99	18.94	40.87	23.12	16.69
weasel	36.56	13.06	26.26	10.14	3.75
otter	27.22	17.92	30.56	23.19	7.78
buffalo	10.00	52.99	8.80	9.38	2.31
sebra	22.63	80.60	1.25	19.09	8.94
giant+panda	37.30	28.10	55.38	38.19	29.58
deer	40.05	55.48	6.51	31.55	10.27
bobcat	42.89	13.47	59.67	20.29	3.75
pig	23.04	28.07	4.19	17.88	46.95
lion	25.70	56.17	11.88	43.12	3.75
mouse	12.42	23.96	6.88	23.26	27.64
polar+bear	16.25	13.75	48.75	19.38	5.00
collie	49.17	0.82	45.98	18.57	79.11
walrus	23.96	60.41	11.19	33.77	18.75
raccoon	48.68	13.01	35.95	28.26	5.00
cow	13.97	51.57	3.04	18.89	72.89
dolphin	60.38	49.62	3.96	14.05	37.98

[50 rows x 85 columns]

```
In [36]: # run k-means on the data, with k = 10.
kmeans = KMeans(n_clusters = 10, init = 'k-means++', n_init = 10)
kmeans.fit(predicate_matrix, classes_with_animal_names)
```

Out[36]: RMeans(n_clusters=10)

```
In [37]: # label all animals
print('Labels for animals:\n',kmeans.labels_)

Labels for animals:
[3 7 2 9 5 3 9 2 5 4 4 8 1 8 3 6 2 1 6 1 0 3 2 6 4 4 1 4 0 3 8 5 0 0 9 3
 3 5 3 8 1 8 4 7 5 2 0 3 2]
```

```
In [38]: # cluster animals to their labels

clustering = {}
for i in range(0,10):
    clustering[i] = []
    for i,j in enumerate(classes_with_animal_names):
        clustering[kmeans.labels_[i]].append(j)

for i in range(0,10):
    print('label =',i,'\n')
    print('animals:',clustering[i],'\n')

label = 0

animals: ['fox', 'bat', 'rat', 'weasel', 'raccoon']

label = 1

animals: ['hippopotamus', 'elephant', 'ox', 'rhinoceros', 'pig']

label = 2

animals: ['killer-whale', 'blue-whale', 'humpback-whale', 'seal', 'walrus', 'dolphin']

label = 3

animals: ['antelope', 'horse', 'moose', 'sheep', 'giraffe', 'buffalo', 'zebra', 'deer', 'cow']

label = 4

animals: ['skunk', 'mole', 'hamster', 'squirrel', 'rabbit', 'mouse']

label = 5

animals: ['dalmatian', 'persian+cat', 'siamese+cat', 'chihuahua', 'giant+panda', 'collie']

label = 6

animals: ['spider+monkey', 'gorilla', 'chimpanzee']

label = 7

animals: ['grizzly+bear', 'polar+bear']

label = 8

animals: ['german+shepherd', 'tiger', 'leopard', 'wolf', 'bobcat', 'lion']

label = 9

animals: ['beaver', 'otter']
```

Conclusion for (a):

Does the clustering make sense?

--- Yes, it makes some sense.

--- From the clustering above, we can see that pets in label 7 are grouped close together as well as other labels.

--- However, it is not perfect. we need to figure out optimal k or find other algorithms to cluster these animals better.

(b) Now hierarchically cluster this data, using scipy.cluster.hierarchy.linkage. Choose Ward's method, and plot the resulting tree using the dendrogram method, setting the orientation parameter to 'right' and labeling each leaf with the corresponding animal name. You will run into a problem: the plot is too cramped because the default figure size is so small. To make it larger, preface your code with the following:

```
from pylab import rcParams
rcParams['figure.figsize'] = 5, 10
```

(or try a different size if this doesn't seem quite right). Does the hierarchical clustering seem sensible to you?

```
In [40]: # use scipy.cluster.hierarchy.linkage to hierarchically cluster this data.
hierarchically_cluster_link = linkage(predicate_matrix, method = 'ward')

# use dendrogram method to plot the resulting tree.
dendrogram(hierarchically_cluster_link, orientation = 'right', labels = classes_with_animal_names)

# resize the figure
rcParams['figure.figsize'] = [40,40]
```



Conclusion for (b):

Does the hierarchical clustering seem sensible to you?

--- Yes, it seems sensible to me.

--- For choosing the linkage, Ward's method is the sensible default. It groups based on reducing the sum of squared distances of each observation from the average observation in a cluster.

(c) Turn in an iPython notebook with a transcript of all this experimentation.

See above for all this experimentation