

## Assignment 4

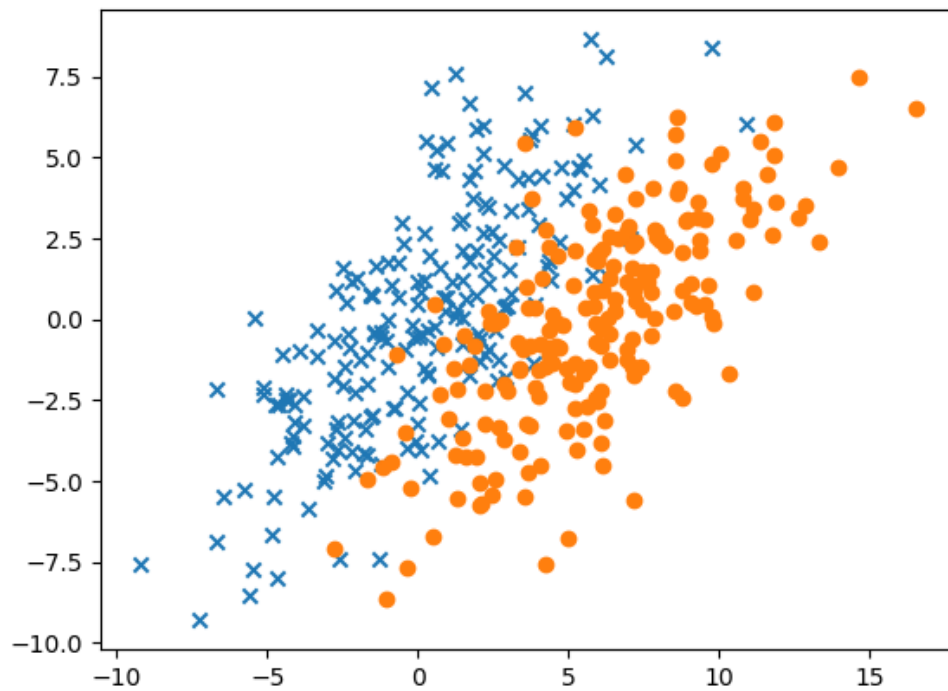
1.

(a)

Code for (a):

```
1 import ...
5
6 # TODO: Run this cell to generate the data
7 num_samples = 400
8 cov = np.array([[1., .7], [.7, 1.]]) * 10
9 mean_1 = [.1, .1]
10 mean_2 = [6., .1]
11
12 x_class1 = np.random.multivariate_normal(mean_1, cov, num_samples // 2)
13 x_class2 = np.random.multivariate_normal(mean_2, cov, num_samples // 2)
14 xy_class1 = np.column_stack((x_class1, np.zeros(num_samples // 2)))
15 xy_class2 = np.column_stack((x_class2, np.ones(num_samples // 2)))
16 data_full = np.row_stack([xy_class1, xy_class2])
17 np.random.shuffle(data_full)
18 data = data_full[:, :2]
19 labels = data_full[:, 2]
20
21 # TODO: Make a scatterplot for the data points showing the true cluster assignments of each p
22 plt.scatter(x_class1[:, 0], x_class1[:, 1], marker="x") # first class, x shape
23 plt.scatter(x_class2[:, 0], x_class2[:, 1], marker="o") # second class, circle shape
24 plt.show()
```

Visualize the generated data points:



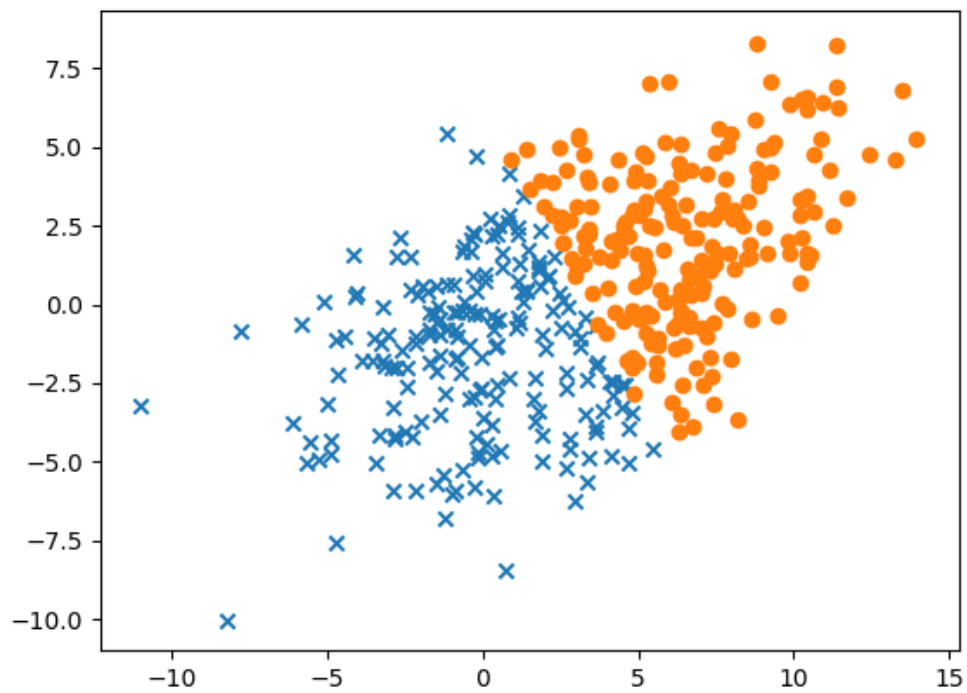
(b)

Code for (b):

```
36 # TODO: K-Means Assignment Step
37 def km_assignment_step(data, Mu):
38     """ Compute K-Means assignment step
39
40     Args:
41         data: a NxD matrix for the data points
42         Mu: a DxK matrix for the cluster means locations
43
44     Returns:
45         R_new: a NxK matrix of responsibilities
46     """
47
48     # Fill this in:
49     N, D = data.shape # Number of datapoints and dimension of datapoint
50     K = Mu.shape[1] # number of clusters
51     r = np.zeros((N, K))
52     for k in range(K):
53         r[:, k] = np.linalg.norm(data - Mu[:, k], axis=1)
54     arg_min = np.argmin(r, axis=1) # argmax/argmin along dimension 1
55     R_new = np.zeros((N, K)) # Set to zeros/ones with shape (N, K)
56     R_new[np.arange(N), arg_min] = 1 # Assign to 1
57     return R_new
58
59 # TODO: K-means Refitting Step
60 def km_refitting_step(data, R, Mu):
61     """ Compute K-Means refitting step.
62
63     Args:
64         data: a NxD matrix for the data points
65         R: a NxK matrix of responsibilities
66         Mu: a DxK matrix for the cluster means locations
67
68     Returns:
69         Mu_new: a DxK matrix for the new cluster means locations
70     """
71
72     N, D = data.shape # Number of datapoints and dimension of datapoint
73     K = R.shape[1] # number of clusters
74     Mu_new = data.T.dot(R) / sum(R)
75     return Mu_new
76
77 # TODO: Run this cell to call the K-means algorithm
78 N, D = data.shape
79 K = 2
80 max_iter = 100
81 class_init = np.random.binomial(1., .5, size=N)
82
83 N, D = data.shape
84 K = 2
85 max_iter = 100
86 class_init = np.random.binomial(1., .5, size=N)
87 R = np.vstack([class_init, 1 - class_init]).T
88
89 Mu = np.zeros([D, K])
90 Mu[:, 1] = 1.
91 R.T.dot(data), np.sum(R, axis=0)
92
93 costs = []
94 for it in range(max_iter):
95     R = km_assignment_step(data, Mu)
96     Mu = km_refitting_step(data, R, Mu)
97     costs.append(cost(data, R, Mu))
98
99 class_1 = np.where(R[:, 0])
100 class_2 = np.where(R[:, 1])
101
102 class_1 = np.where(R[:, 0])
103 class_2 = np.where(R[:, 1])
104
105 class_1_labels = labels[class_1[0]]
106 match_1 = max(sum(class_1_labels), len(class_1[0]) - sum(class_1_labels))
107 class_2_labels = labels[class_2[0]]
108 match_2 = max(sum(class_2_labels), len(class_2[0]) - sum(class_2_labels))
109 accuracy = (match_1 + match_2) / N
110 print("Misclassification error is ", 1 - accuracy)
111
112 # TODO: Make a scatterplot for the data points showing the K-Means cluster assignments of each point
113 plt.scatter(data[class_1[0], np.zeros(class_1[0].shape[0], dtype=int)],
114             data[class_1[0], np.ones(class_1[0].shape[0], dtype=int)], marker="x") # first class, x shape
115 plt.scatter(data[class_2[0], np.zeros(class_2[0].shape[0], dtype=int)],
116             data[class_2[0], np.ones(class_2[0].shape[0], dtype=int)], marker="o") # second class, circle shape
117 plt.show()
118 plt.plot(np.arange(max_iter), costs)
119 plt.show()
```

---

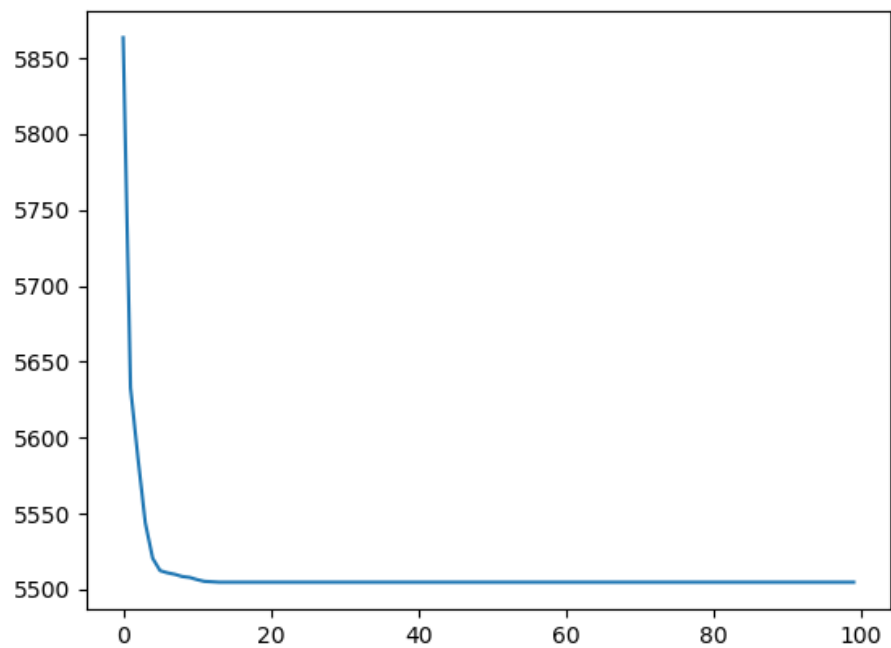
The resulting cluster assignments:



We have the misclassification rate as:

Misclassification error is 0.27

The cost vs. each iteration:



(c)

Code for (c):

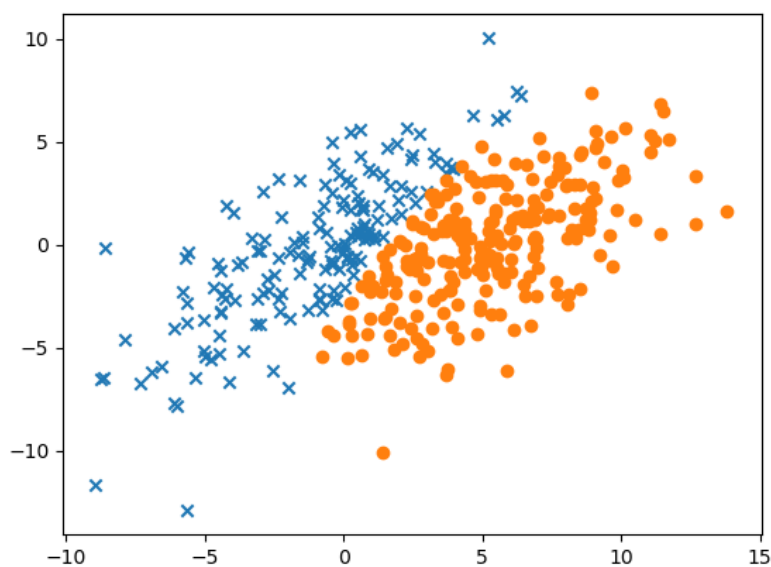
```
40 def log_likelihood(data, Mu, Sigma, Pi):
41     """ Compute log likelihood on the data given the Gaussian Mixture Parameters.
42
43     Args:
44         data: a NxD matrix for the data points
45         Mu: a DxK matrix for the means of the K Gaussian Mixtures
46         Sigma: a list of size K with each element being DxD covariance matrix
47         Pi: a vector of size K for the mixing coefficients
48
49     Returns:
50         L: a scalar denoting the log likelihood of the data given the Gaussian Mixture
51     """
52     # Fill this in:
53     N, D = data.shape[0], data.shape[1] # Number of datapoints and dimension of datapoint
54     K = Mu.shape[1] # number of mixtures
55     L, T = 0., 0.
56     for n in range(N):
57         T = 0
58         for k in range(K):
59             T += Pi[k] * normal_density(data[n], Mu[:, k], Sigma[k])
60             # Compute the likelihood from the k-th Gaussian weighted by the mixing coefficients
61         L += np.log(T)
62     return L
63
64 # TODO: Gaussian Mixture Expectation Step
65 def gm_e_step(data, Mu, Sigma, Pi):
66     """ Gaussian Mixture Expectation Step.
67
68     Args:
69         data: a NxD matrix for the data points
70         Mu: a DxK matrix for the means of the K Gaussian Mixtures
71         Sigma: a list of size K with each element being DxD covariance matrix
72         Pi: a vector of size K for the mixing coefficients
73
74     Returns:
75         Gamma: a NxK matrix of responsibilities
76     """
77     # Fill this in:
78     N, D = data.shape[0], data.shape[1] # Number of datapoints and dimension of datapoint
79     K = Mu.shape[1] # number of mixtures
80     Gamma = np.zeros((N,K)) # zeros of shape (N,K), matrix of responsibilities
81     for n in range(N):
82         for k in range(K):
83             Gamma[n, k] = Pi[k] * normal_density(data[n], Mu[:, k], Sigma[k])
84             Gamma[n, :] /= np.sum(Gamma, axis=1)[n]
85             # Normalize by sum across second dimension (mixtures)
86     return Gamma
87
88 # TODO: Gaussian Mixture Maximization Step
89 def gm_m_step(data, Gamma):
90     """ Gaussian Mixture Maximization Step.
91
92     Args:
93         data: a NxD matrix for the data points
94         Gamma: a NxK matrix of responsibilities
95
96     Returns:
97         Mu: a DxK matrix for the means of the K Gaussian Mixtures
98         Sigma: a list of size K with each element being DxD covariance matrix
99         Pi: a vector of size K for the mixing coefficients
100     """
101     # Fill this in:
102     N, D = data.shape[0], data.shape[1] # Number of datapoints and dimension of datapoint
103     K = Gamma.shape[1] # number of mixtures
104     Nk = sum(Gamma) # Sum along first axis
105     Mu = data.T.dot(Gamma) / Nk
106     Sigma = [0] * K
107     for k in range(K):
108         gamma_matrix = np.diag(Gamma[:, k])
109         Sigma[k] = (data - Mu[:, k]).T.dot(gamma_matrix).dot(data - Mu[:, k]) / Nk[k]
110     Pi = Nk / N
111     return Mu, Sigma, Pi
112
113
114 N, D = data.shape
115 K = 2
116 Mu = np.zeros([D, K])
117 Mu[:, 1] = 1.
118 Sigma = [np.eye(2), np.eye(2)]
119 Pi = np.ones(K) / K
120 Gamma = np.zeros([N, K]) # Gamma is the matrix of responsibilities
121
122 max_iter = 200
123
124 costs = []
125 log_likelihoods = []
126 for it in range(max_iter):
127     Gamma = gm_e_step(data, Mu, Sigma, Pi)
128     Mu, Sigma, Pi = gm_m_step(data, Gamma)
129     costs.append(cost(data, Gamma, Mu))
130     log_likelihoods.append(log_likelihood(data, Mu, Sigma, Pi))
131 # print(it, log_likelihood(data, Mu, Sigma, Pi)) # This function makes the computation longer, but good for debugging
132 # print(Gamma)
```

```

137 class_1 = np.where(Gamma[:, 0] >= .5)
138 class_2 = np.where(Gamma[:, 1] >= .5)
139
140 class_1_labels = labels[class_1[0]]
141 match_1 = max(sum(class_1_labels), len(class_1[0]) - sum(class_1_labels))
142 class_2_labels = labels[class_2[0]]
143 match_2 = max(sum(class_2_labels), len(class_2[0]) - sum(class_2_labels))
144 accuracy = (match_1 + match_2) / N
145 print("Misclassification error is ", 1 - accuracy)
146
147 # TODO: Make a scatterplot for the data points showing the Gaussian Mixture cluster assignments of each point
148 plt.scatter(data[class_1[0], np.zeros(class_1[0].shape[0], dtype=int)],
149            data[class_1[0], np.ones(class_1[0].shape[0], dtype=int)], marker="x") # first class, x shape
150 plt.scatter(data[class_2[0], np.zeros(class_2[0].shape[0], dtype=int)],
151            data[class_2[0], np.ones(class_2[0].shape[0], dtype=int)], marker="o") # second class, circle shape
152 plt.show()
153
154 plt.plot(np.arange(max_iter), log_likelihoods)
155 plt.show()

```

The resulting cluster assignments:



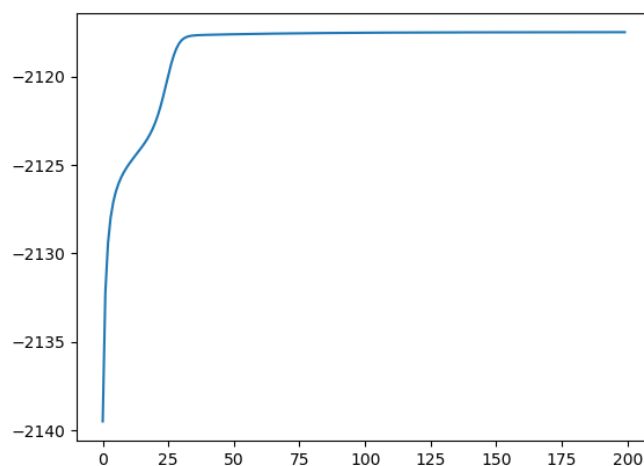
We have the misclassification rate as:

```

Misclassification error is 0.13249999999999995

```

Log-likelihood vs. number of iterations:



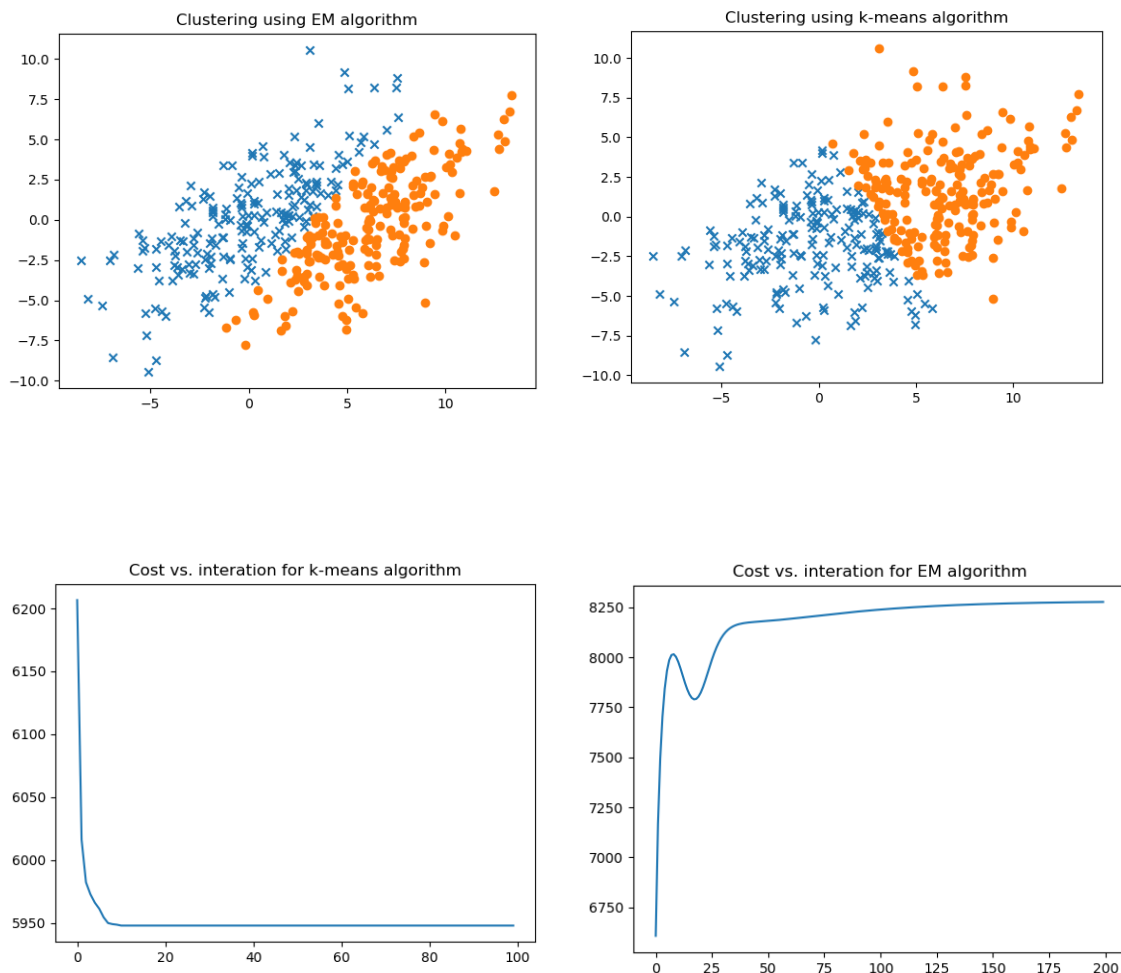
(d)

<a> We first compare the misclassification error for the two methods. As EM algorithm compares the probabilities of each cluster, it should get better accuracy. Our test output proved this. Error for EM algorithm is around 10%, which is smaller than k-means (around 25%).

<b> Comparing the cost vs. iteration plot for k-means and log-likelihood vs. iteration for EM algorithm, we see that EM algorithm takes more iteration to converge, since k-means converge at iteration around 10, but EM takes around 30. This may be a bottleneck for EM algorithm. In real practice, EM may take running time.

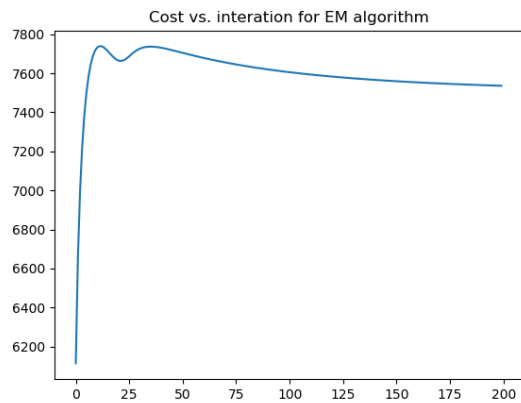
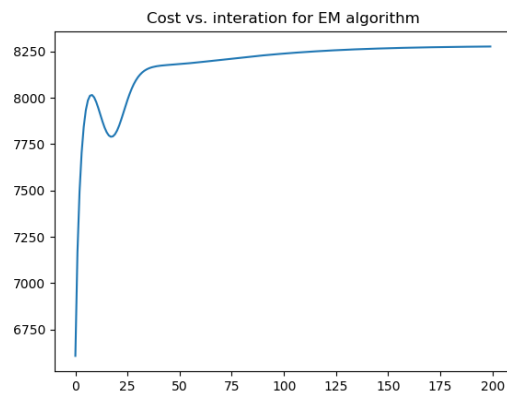
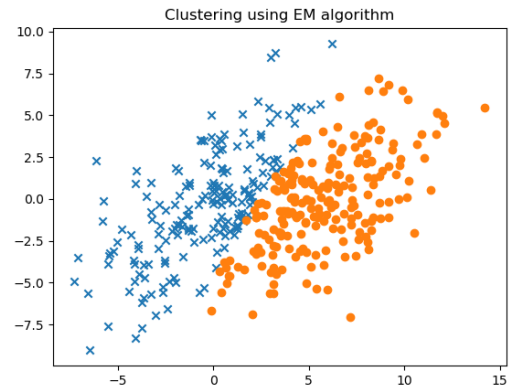
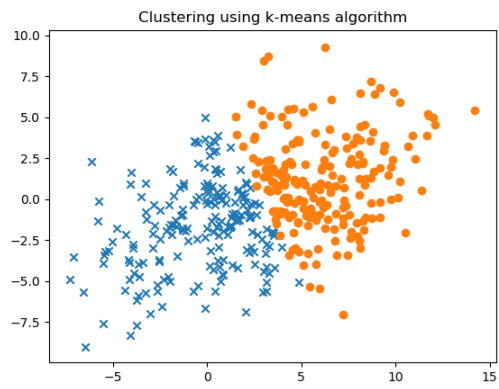
<c> We run the two algorithms on the same generated data for 5 times, we get the following resulting graph:

1<sup>st</sup> time:



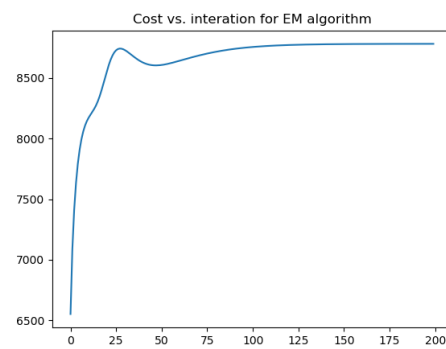
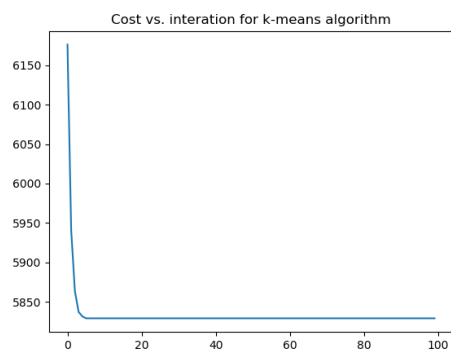
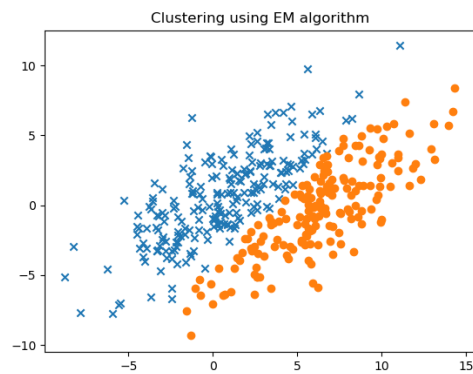
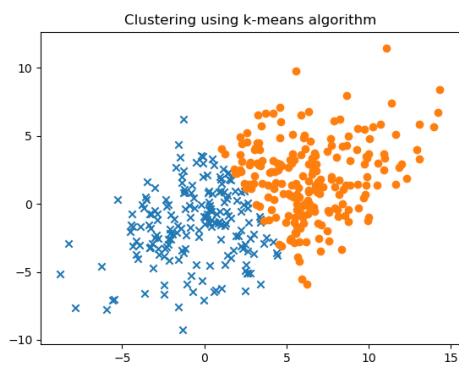
Misclassification error for k-means algorithm is 0.28  
Misclassification error for EM algorithm is 0.10499999999999998

2<sup>nd</sup> time:



Misclassification error for k-means algorithm is 0.245  
Misclassification error for EM algorithm is 0.11750000000000005

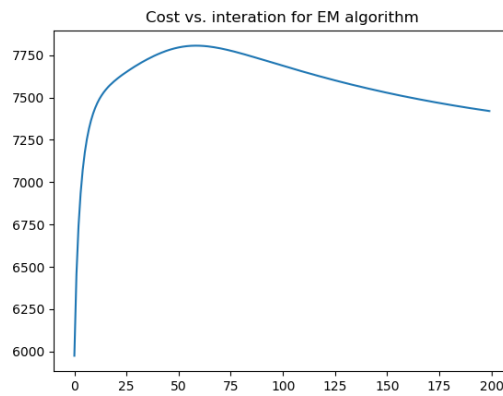
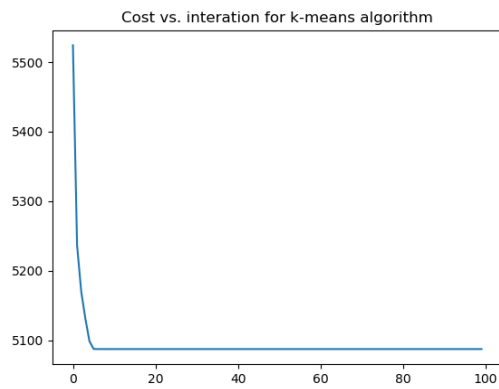
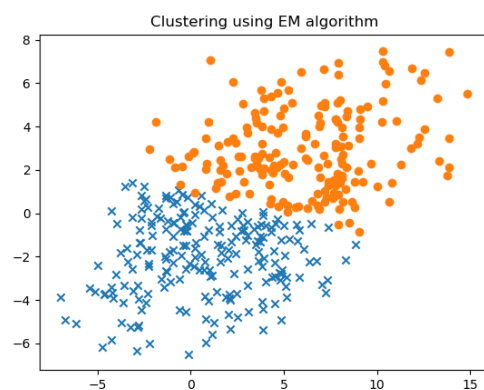
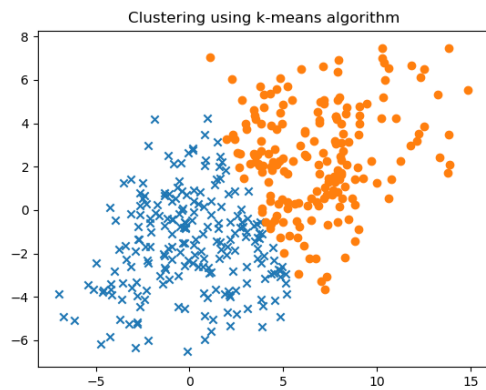
3<sup>rd</sup> time:



---

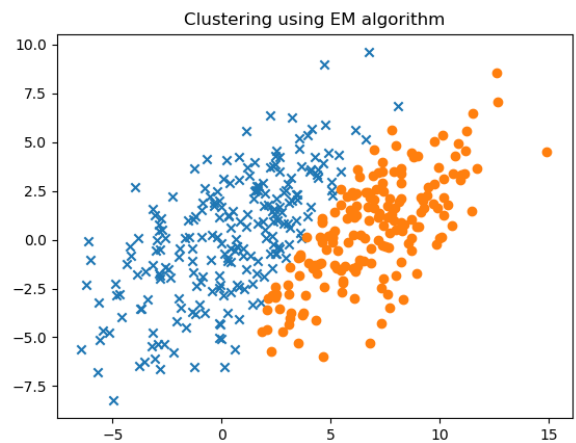
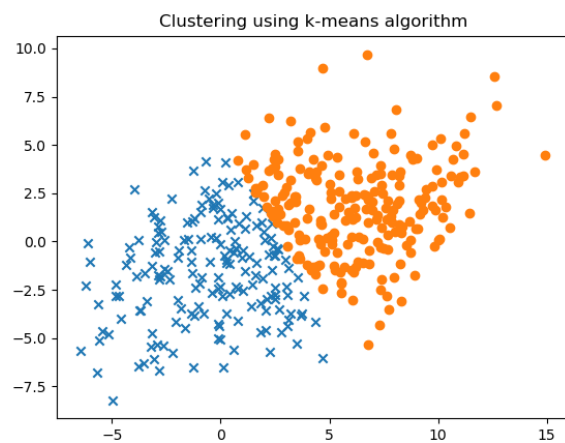
Misclassification error for k-means algorithm is 0.2524999999999995  
Misclassification error for EM algorithm is 0.0899999999999997

4<sup>th</sup> time:

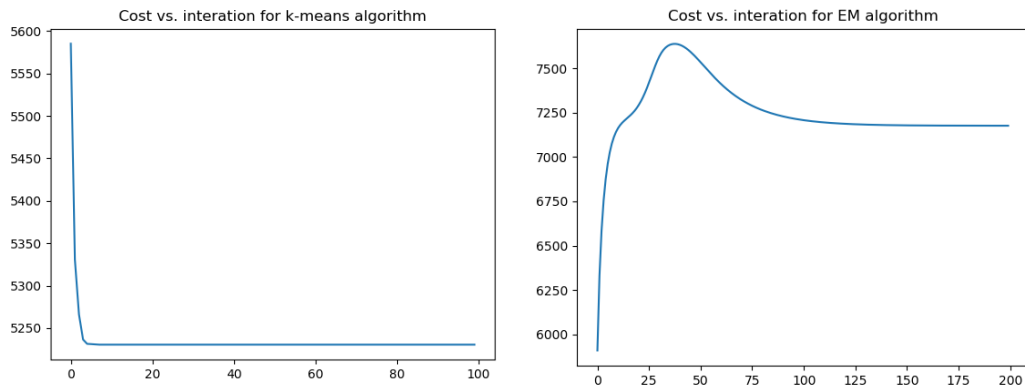


Misclassification error for k-means algorithm is 0.30000000000000004  
Misclassification error for EM algorithm is 0.41000000000000003

5<sup>th</sup> time:







Misclassification error for k-means algorithm is 0.24  
 Misclassification error for EM algorithm is 0.08999999999999997

K-means algorithm's performance is stable as the classification errors are all around 25%. However, in the 4<sup>th</sup> running, the classification error for the EM algorithm rises up to 41%, which does not follow the trend of the other 4 experiments. Cost functions converge for both algorithms. We conclude that the performance for k-means is stable but a bit low in accuracy, EM algorithm usually has higher accuracy, but depends on data realizations. For some “unfortunate” data, EM perform worse than k-means.

2.

2.1

Code for 2.1:

```

1 import ...
2
3
4
5
6 def qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy, init_beta=None, k_exp_sched=None):
7     """ Runs tabular Q learning algorithm for stochastic environment.
8
9     Args:
10         env: instance of environment object
11         num_iters (int): Number of episodes to run Q-learning algorithm
12         alpha (float): The learning rate between [0,1]
13         gamma (float): Discount factor, between [0,1]
14         epsilon (float): Probability in [0,1] that the agent selects a random move instead of
15             selecting greedily from Q value
16         max_steps (int): Maximum number of steps in the environment per episode
17         use_softmax_policy (bool): Whether to use softmax policy (True) or Epsilon-Greedy (False)
18         init_beta (float): If using stochastic policy, sets the initial beta as the parameter for the softmax
19         k_exp_sched (float): If using stochastic policy, sets hyperparameter for exponential schedule
20             on beta
21
22     Returns:
23         q_hat: A Q-value table shaped [num_states, num_actions] for environment with with num_states
24             number of states (e.g. num rows * num columns for grid) and num_actions number of possible
25             actions (e.g. 4 actions up/down/left/right)
26         steps_vs_iters: An array of size num_iters. Each element denotes the number
  
```

```

21
22 Returns:
23     q_hat: A Q-value table shaped [num_states, num_actions] for environment with with num_states
24         number of states (e.g. num rows * num columns for grid) and num_actions number of possible
25         actions (e.g. 4 actions up/down/left/right)
26     steps_vs_iters: An array of size num_iters. Each element denotes the number
27         of steps in the environment that the agent took to get to the goal
28         (capped to max_steps)
29
30     """
31     action_space_size = env.num_actions
32     state_space_size = env.num_states
33     q_hat = np.zeros(shape=(state_space_size, action_space_size))
34     steps_vs_iters = np.zeros(num_iters)
35
36     for i in range(num_iters):
37         # TODO: Initialize current state by resetting the environment
38         curr_state = env.reset()
39         num_steps = 0
40         done = False
41
42         # TODO: Keep looping while environment isn't done and less than maximum steps
43         while done is False and num_steps < max_steps:
44             num_steps += 1
45
46             # Choose an action using policy derived from either softmax Q-value
47             # or epsilon greedy
48             if use_softmax_policy:
49                 assert(init_beta is not None)
50                 assert(k_exp_sched is not None)
51                 # TODO: Boltzmann stochastic policy (softmax policy)
52                 beta = beta_exp_schedule(init_beta, i, k_exp_sched)
53                 # Call beta_exp_schedule to get the current beta value
54                 action = softmax_policy(q_hat, beta, curr_state)
55             else:
56                 # TODO: Epsilon-greedy
57                 action = epsilon_greedy(q_hat, epsilon, curr_state, action_space_size)
58
59             # TODO: Execute action in the environment and observe the next state, reward, and done flag
60             next_state, reward, done = env.step(action)
61
62             # TODO: Update Q_value
63             if next_state != curr_state:
64                 new_value = reward + gamma * np.max(q_hat[next_state]) - q_hat[curr_state, action]
65                 # TODO: Use Q-learning rule to update q_hat for the curr_state and action:
66                 # i.e.,  $Q(s,a) \leftarrow Q(s,a) + \alpha * [reward + \gamma * \max_{a'}(Q(s',a')) - Q(s,a)]$ 
67                 q_hat[curr_state, action] = q_hat[curr_state, action] + alpha * new_value
68
69             # TODO: Update the current state to be the next state
70             curr_state = next_state
71
72     steps_vs_iters[i] = num_steps
73     return q_hat, steps_vs_iters
74
75 def epsilon_greedy(q_hat, epsilon, state, action_space_size):
76     """ Chooses a random action with p_rand_move probability,
77         otherwise choose the action with highest Q value for
78         current observation
79
80     Args:
81         q_hat: A Q-value table shaped [num_rows, num_col, num_actions] for
82             grid environment with num_rows rows and num_col columns and num_actions
83
84         number of possible actions
85         epsilon (float): Probability in [0,1] that the agent selects a random
86         move instead of selecting greedily from Q value
87         state: A 2-element array with integer element denoting the row and column
88         that the agent is in
89         action_space_size (int): number of possible actions
90
91     Returns:
92         action (int): A number in the range [0, action_space_size-1]
93         denoting the action the agent will take
94
95     """
96     # TODO: Implement your code here
97     # Hint: Sample from a uniform distribution and check if the sample is below
98     # a certain threshold
99     q_state = q_hat[state]
100     q_max_move = np.argmax(q_state)
101     probabilities = [epsilon, 1 - epsilon]
102     random_move = np.random.choice(action_space_size)
103     action = np.random.choice([random_move, q_max_move], p=probabilities)
104
105     if np.all(q_state == 0):
106         action = random_move
107     return action

```

```

110 def softmax_policy(q_hat, beta, state):
111     """ Choose action using policy derived from Q, using
112         softmax of the Q values divided by the temperature.
113
114     Args:
115         q_hat: A Q-value table shaped [num_rows, num_col, num_actions] for
116             grid environment with num_rows rows and num_col columns
117         beta (float): Parameter for controlling the stochasticity of the action
118         obs: A 2-element array with integer element denoting the row and column
119             that the agent is in
120
121     Returns:
122         action (int): A number in the range [0, action_space_size-1]
123             denoting the action the agent will take
124
125     """
126     # TODO: Implement your code here
127     # Hint: use the stable_softmax function defined below
128     q_state = q_hat[state]
129     max_q = np.max(q_state)
130     z = np.exp(beta * (q_state - max_q))
131     probabilities = z / np.sum(z)
132     actions = np.arange(q_hat.shape[1])
133     action = np.random.choice(actions, p=probabilities)
134     return action

```

We resolve the issue of “ties” in  $Q\_hat$  matrix by adding code in epsilon-greedy policy. For softmax policy, the algorithm will randomly assign a move when we have ties in several 0’s in  $Q\_hat$ .

## 2.2.1

(a)

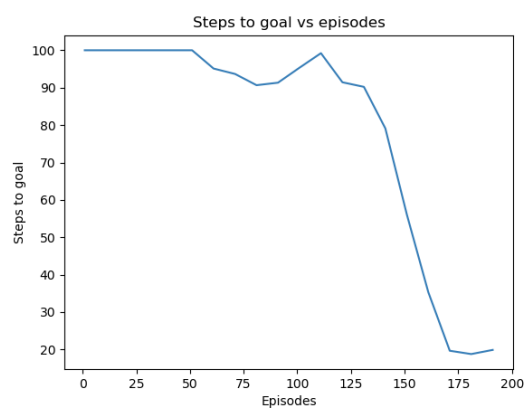
Code for 2.2.1(a):

```

1 import ...
4
5 # TODO: Fill this in
6 num_iters = 200
7 alpha = 1.0
8 gamma = 0.9
9 epsilon = 0.1
10 max_steps = 100
11 use_softmax_policy = False
12
13 # TODO: Instantiate the MazeEnv environment with default arguments
14 env = MazeEnv()
15
16 # TODO: Run Q-learning:
17 q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_p
18
19 # TODO: Plot the steps vs iterations
20 plot_steps_vs_iters(steps_vs_iters)
21
22
23 # TODO: plot the policy from the Q value
24 plot_policy_from_q(q_hat, env)

```

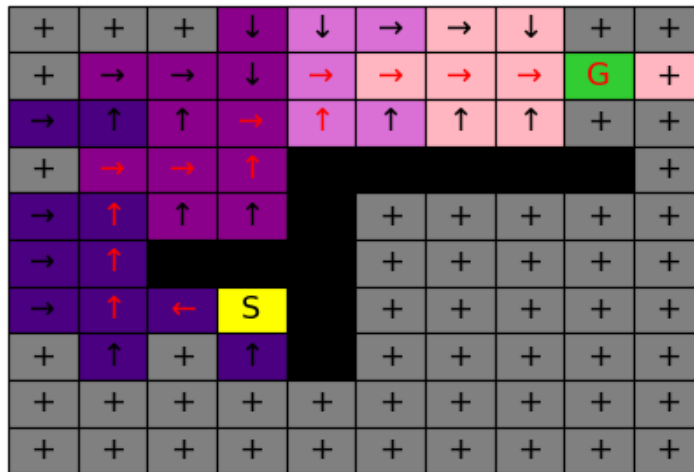
Plot the steps to goal vs training iterations (episodes):



---

Steps to goal converges according to our implementation.

Visualize the learned greedy policy from the Q values:



Following the red arrows, we see that the agent found the optimal path to goal.

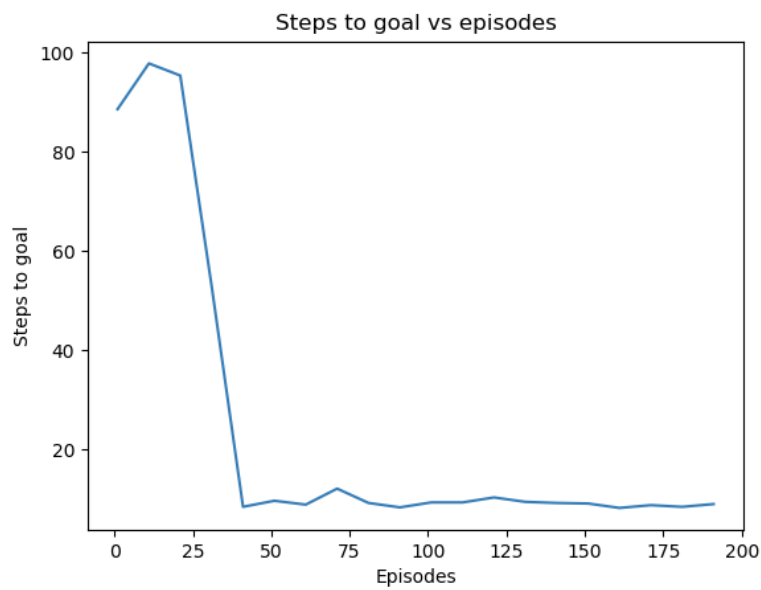
### 2.2.1(b)

Code for 2.2.1(b):

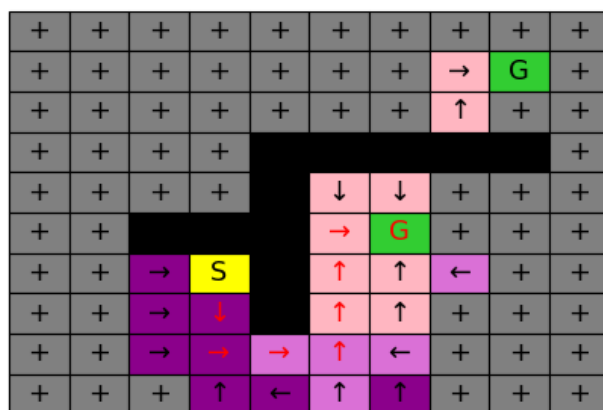
```
1 import ...
4
5
6 # TODO: Fill this in (same as before)
7 num_iters = 200
8 alpha = 1.0
9 gamma = 0.9
10 epsilon = 0.1
11 max_steps = 100
12 use_softmax_policy = False
13
14 # TODO: Set the goal
15 goal_locs = [[1, 8], [5, 6]]
16 env = MazeEnv(goals=goal_locs)
17
18 # TODO: Run Q-learning:
19 q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_p
20
21 # TODO: Plot the steps vs iterations
22 plot_steps_vs_iters(steps_vs_iters)
23
24 # TODO: plot the policy from the Q values
25 plot_policy_from_q(q_hat, env)
26
```

---

Plot the steps to goal vs training iterations (episodes):



Visualize the learned greedy policy from the Q values:



Following the red arrows, the agent gets to the nearest goal (5,6). Since the total distance is shorter this time, steps to goal converges faster than 2.2.1(a).

2.2.2

2.2.2(a)

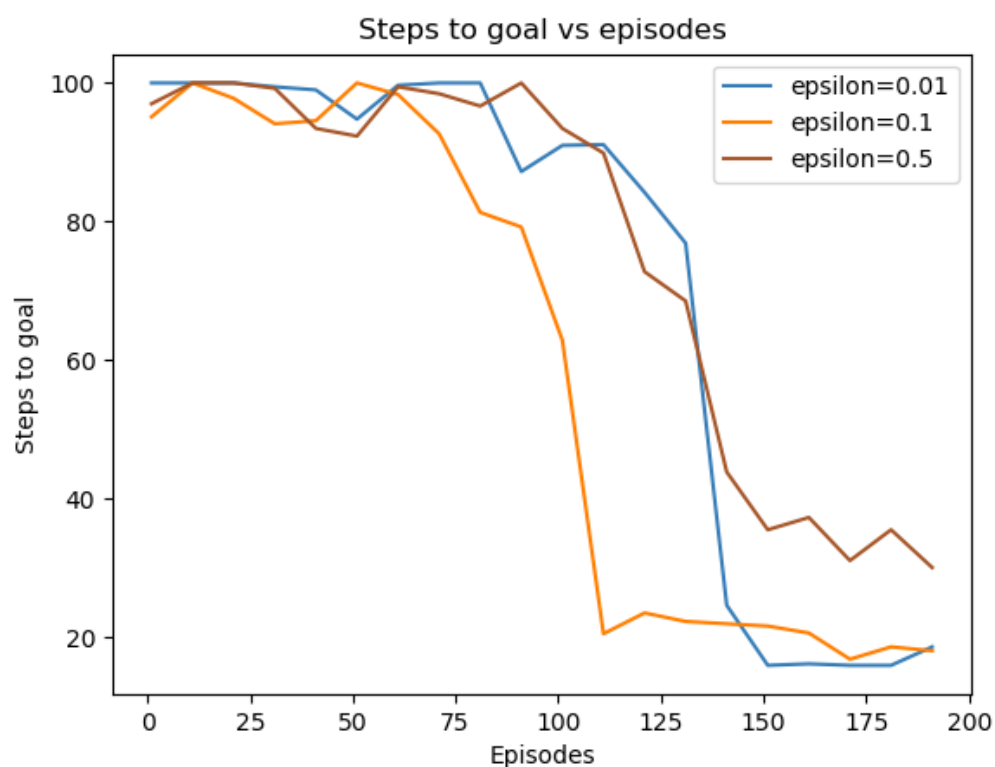
Code for 2.2.2(a):

```

1 import ...
4
5
6 # TODO: Fill this in (same as before)
7 num_iters = 200
8 alpha = 1.0
9 gamma = 0.9
10 epsilon = 0.1
11 max_steps = 100
12 use_softmax_policy = False
13
14 # TODO: set the epsilon lists in increasing order:
15 epsilon_list = [0.01, 0.1, 0.5]
16
17 env = MazeEnv()
18
19 steps_vs_iters_list = []
20 for epsilon in epsilon_list:
21     q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy)
22     steps_vs_iters_list.append(steps_vs_iters)
23 # TODO: Plot the results
24 label_list = ["epsilon={}".format(eps) for eps in epsilon_list]
25 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)

```

Plot the steps to goal vs training iterations (episodes):



It seems that  $\epsilon = 0.1$  works best, since it converges the fastest. This indicates that too low and too high exploration rates will influence performance negatively.

2.2.2(b)

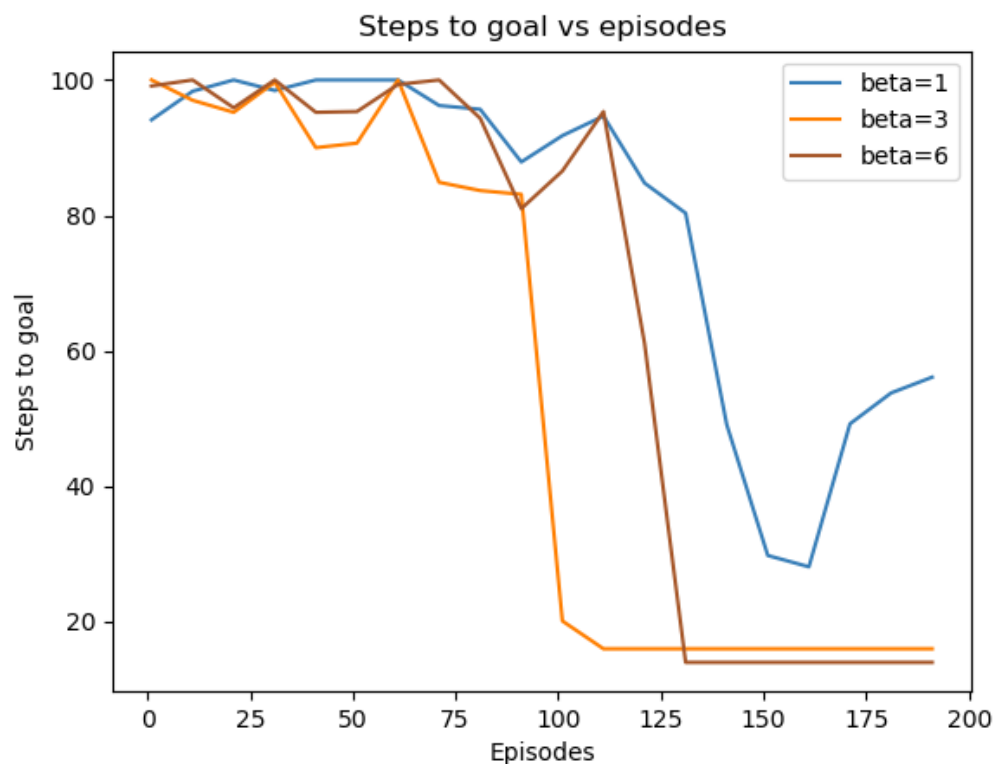
Code for 2.2.2(b):

```

1 import ...
2
3 # TODO: Fill this in for Static Beta with softmax of Q-values
4 num_iters = 200
5 alpha = 1.0
6 gamma = 0.9
7 epsilon = 0.1
8 max_steps = 100
9
10
11
12 # TODO: Set the beta
13 beta_list = [1, 3, 6]
14 use_softmax_policy = True
15 k_exp_schedule = 0.0 # (float) choose k such that we have a constant beta during training
16
17 env = MazeEnv()
18 steps_vs_iters_list = []
19 for beta in beta_list:
20     q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon,
21                                   max_steps, use_softmax_policy, beta, k_exp_schedule)
22     steps_vs_iters_list.append(steps_vs_iters)
23
24 label_list = ["beta={}".format(beta) for beta in beta_list]
25 # TODO:
26 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)

```

Plot the steps to goal vs training iterations (episodes):



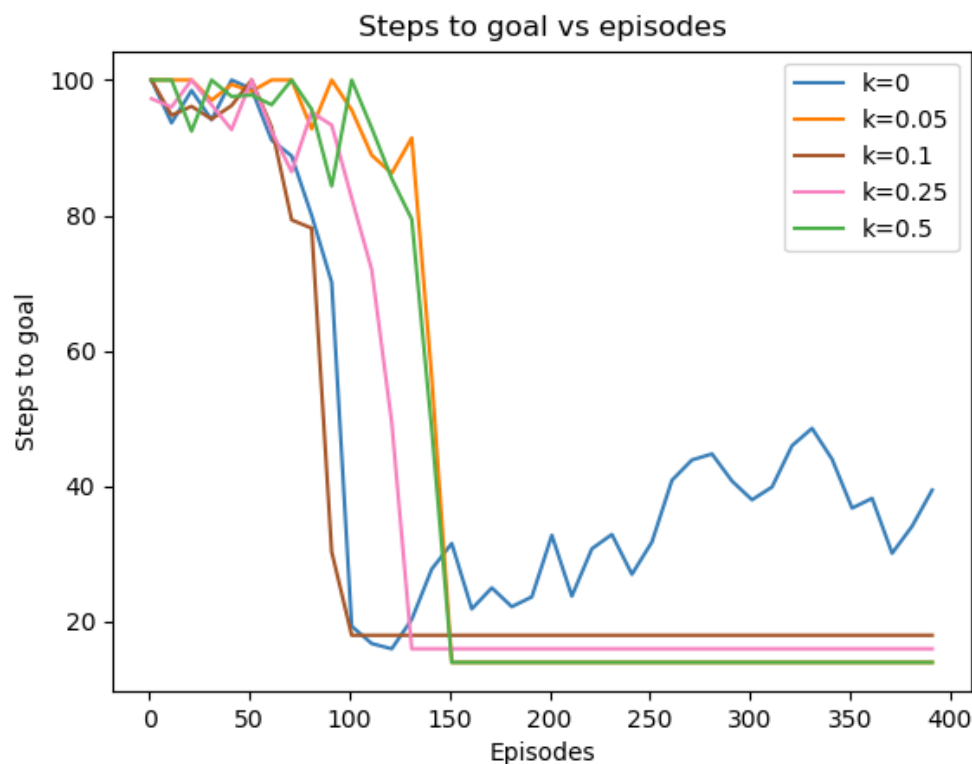
Beta = 3 has the best performance according to our experiments. Again, like the previous problem, too small and too large beta also negatively affect performance.

### 2.2.2(c)

Code for 2.2.2(c):

```
1 import ...
4 # TODO: Fill this in for Dynamic Beta
5 num_iters = 600
6 alpha = 1.0
7 gamma = 0.9
8 epsilon = 0.1
9 max_steps = 100
10
11 # TODO: Set the beta
12 beta = 1.0
13 use_softmax_policy = True
14 k_exp_schedule_list = [0, 0.05, 0.1, 0.25, 0.5]
15 env = MazeEnv()
16
17 steps_vs_iters_list = []
18 for k_exp_schedule in k_exp_schedule_list:
19     q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon,
20                                   max_steps, use_softmax_policy, beta, k_exp_schedule)
21     steps_vs_iters_list.append(steps_vs_iters)
22
23 # TODO: Plot the steps vs iterations
24 label_list = ["k={}".format(k_exp_schedule) for k_exp_schedule in k_exp_schedule_list]
25 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)
```

Plot the steps to goal vs training iterations (episodes):



Clearly from the plot, it is better to have  $k \neq 0$ , since in the case  $k = 0$ , steps to goal does not converge even after 400 iterations. It seems that as long as  $k \neq 0$ , the algorithm works well. If we want to converge as quick as possible, we can choose  $k = 0.1$ . This



also indicates that the convergence rate depends on the moderately chosen values of hyperparameters, since small  $k$  or large  $k$  won't converge as fast as  $k = 0.1$ , the medium value.

### 2.2.3

Code for 2.2.3(a):

```
142 class ProbabilisticMazeEnv(MazeEnv):
143     """ (Q2.3) Hints: you can refer the implementation in MazeEnv
144     """
145
146     def __init__(self, goals=[[2, 8]], p_random=0.05):
147         """ Probabilistic Maze Environment
148
149         Args:
150             goals (list): list of goals coordinates
151             p_random (float): random action rate
152         """
153
154         super().__init__(goals=goals)
155         self.p_random = p_random
156
157     def step(self, a):
158         done, reward = False, 0.0
159         next_obs = copy.copy(self.obs)
160
161         random_action = np.random.choice(self.num_actions)
162         a = np.random.choice([a, random_action], p=[1-self.p_random, self.p_random])
163
164         random_action = np.random.choice(self.num_actions)
165         a = np.random.choice([a, random_action], p=[1-self.p_random, self.p_random])
166
167         if a == 0:
168             next_obs[0] = next_obs[0] - 1
169         elif a == 1:
170             next_obs[1] = next_obs[1] + 1
171         elif a == 2:
172             next_obs[1] = next_obs[1] - 1
173         elif a == 3:
174             next_obs[0] = next_obs[0] + 1
175         else:
176             raise Exception("Action is Not Valid")
177
178         if self.is_valid_obs(next_obs):
179             self.obs = next_obs
180
181         if self.map[self.obs[0], self.obs[1]] == -1:
182             reward = self.reward
183             done = True
184
185         state = self.get_state_from_coords(self.obs[0], self.obs[1])
186
187         return state, reward, done
```

We override the initializer and the step() method in the subclass ProbabilisticMazeEnv.

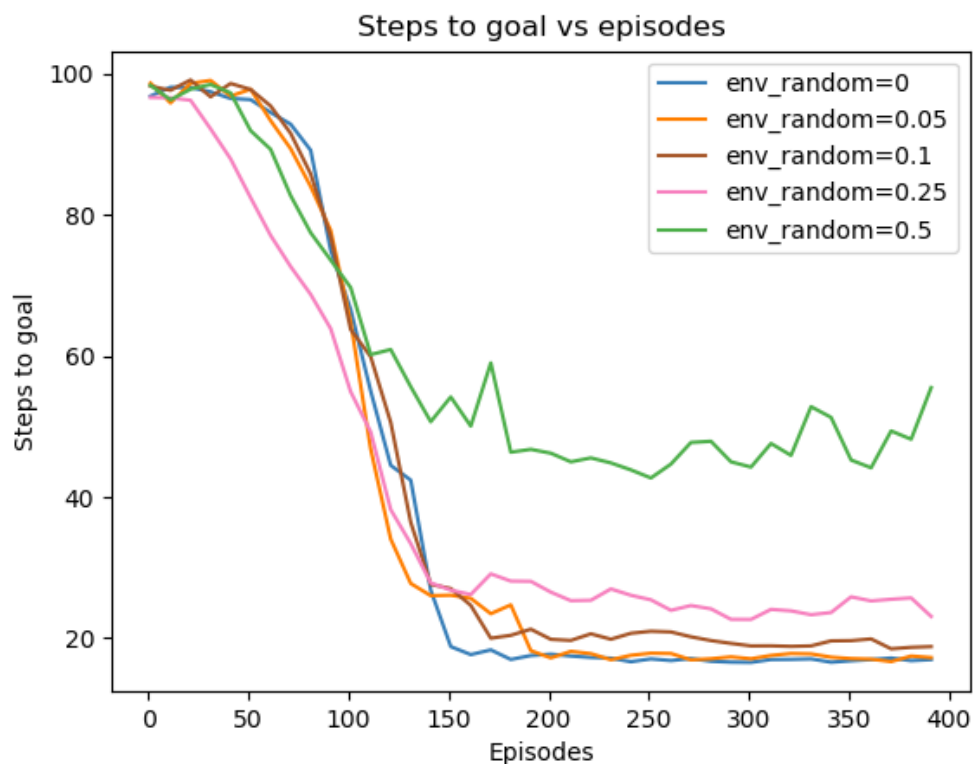
Code for 2.2.3(b):

```

1 import ...
5 # TODO: Use the same parameters as in the first part, except change alpha
6 num_iters = 400
7 alpha = 0.5
8 gamma = 0.9
9 epsilon = 0.1
10 max_steps = 100
11 use_softmax_policy = False
12 # Set the environment probability of random
13 env_p_rand_list = [0, 0.05, 0.1, 0.25, 0.5]
14 steps_vs_iters_list = []
15 for env_p_rand in env_p_rand_list:
16     # Instantiate with ProbabilisticMazeEnv
17     env = ProbabilisticMazeEnv(p_random=env_p_rand)
18
19     # Note: We will repeat for several runs of the algorithm to make the result less noisy
20     avg_steps_vs_iters = np.zeros(num_iters)
21     for i in range(10):
22         q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_s
23         avg_steps_vs_iters += steps_vs_iters
24     avg_steps_vs_iters /= 10
25     steps_vs_iters_list.append(avg_steps_vs_iters)
26 label_list = ["env_random={}".format(env_p_rand) for env_p_rand in env_p_rand_list]
27 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)

```

Plot the steps to goal vs training iterations (episodes):



Take  $\text{env\_random} = 0.05$  or  $0.1$  if we want the fastest convergence. It should not be too large, since  $\text{env\_random} = 0.5$  does not converge after 400 iterations.

---

### 2.3

As a summary, we realize that if we want the fastest convergence rate, we have to set the parameters carefully, not too large and not too small. In general, it's better for them to be non-zero to make the algorithm converge.