

```
1  """
2  Clusters Bob Ross paintings by features.
3
4  By Walter Hickey <walter.hickey@fivethirtyeight.com>
5
6  See http://fivethirtyeight.com/features/a-statistical-analysis-of-the-work-of-bob-ross/
7  """
8
9  import numpy as np
10 from scipy.cluster.vq import vq, kmeans, whiten
11 import math
12 import csv
13
14 def main():
15
16     # load data into vectors of 1s and 0s for each tag
17     with open('elements-by-episode.csv', 'r') as csvfile:
18         reader = csv.reader(csvfile)
19         reader.next() # skip header
20         data = []
21         for row in reader:
22             data.append(map(lambda x: int(x), row[2:])) # exclude EPISODE and TITLE
23 columns
24
25     # convert to numpy matrix
26     matrix = np.array(data)
27
28     # remove columns that have been tagged less than 5 times
29     columns_to_remove = []
30     for col in range(np.shape(matrix)[1]):
31         if sum(matrix[:,col]) <= 5:
32             columns_to_remove.append(col)
33     matrix = np.delete(matrix, columns_to_remove, axis=1)
34
35     # normalize according to stddev
36     whitened = whiten(matrix)
37     output = kmeans(whitened, 10)
38
39     print "episode", "distance", "cluster"
40
41     # determine distance between each of 403 vectors and each centroid, find closest neighbor
42     for i, v in enumerate(whitened):
43
44         # distance between centroid 0 and feature vector
45         distance = math.sqrt(sum((v - output[0][0]) ** 2))
46
47         # group is the centroid it is closest to so far, set initially to centroid 0
48         group = 0
49         closest_match = (distance, group)
50
51         # test the vector i against the 10 centroids, find nearest neighbor
52         for x in range(0, 10):
53             dist_x = math.sqrt(sum((v - output[0][x]) ** 2))
54             if dist_x < closest_match[0]:
55                 closest_match = (dist_x, x)
56
57         print i+1, closest_match[0], closest_match[1]
58
59 if __name__ == "__main__":
60     main()
```

