

CSE 569 Fundamentals of Statistical Learning HW4 Solutions

Akshay Shah, 1215153316

Fall 2019

1 Solution

1.1 Code

```
1
2 # coding: utf-8
3 import math
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from sklearn.cluster import KMeans
7
8 def get_x(fname):
9     f = open(fname)
10    return np.array([int(x) for x in f.readlines()])
11
12 def initialize_params():
13     theta = np.array([np.random.uniform(0,1) for i in range(3)])
14     pi = np.array([np.random.uniform(0,1) for i in range(3)])
15     return theta, pi
16
17 def E_step(data,pi,theta,no_of_clusters,no_of_tosses,no_of_experiments):
18     si_mat = np.zeros((no_of_experiments,no_of_clusters))
19     for i in range(no_of_experiments):
20         factorial_cmp = math.factorial(no_of_tosses) // (math.factorial(no_of_tosses -
21         no_of_clusters) * math.factorial(no_of_clusters))
22         X = pi * factorial_cmp * (theta**data[i][0] * (1-theta)**(no_of_tosses - data[i][0]))
23         si_mat[i] = X / np.sum(X)
24     return si_mat
25
26 def M_step(si_mat,data,no_of_experiments):
27     n_k = np.sum(si_mat,axis=0)
28     theta = (1/(20 * n_k)) * np.sum(si_mat * data,axis=0)
29     pi = n_k/no_of_experiments
30     return pi,theta
31
32 def log_likelihood(data,pi,theta,no_of_clusters,no_of_tosses,no_of_experiments):
33     ll = 0
34     for i in range(no_of_experiments):
35         factorial_cmp = math.factorial(no_of_tosses) // (math.factorial(no_of_tosses - data[i]
36         ][0]) * math.factorial(data[i][0]))
37         theta_cmp = theta**data[i][0] * ((1-theta)**(no_of_tosses - data[i][0]))
38         ll += (np.log(np.sum(factorial_cmp * theta_cmp * pi,axis=0)))
39     return ll
40
41 #Plotting neg log likelihood
```

```

40 def plot_neg_likelihood_vs_iterations(iterations,likelihood_list):
41     plt.plot(list(range(iterations-1)),likelihood_list)
42     plt.title("No. of Iterations vs Negative Log Likelihood ")
43     plt.xlabel("Iterations")
44     plt.ylabel("Negative Log Likelihood")
45     plt.show()
46
47 def EM(x,pi,theta,no_of_clusters,no_of_experiments,no_of_tosses):
48     iterations = 0
49     ll = log_likelihood(x,pi,theta,no_of_clusters,no_of_tosses,no_of_experiments)
50     lls = [ll]
51     while 1:
52         si = E_step(x,pi,theta,no_of_clusters,no_of_tosses,no_of_experiments)
53         pi,theta = M_step(si,x,no_of_experiments)
54         ll = log_likelihood(x,pi,theta,no_of_clusters,no_of_tosses,no_of_experiments)
55         iterations+=1
56         if abs(lls[-1] - ll) < 0.00001:
57             break
58         lls.append(ll)
59
60     lls = [-x for x in lls[1:]]
61     return iterations,pi,theta,lls
62
63 def print_params(pi,theta):
64     for i in range(3):
65         print("Prior probability for coin%s:%s"%(i+1,pi[i]))
66         print("Probability of head for coin%s:%s"%(i+1,theta[i]))
67
68 no_of_clusters = 3
69 no_of_tosses = 20
70 no_of_experiments = 1000
71
72 x = get_x("Binomial_20_flips.txt").reshape(-1,1)
73
74 #Kmeans experiment
75 kmeans = KMeans(n_clusters=3, random_state=0).fit(x)
76 cluster_assignments = kmeans.labels_
77
78 coin1_points = cluster_assignments[cluster_assignments == 0]
79 coin2_points = cluster_assignments[cluster_assignments == 1]
80 coin3_points = cluster_assignments[cluster_assignments == 2]
81
82 coin1_pi = coin1_points.shape[0] / no_of_experiments
83 coin2_pi = coin2_points.shape[0] / no_of_experiments
84 coin3_pi = coin3_points.shape[0] / no_of_experiments
85
86 theta1 = kmeans.cluster_centers_[0][0] / no_of_tosses
87 theta2 = kmeans.cluster_centers_[1][0] / no_of_tosses
88 theta3 = kmeans.cluster_centers_[2][0] / no_of_tosses
89
90 pi_kmeans = [coin1_pi,coin2_pi,coin3_pi]
91 theta_kmeans = [theta1,theta2,theta3]
92
93 print("Estimated Parameters by Kmeans clustering algorithm")
94 print_params(pi_kmeans,theta_kmeans)
95
96 # pi = np.array([0.333, 0.333, 0.333])
97 # theta = np.array([0.5, 0.6, 0.4])
98
99 #Non Kmeans initialization

```

```

100 pi,theta = initialize_params()
101 iterations,pi,theta,likelihood_list = EM(x,pi,theta,no_of_clusters,no_of_experiments,no_of_tosses
    )
102 print("#####")
103 print("Estimated Parameters by EM with random initialization")
104 print_params(pi,theta)
105 plot_neg_likelihood_vs_iterations(iterations,likelihood_list)
106
107 #Kmeans initialization
108 pi,theta = np.array(pi_kmeans),np.array(theta_kmeans)
109 iterations,pi,theta,likelihood_list = EM(x,pi,theta,no_of_clusters,no_of_experiments,no_of_tosses
    )
110 print("#####")
111 print("Estimated Parameters by EM with Kmeans initialization")
112 print_params(pi,theta)
113
114 plot_neg_likelihood_vs_iterations(iterations,likelihood_list)

```

Listing 1: Kmeans and Binomial Mixture Model

1.2 Outputs

Estimated Parameters by Kmeans clustering algorithm

Prior probability for coin1:**0.341**

Probability of head for coin1:**0.6653958944281524**

Prior probability for coin2:**0.212**

Probability of head for coin2:**0.2646226415094358**

Prior probability for coin3:**0.447**

Probability of head for coin3:**0.8731543624161094**

Estimated Parameters by EM with random initialization

Prior probability for coin1:**0.3146525869543541**

Probability of head for coin1:**0.6656711767117973**

Prior probability for coin2:**0.4827779465602027**

Probability of head for coin2:**0.852239260664435**

Prior probability for coin3:**0.2025694664854404**

Probability of head for coin3:**0.2586891442935458**

Estimated Parameters by EM with Kmeans initialization

Prior probability for coin1:**0.31530369152313953**

Probability of head for coin1:**0.6659352656434103**

Prior probability for coin2:**0.20261831177194986**

Probability of head for coin2:**0.2587294546111306**

Prior probability for coin3:**0.482077996704911**

Probability of head for coin3:**0.8523617131186518**

1.3 Plots

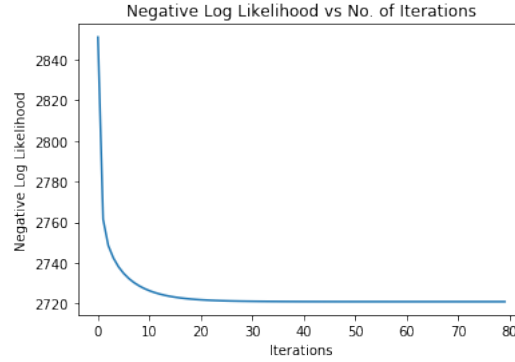


Figure 1: EM without KMeans

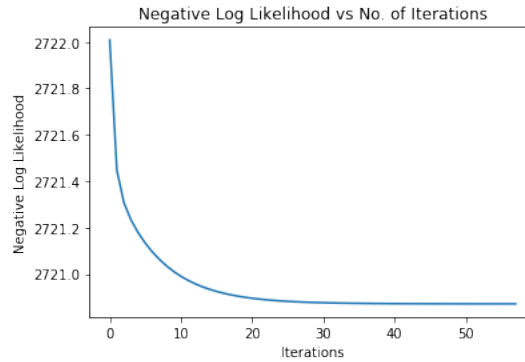


Figure 2: EM with Kmeans Initialization

1.4 Discussion

The EM algorithm is deterministic, meaning that upon repeated initializations, a given set of starting values will necessarily converge to the same solution. Hence, the start values become very important to obtain good final solution. The results from the experiment of EM with Kmeans parameter initialization got us very similar results to that of pure Kmeans implementation as well as closer than random values initialization. Additionally, we can also infer that although we get well formed clusters with both approaches, it might not be theoretically possible to get closer to actual ground truth with respect to the main coin experiment. Moreover, it is seen that EM with Kmeans converges much faster than EM with random initialization of parameters which can be seen in the plots which also show that curves depicted show similar trend of negative log likelihood barring actual values. We can see that negative log likelihood flattens out after 10-15 iterations in EM experiment. Overall key reflections of this exercise is we see how probability densities can be used to cluster set of points be it single cluster distribution or mixture density (in this case a mixture of binomial distribution) as depicted in the mixture model via EM algorithm. But it can be easily extended to mixture of gaussians with respective parameters update in the main algorithm. Mixture model allows to overcome

the limitations of Kmeans such as hard cluster assignment which helps in real time use cases where you can have different datapoint aligned with different hidden topics/latent semantics, a prime method such as Latent Dirichlet Allocation helps to cluster hidden topics from set of documents. Thus we can see that the EM algorithm can be generalized to different domains based on the use case.