```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4

Enter your authorization code:
..........
Mounted at /content/drive

# 1. Unsupervised Learning

```
1 %matplotlib inline
2 import scipy
3 import numpy as np
4 import itertools
5 import matplotlib.pyplot as plt
6 import time
```

## 1. Generating the data

First, we will generate some data for this problem. Set the number of points $N = 400$, their dimension $D = 2$, and the number of clusters $K = 2$, and generate data from the distribution $p(x|z = k) = \mathcal{N}(\mu_k, \Sigma_k)$. Sample $200$ data points for $k = 1$ and 200 for $k = 2$, with

$$\mu_1 = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}, \mu_2 = \begin{bmatrix} 6.0 \\ 0.1 \end{bmatrix} \text{ and } \Sigma_1 = \Sigma_2 = \begin{bmatrix} 10 & 7 \\ 7 & 10 \end{bmatrix}$$

Here, $N = 400$. Since you generated the data, you already know which sample comes from which class. Run the cell in the IPython notebook to generate the data.

```
1 # TODO: Run this cell to generate the data
2 num_samples = 400
3 cov = np.array([[1., .7], [.7, 1.]]) * 10
4 mean_1 = [.1, .1]
5 mean_2 = [6., .1]
6
7 x_class1 = np.random.multivariate_normal(mean_1, cov, num_samples // 2)
8 x_class2 = np.random.multivariate_normal(mean_2, cov, num_samples // 2)
9 xy_class1 = np.column_stack((x_class1, np.zeros(num_samples // 2)))
10 xy_class2 = np.column_stack((x_class2, np.ones(num_samples // 2)))
11 data_full = np.row_stack([xy_class1, xy_class2])
12 np.random.shuffle(data_full)
13 data = data_full[:, :2]
14 labels = data_full[:, 2]
```

Make a scatter plot of the data points showing the true cluster assignment of each point using different color codes and shape (x for first class and circles for second class):

```
1 # TODO: Make a scatterplot for the data points showing the true cluster assignments of each point
2 # plt.plot(...) # first class, x shape
3 # plt.plot(...) # second class, circle shape
4 plt.plot(x_class1[:, 0], x_class1[:, 1], 'x', c='red')
5 plt.plot(x_class2[:, 0], x_class2[:, 1], 'o', c='green')
6 plt.show()
```

## 2. Implement and Run K-Means algorithm

Now, we assume that the true class labels are not known. Implement the k-means algorithm for this problem. Write two functions: `km_assignment_step`, and `km_refitting_step` as given in the lecture (Here, `km_` means k-means). Identify the correct arguments, and the order to run them. Initialize the algorithm with

$$\hat{\mu}_1 = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix}, \hat{\mu}_2 = \begin{bmatrix} 1.0 \\ 1.0 \end{bmatrix}$$

and run it until convergence. Show the resulting cluster assignments on a scatter plot either using different color codes or shape or both. Also plot the cost vs. the number of iterations. Report your misclassification error.

```python
1 def cost(data, R, Mu):
2     N, D = data.shape
3     K = Mu.shape[1]
4     J = 0
5     for k in range(K):
6         J += np.sum(np.dot(np.linalg.norm(data - np.array([Mu[:, k], ] * N), axis=1)**2, R))
7     return J
```

```python
1 # TODO: K-Means Assignment Step
2 def km_assignment_step(data, Mu):
3     """ Compute K-Means assignment step
4
5     Args:
6         data: a NxD matrix for the data points
7         Mu: a DxK matrix for the cluster means locations
8
9     Returns:
10         R_new: a NxK matrix of responsibilities
11     """
12     N, D = data.shape  # Number of datapoints and dimension of datapoint
13     K = Mu.shape[1]  # number of clusters
14     r = np.zeros([N, K])
15     #a matrix of NxK dimension for the distances of the the N datapoints to each of the K cluster centers.
16     for k in range(K):
17         r[:, k] = np.linalg.norm(data - np.array([Mu[:, k], ] * N), axis=1)**2
18     arg_min = np.argmin(r, axis=1) # argmax/argmin along dimension 1
19     R_new = np.zeros([N, K]) # Set to zeros/ones with shape (N, K)
20     R_new[np.arange(N),arg_min] = 1 # Assign to 1
21     return R_new
```

```python
1 # TODO: K-means Refitting Step
2 def km_refitting_step(data, R, Mu):
3     """ Compute K-Means refitting step.
4
5     Args:
6         data: a NxD matrix for the data points
7         R: a NxK matrix of responsibilities
8         Mu: a DxK matrix for the cluster means locations
9
10     Returns:
11         Mu_new: a DxK matrix for the new cluster means locations
12     """
13     N, D = data.shape
14     K = Mu.shape[1]
15     Mu_new = np.zeros((D, K))
16     for k in range(K):
17         index = np.where(R[:, k] == 1)
18         Mu_new[:, k] = np.mean(data[index, :], axis = 1)
19     return Mu_new
```

```python
1 N, D = data.shape
2 K = 2
3 max_iter = 100
4 class_init = np.random.binomial(1., .5, size=N)
5 R = np.vstack([class_init, 1 - class_init]).T
6
7 Mu = np.zeros([D, K])
8 Mu[:, 1] = 1.
9 R.T.dot(data), np.sum(R, axis=0)
10
11 iteration = []
12 clst = []
13 for it in range(max_iter):
14     R = km_assignment_step(data, Mu)
15     Mu = km_refitting_step(data, R, Mu)
16     c = cost(data, R, Mu)
17     iteration.append(it)
18     clst.append(c)
19     print(it, c)
20
21 class_1 = np.where(R[:, 0])
22 class_2 = np.where(R[:, 1])
```

```
 0  37292.58165797178
 1  37250.88887891708
 2  37167.849589052785
 3  37036.73687303193
 4  36928.79657380079
 5  36844.57439989758
 6  36773.98217006367
 7  36739.05559204959
 8  36726.52159971727
 9  36726.52159971727
10  36726.52159971727
11  36726.52159971727
12  36726.52159971727
13  36726.52159971727
14  36726.52159971727
15  36726.52159971727
16  36726.52159971727
17  36726.52159971727
18  36726.52159971727
19  36726.52159971727
20  36726.52159971727
21  36726.52159971727
22  36726.52159971727
23  36726.52159971727
24  36726.52159971727
25  36726.52159971727
26  36726.52159971727
27  36726.52159971727
28  36726.52159971727
29  36726.52159971727
30  36726.52159971727
31  36726.52159971727
32  36726.52159971727
33  36726.52159971727
34  36726.52159971727
35  36726.52159971727
36  36726.52159971727
37  36726.52159971727
38  36726.52159971727
39  36726.52159971727
40  36726.52159971727
41  36726.52159971727
42  36726.52159971727
43  36726.52159971727
44  36726.52159971727
45  36726.52159971727
46  36726.52159971727
47  36726.52159971727
48  36726.52159971727
49  36726.52159971727
50  36726.52159971727
51  36726.52159971727
52  36726.52159971727
53  36726.52159971727
54  36726.52159971727
55  36726.52159971727
56  36726.52159971727
57  36726.52159971727
58  36726.52159971727
59  36726.52159971727
60  36726.52159971727
61  36726.52159971727
62  36726.52159971727
63  36726.52159971727
64  36726.52159971727
65  36726.52159971727
66  36726.52159971727
67  36726.52159971727
68  36726.52159971727
69  36726.52159971727
70  36726.52159971727
71  36726.52159971727
72  36726.52159971727
73  36726.52159971727
74  36726.52159971727
75  36726.52159971727
76  36726.52159971727
77  36726.52159971727
78  36726.52159971727
79  36726.52159971727
80  36726.52159971727
81  36726.52159971727
82  36726.52159971727
83  36726.52159971727
84  36726.52159971727
85  36726.52159971727
86  36726.52159971727
87  36726.52159971727
88  36726.52159971727
89  36726.52159971727
90  36726.52159971727
91  36726.52159971727
92  36726.52159971727
93  36726.52159971727
94  36726.52159971727
95  36726.52159971727
```
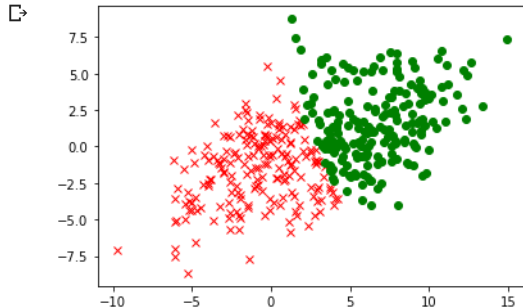
```
96 36726.52159971727
97 36726.52159971727
98 36726.52159971727
99 36726.52159971727
```

```
1 # TODO: Make a scatterplot for the data points showing the K-Means cluster assignments of each point
2 # plt.plot(...) # first class, x shape
3 # plt.plot(...) # second class, circle shape
4 plt.plot(data[class_1, 0], data[class_1,1], 'x', color='red')
5 plt.plot(data[class_2, 0], data[class_2,1], 'o', color='green')
6 plt.show()
```
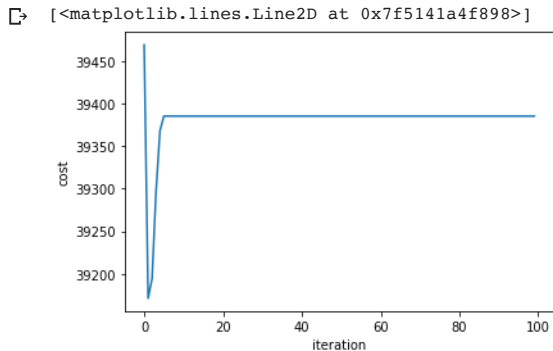


```
1 plt.ylabel('cost')
2 plt.xlabel('iteration')
3 plt.plot(iterations, costs)
```

[<matplotlib.lines.Line2D at 0x7f5141a4f898>]



```
1 labels_predict = np.argwhere(R == 1)[:, 1]
2 print("Error rate: ", np.mean(labels != labels_predict))
```

## 3. Implement EM algorithm for Gaussian mixtures

Next, implement the EM algorithm for Gaussian mixtures. Write three functions: `log_likelihood`, `gm_e_step`, and `gm_m_step` as given in the lecture. Identify the correct arguments, and the order to run them. Initialize the algorithm with means as in Qs 2.1 k-means initialization, covariances with $\hat{\Sigma}_1 = \hat{\Sigma}_2 = I$, and $\hat{\pi}_1 = \hat{\pi}_2$.

In addition to the update equations in the lecture, for the M (Maximization) step, you also need to use this following equation to update the covariance $\Sigma_k$:

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^{N} r_k^{(n)} (\mathbf{x}^{(n)} - \hat{\mu}_k)(\mathbf{x}^{(n)} - \hat{\mu}_k)^\top$$

Run the algorithm until convergence and show the resulting cluster assignments on a scatter plot either using different color codes or shape or both. Also plot the log-likelihood vs. the number of iterations. Report your misclassification error.

```
1 def normal_density(x, mu, Sigma):
2     return np.exp(-.5 * np.dot(x - mu, np.linalg.solve(Sigma, x - mu))) \
3         / np.sqrt(np.linalg.det(2 * np.pi * Sigma))
```

```
1 def log_likelihood(data, Mu, Sigma, Pi):
2     """ Compute log likelihood on the data given the Gaussian Mixture Parameters.
3
4     Args:
5         data: a NxD matrix for the data points
6         Mu: a DxK matrix for the means of the K Gaussian Mixtures
7         Sigma: a list of size K with each element being DxD covariance matrix
8         Pi: a vector of size K for the mixing coefficients
9
10     Returns:
11         L: a scalar denoting the log likelihood of the data given the Gaussian Mixture
```

```
12      """
13      # Fill this in:
14      # N, D = ...  # Number of datapoints and dimension of datapoint
15      # K = ... # number of mixtures
16      # L, T = 0., 0.
17      # for n in range(N):
18      #     for k in range(K):
19      #         # T += ... # Compute the likelihood from the k-th Gaussian weighted by the mixing coefficients
20      #     L += np.log(T)
21      # return L
22
23      N, D = data.shape  # Number of datapoints and dimension of datapoint
24      K = Mu.shape[1] # number of mixtures
25      L, T = 0., 0.
26      for n in range(N):
27          T = 0.
28          for k in range(K):
29              T += Pi[k] * normal_density(data[n], Mu[:, k], Sigma[k]) # Compute the likelihood from the k-th Gaussian weighted
30          L += np.log(T)
31      return L
```

```
1 # TODO: Gaussian Mixture Expectation Step
2 def gm_e_step(data, Mu, Sigma, Pi):
3     """ Gaussian Mixture Expectation Step.
4
5     Args:
6         data: a NxD matrix for the data points
7         Mu: a DxK matrix for the means of the K Gaussian Mixtures
8         Sigma: a list of size K with each element being DxD covariance matrix
9         Pi: a vector of size K for the mixing coefficients
10
11     Returns:
12         Gamma: a NxK matrix of responsibilities
13     """
14     # Fill this in:
15     # N, D = ... # Number of datapoints and dimension of datapoint
16     # K = ... # number of mixtures
17     # Gamma = ... # zeros of shape (N,K), matrix of responsibilities
18     # for n in range(N):
19     #     for k in range(K):
20     #         # Gamma[n, k] = ....
21     #     # Gamma[n, :] /= ... # Normalize by sum across second dimension (mixtures)
22     # return Gamma
23     N, D = data.shape # Number of datapoints and dimension of datapoint
24     K = Mu.shape[1] # number of mixtures
25     Gamma = np.zeros((N,K)) # zeros of shape (N,K), matrix of responsibilities
26     for n in range(N):
27         for k in range(K):
28             Gamma[n, k] = Pi[k] * normal_density(data[n], Mu[:,k], Sigma[k])
29         Gamma[n, :] /= np.sum(Gamma[n, :]) # Normalize by sum across second dimension (mixtures)
30     return Gamma
```

```
1 # TODO: Gaussian Mixture Maximization Step
2 def gm_m_step(data, Gamma):
3     """ Gaussian Mixture Maximization Step.
4
5     Args:
6         data: a NxD matrix for the data points
7         Gamma: a NxK matrix of responsibilities
8
9     Returns:
10         Mu: a DxK matrix for the means of the K Gaussian Mixtures
11         Sigma: a list of size K with each element being DxD covariance matrix
12         Pi: a vector of size K for the mixing coefficients
13     """
14     N, D = data.shape # Number of datapoints and dimension of datapoint
15     K = Gamma.shape[1]  # number of mixtures
16     Nk = np.sum(Gamma, axis=0) # Sum along first axis
17     Mu = np.zeros([D, K])
18     Sigma = np.zeros((K,D,D))
19
20     for k in range(K):
21         for n in range(N):
22             Mu[:,k] += Gamma[n,k] * data[n,:]/Nk[k]
23     for k in range(K):
24         for n in range(N):
25             temp = (data[n,:] - Mu[:,k]).reshape((D,1))
26             Sigma[k] += Gamma[n,k] * np.dot(temp, temp.T)/Nk[k]
27     Pi = Nk / N
28     return Mu, Sigma, Pi
```

```
1 # TODO: Run this cell to call the Gaussian Mixture EM algorithm
2 N, D = data.shape
3 K = 2
4 Mu = np.zeros([D, K])
```

```
 5 Mu[:, 1] = 1.
 6 Sigma = [np.eye(2), np.eye(2)]
 7 Pi = np.ones(K) / K
 8 Gamma = np.zeros([N, K]) # Gamma is the matrix of responsibilities
 9
10 max_iter  = 200
11
12 iterations = []
13 costs = []
14 for it in range(max_iter):
15     Gamma = gm_e_step(data, Mu, Sigma, Pi)
16     Mu, Sigma, Pi = gm_m_step(data, Gamma)
17     loglike = log_likelihood(data, Mu, Sigma, Pi)
18     iterations.append(it)
19     costs.append(loglike)
20     print(it, loglike) # This function makes the computation longer, but good for debugging
21
22 class_1 = np.where(Gamma[:, 0] >= .5)
23 class_2 = np.where(Gamma[:, 1] >= .5)
```

⤷

```
0  -2108.569030648646
1  -2103.523067787013
2  -2101.578628940565
3  -2100.39571336378
4  -2099.4236632574143
5  -2098.464055686756
6  -2097.408395697052
7  -2096.1536982077464
8  -2094.5612603236423
9  -2092.519545762133
10  -2090.247433322861
11  -2088.2949298098933
12  -2086.9395731095738
13  -2086.1168033743156
14  -2085.660531138905
15  -2085.421344453549
16  -2085.298761138278
17  -2085.235414098758
18  -2085.2015105004816
19  -2085.1822768473003
20  -2085.1704904611533
21  -2085.1626047237182
22  -2085.1568538546535
23  -2085.1523420208
24  -2085.148604067255
25  -2085.1453911308645
26  -2085.1425645465742
27  -2085.140042670383
28  -2085.137773932699
29  -2085.13572305107
30  -2085.1338639070595
31  -2085.1321758226945
32  -2085.130641581603
33  -2085.1292463507466
34  -2085.1279770703873
35  -2085.1268220897928
36  -2085.1257709341003
37  -2085.1248141428205
38  -2085.1239431490208
39  -2085.1231501829784
40  -2085.122428191501
41  -2085.121770768269
42  -2085.1211720923784
43  -2085.1206268735164
44  -2085.1201303027024
45  -2085.1196780078126
46  -2085.1192660133675
47  -2085.1188907041155
48  -2085.1185487919975
49  -2085.1182372862168
50  -2085.117953466067
51  -2085.117694856276
52  -2085.117459204643
53  -2085.117244461728
54  -2085.117048762411
55  -2085.116870409161
56  -2085.116707856822
57  -2085.1165596988235
58  -2085.1164246546223
59  -2085.1163015583297
60  -2085.116189348338
61  -2085.1160870579333
62  -2085.1159938067362
63  -2085.1159087929304
64  -2085.115831286205
65  -2085.115760621329
66  -2085.115696192309
67  -2085.1156374470756
68  -2085.115583882667
69  -2085.1155350408185
70  -2085.11549050397
71  -2085.115449891633
72  -2085.115412857078
73  -2085.1153790843246
74  -2085.115348285397
75  -2085.115320197835
76  -2085.115294582412
77  -2085.1152712210865
78  -2085.1152499150935
79  -2085.1152304832535
80  -2085.1152127604037
81  -2085.115196595972
82  -2085.1151818526905
83  -2085.115168405407
84  -2085.115156140023
85  -2085.1151449524973
86  -2085.1151347479804
87  -2085.115125439978
88  -2085.1151169496216
89  -2085.115109205
90  -2085.115102140534
91  -2085.115095696413
92  -2085.1150898181154
93  -2085.1150844558997
94  -2085.115079564413
95  -2085.1150751022924
```

```
 96 -2085.1150710318134
 97 -2085.1150673185743
 98 -2085.1150639311895
 99 -2085.115060841046
100 -2085.115058022042
101 -2085.1150554503643
102 -2085.1150531043013
103 -2085.1150509640465
104 -2085.1150490115356
105 -2085.1150472302825
106 -2085.1150456052655
107 -2085.1150441227674
108 -2085.1150427702855
109 -2085.1150415364123
110 -2085.1150404107357
111 -2085.115039383769
112 -2085.1150384468465
113 -2085.115037592078
114 -2085.1150368122526
115 -2085.1150361007944
116 -2085.11503545171
117 -2085.1150348595315
118 -2085.115034319264
119 -2085.115033826356
120 -2085.1150333766577
121 -2085.115032966381
122 -2085.1150325920626
123 -2085.1150322505555
124 -2085.115031938981
125 -2085.115031654714
126 -2085.115031395365
127 -2085.115031158744
128 -2085.115030942858
129 -2085.115030745898
130 -2085.115030566197
131 -2085.115030402244
132 -2085.11503025266
133 -2085.115030116183
134 -2085.1150299916653
135 -2085.115029878061
136 -2085.1150297744116
137 -2085.1150296798455
138 -2085.115029593564
139 -2085.1150295148445
140 -2085.115029443021
141 -2085.115029377494
142 -2085.115029317707
143 -2085.1150292631587
144 -2085.11502921339
145 -2085.11502916798
146 -2085.115029126552
147 -2085.1150290887545
148 -2085.1150290542655
149 -2085.1150290228
150 -2085.1150289940915
151 -2085.115028967899
152 -2085.115028943999
153 -2085.115028922195
154 -2085.1150289023008
155 -2085.1150288841523
156 -2085.115028867589
157 -2085.1150288524777
158 -2085.1150288386934
159 -2085.1150288261147
160 -2085.1150288146396
161 -2085.115028804169
162 -2085.115028794612
163 -2085.1150287858973
164 -2085.115028777943
165 -2085.1150287706896
166 -2085.115028764068
167 -2085.1150287580294
168 -2085.1150287525174
169 -2085.1150287474907
170 -2085.1150287429027
171 -2085.115028738715
172 -2085.1150287348955
173 -2085.115028731412
174 -2085.1150287282317
175 -2085.115028725332
176 -2085.1150287226865
177 -2085.115028720271
178 -2085.115028718067
179 -2085.115028716056
180 -2085.115028714223
181 -2085.115028712549
182 -2085.115028711023
183 -2085.1150287096307
184 -2085.115028708359
185 -2085.115028707199
186 -2085.115028706141
187 -2085.1150287051746
188 -2085.1150287042956
189 -2085.1150287034893
190 -2085.1150287027563
191 -2085.11150287020887
```

```
192 -2085.115028701477
193 -2085.115028700922
194 -2085.115028700415
195 -2085.115028699949
196 -2085.1150286995266
197 -2085.1150286991387
198 -2085.1150286987895
199 -2085.115028698466
```
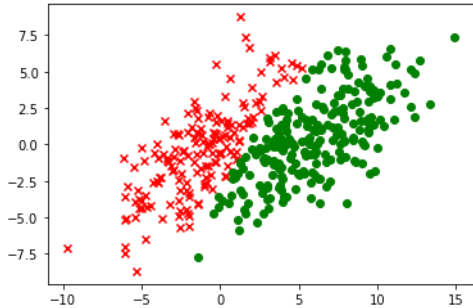
```
1 # TODO: Make a scatterplot for the data points showing the Gaussian Mixture cluster assignments of each point
2
3 plt.scatter(data[class_1][:, 0], data[class_1][:, 1], marker='x', color='red')
4 plt.scatter(data[class_2][:, 0], data[class_2][:, 1], marker='o', color = 'green')
```
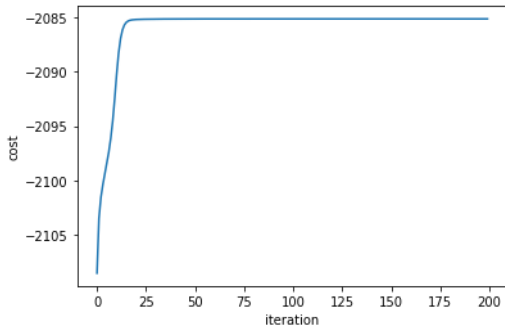
<matplotlib.collections.PathCollection at 0x7f513ea6bc50>



```
1 plt.ylabel('cost')
2 plt.xlabel('iteration')
3 plt.plot(iterations, costs)
```

[<matplotlib.lines.Line2D at 0x7f5141c02a90>]



```
1 labels_predict_EM = np.argwhere(Gamma >= 0.5)[:, 1]
2 print("Error rate: ", np.mean(labels != labels_predict_EM))
```

Error rate:  0.11

```
1 k_means_error_list = []
2 EM_error_list = []
3 k_means_time_list = []
4 EM_time_list = []
5
6 for i in range(5):
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    N, D = data.shape
22    K = 2
23    max_iter = 200
24    class_init = np.random.binomial(1., .5, size=N)
25    R = np.vstack([class_init, 1 - class_init]).T
26
27    start_time=time.clock()
28    Mu = np.zeros([D, K])
29    Mu[:, 1] = 1.
30    R.T.dot(data), np.sum(R, axis=0)
31    for it in range(max_iter):
```

```
32          R = km_assignment_step(data, Mu)
33          Mu = km_refitting_step(data, R, Mu)
34          c = cost(data, R, Mu)
35      labels_predict = np.argwhere(R == 1)[:, 1]
36      k_means_error_list.append(np.mean(labels != labels_predict))
37      time_use=time.clock()-start_time
38      k_means_time_list.append(time_use)
39
40      start_time = time.clock()
41      Sigma = [np.eye(2), np.eye(2)]
42      Pi = np.ones(K) / K
43      Gamma = np.zeros([N, K]) # Gamma is the matrix of responsibilities
44      for it in range(max_iter):
45          Gamma = gm_e_step(data, Mu, Sigma, Pi)
46          Mu, Sigma, Pi = gm_m_step(data, Gamma)
47          loglike = log_likelihood(data, Mu, Sigma, Pi)
48      labels_predict = np.argwhere(Gamma >= 0.5)[:, 1]
49      EM_error_list.append(np.mean(labels != labels_predict))
50      time_use=time.clock()-start_time
51      EM_time_list.append(time_use)
```

```
1 print("BELOW IS THE RESULT FOR FIVE DIFFERENT DATA REALIZATION")
2 print("==============================================================")
3 print("THE K_MEANS REUSLT IS:", k_means_error_list)
4 print("==============================================================")
5 print("THE K_MEANS TIME IS:", k_means_time_list)
6 print("==============================================================")
7 print("THE EM_ERROR REUSLT IS:", EM_error_list)
8 print("==============================================================")
9 print("THE EM_ERROR REUSLT IS:", EM_time_list)
```

```
⤷  BELOW IS THE RESULT FOR FIVE DIFFERENT DATA REALIZATION
    ==============================================================
    THE K_MEANS REUSLT IS: [0.2325, 0.2675, 0.2675, 0.2675, 0.2375]
    ==============================================================
    THE K_MEANS TIME IS: [0.1449140000000284, 0.1250660000000039, 0.12583499999999503, 0.12067200000001321, 0.12221499999998287]
    ==============================================================
    THE EM_ERROR REUSLT IS: [0.065, 0.0875, 0.1025, 0.1075, 0.0975]
    ==============================================================
    THE EM_ERROR REUSLT IS: [11.571901000000025, 11.817532999999969, 11.66621299999997, 11.625350999999966, 11.75290799999999]
```

## 4. Comment on findings + additional experiments

Comment on the results:

- Compare the performance of k-Means and EM based on the resulting cluster assignments. **Answer:** Compare the performance of k-Means and EM besed on resulting cluster assignments, we could see that they are different. Also, notice the misclassification rate for k-Means is significantly larger than the misclassification rate for EM, we could draw the conclusion that EM has a better clusterring result than k-Means.

- Compare the performance of k-Means and EM based on their convergence rate. What is the bottleneck for which method? **Answer:** From the result above we could see that for k-Means the convergenece time is between 0.1 point to 0.2 point. For EM method, the convergence time is roughly over 10 seconds. Therefore, k-Means converges way faster then EM, and the bottleneck for EM method is EM is more accurate, but takes longer computataion time, and for k-Means the bottleneck is it indeed takes relatively shorter time but the prediction is not as accurate as EM's.

- Experiment with 5 different data realizations (generate new data), run your algorithms, and summarize your findings. Does the algorithm performance depend on different realizations of data? **Answer:** From

## 2. Reinforcement Learning

There are 3 files:

1. `maze.py`: defines the `MazeEnv` class, the simulation environment which the Q-learning agent will interact in.
2. `qlearning.py`: defines the `qlearn` function which you will implement, along with several helper functions. Follow the instructions in the file.
3. `plotting_utils.py`: defines several plotting and visualization utilities. In particular, you will use `plot_steps_vs_iters`, `plot_several_steps_vs_iters`, `plot_policy_from_q`

```
1 # from qlearning import qlearn
2 # from maze import MazeEnv, ProbabilisticMazeEnv
3 # from plotting_utils import plot_steps_vs_iters, plot_several_steps_vs_iters, plot_policy_from_q
```

```
1 import numpy as np
2 import matplotlib
3 import matplotlib.pyplot as plt
4 # from qlearning import *
5 # from maze import *
6
7 #  UTILITY FUNCTIONS
```

```
 8
 9
10 color_cycle = ['#377eb8', '#ff7f00', '#a65628',
11                '#f781bf','#4daf4a', '#984ea3',
12                '#999999', '#e41a1c', '#dede00']
13
14 def plot_steps_vs_iters(steps_vs_iters, block_size=10):
15     num_iters = len(steps_vs_iters)
16     block_size = 10
17     num_blocks = num_iters // block_size
18     smooted_data = np.zeros(shape=(num_blocks, 1))
19     for i in range(num_blocks):
20         lower = i * block_size
21         upper = lower + 9
22         smooted_data[i] = np.mean(steps_vs_iters[lower:upper])
23
24     plt.figure()
25     plt.title("Steps to goal vs episodes")
26     plt.ylabel("Steps to goal")
27     plt.xlabel("Episodes")
28     plt.plot(np.arange(1,num_iters,block_size), smooted_data, color=color_cycle[0])
29
30     return
31
32 def plot_several_steps_vs_iters(steps_vs_iters_list, label_list, block_size=10):
33     smooted_data_list = []
34     for steps_vs_iters in steps_vs_iters_list:
35         num_iters = len(steps_vs_iters)
36         block_size = 10
37         num_blocks = num_iters // block_size
38         smooted_data = np.zeros(shape=(num_blocks, 1))
39         for i in range(num_blocks):
40             lower = i * block_size
41             upper = lower + 9
42             smooted_data[i] = np.mean(steps_vs_iters[lower:upper])
43         smooted_data_list.append(smooted_data)
44
45     plt.figure()
46     plt.title("Steps to goal vs episodes")
47     plt.ylabel("Steps to goal")
48     plt.xlabel("Episodes")
49     index = 0
50     for label, smooted_data in zip(label_list, smooted_data_list):
51         plt.plot(np.arange(1,num_iters,block_size), smooted_data, label=label, color=color_cycle[index])
52         index += 1
53     plt.legend()
54
55     return
56
57
58 # this function sets color values for
59 # Q table cells depending on expected reward value
60 def get_color(value, min_val, max_val):
61
62     switcher={
63                 0:'gray',
64                 1:'indigo',
65                 2:'darkmagenta',
66                 3:'orchid',
67                 4:'lightpink',
68             }
69
70     step = (max_val-min_val)/5
71     i = 0
72     color='lightpink'
73
74     for limit in np.arange(min_val, max_val, step):
75         if limit <= value < limit+step:
76             color = switcher.get(i)
77         i+=1
78     return color
79
80
81
82 # get first cell out of the start state
83 def get_next_cell(x1,x2,heatmap,policy_table,xlim=9,ylim=9):
84     up_reward=-10000
85     down_reward=-10000
86     left_reward=-10000
87     right_reward=-10000
88
89     if (x1<ylim):
90         if (policy_table[x1-1][x2]!=3):
91             up_reward = heatmap[x1-1][x2]
92     else:
93         up_reward = -1000
```

```
 94
 95        if (x1>0):
 96            if (policy_table[x1+1][x2]!=0):
 97                down_reward = heatmap[x1+1][x2]
 98        else:
 99            down_reward = -1000
100
101        if (x2>0):
102            if (policy_table[x1][x2-1]!=1):
103                left_reward = heatmap[x1][x2-1]
104
105        else:
106            left_reward = -1000
107
108        if (x2<xlim):
109            if (policy_table[x1][x2+1]!=2):
110                right_reward = heatmap[x1][x2+1]
111
112        else:
113            right_reward = -1000
114
115        rewards = np.array([up_reward, down_reward, left_reward, right_reward])
116        idx = np.argmax(rewards)
117        next_cell = [(x1-1,x2), (x1+1,x2), (x1,x2-1), (x1,x2+1)][idx]
118        choice = ['up', 'down', 'left', 'right']
119        #print ('picking ',choice[idx])
120        return next_cell
121
122
123
124
125 # get coordinates of the cells
126 # on the way from the start to goal state
127 def get_path(x1,x2, policy_table):
128        x_coords = [x1]
129        y_coords = [x2]
130        x1_new = x1
131        x2_new = x2
132
133        i=0
134        num_steps = 0
135        total_cells = len(policy_table)*len(policy_table[0])
136        while (policy_table[x1][x2]!='G') and num_steps < total_cells:
137            if (policy_table[x1][x2]==1): # right
138                x2_new=x2+1
139                #print(i, ' - moving right')
140
141            elif (policy_table[x1][x2]==0):
142                x1_new=x1-1
143                #print(i, ' - moving up')
144
145            elif (policy_table[x1][x2]==3):
146                x1_new=x1+1
147                #print(i, ' - moving down')
148
149            elif (policy_table[x1][x2]==2):
150                x2_new=x2-1
151                #print(i, ' - moving left')
152
153            x1 = x1_new
154            x2 = x2_new
155            x_coords.append(x1)
156            y_coords.append(x2)
157            num_steps += 1
158        return x_coords, y_coords
159
160 # plot Q table
161 # optimal path is highlighted and cells colored by their values
162 def plot_table(env, table_data, heatmap, goal_states, start_state, max_val, min_val, x_coords, y_coords):
163        fig = plt.figure(dpi=80)
164        ax = fig.add_subplot(1,1,1)
165        plt.figure(figsize=(10,10))
166
167        width = len(table_data[0])
168        height = len(table_data)
169
170        new_table = []
171
172        for i in range(height):
173            new_row = []
174
175            for j in range(width):
176                if env.map[i][j] == 0:
177                    new_row.append('')
178                else:
179                    digit = table_data[i][j]
```

```
180                     if (digit==0):
181                         new_row.append('\u2191') # up
182                     elif (digit==1):
183                         new_row.append('\u2192') # right
184                     elif (digit==2):
185                         new_row.append('\u2190') # left
186                     elif (digit==3):
187                         new_row.append('\u2193') # down
188                     elif (digit=='G'):
189                         new_row.append('G') # goal state
190                     elif (digit=='S'):
191                         new_row.append('S') # goal state
192                     elif (digit==-1):
193                         new_row.append('+') # All four directions
194                     else:
195                         new_row.append('x') # unknown
196
197             new_table.append(new_row)
198
199         table = ax.table(cellText=new_table, loc='center',cellLoc='center')
200
201         table.scale(1,2)
202
203         for i in range(height):
204             new_row = []
205
206             for j in range(width):
207                 if new_table[i][j] == '':
208                     table[i, j].set_facecolor('black')
209                 else:
210                     table[i, j].set_facecolor(get_color(heatmap[i][j],min_val,max_val))
211
212         for goal_state in goal_states:
213             table[(goal_state[0], goal_state[1])].set_facecolor("limegreen")
214         table[(start_state[0], start_state[1])].set_facecolor("yellow")
215         ax.axis('off')
216         table.set_fontsize(16)
217
218         for i in range(len(x_coords)):
219             table[(x_coords[i], y_coords[i])].get_text().set_color('red')
220         plt.show()
221
222
223 # this function takes 3D Q table as an input
224 # and outputs optimal trajectory table (policy table)
225 # and corresponding excpected reward values of different cells (heatmap)
226 def get_policy_table(q_hat_3D, start_state, goal_states):
227     policy_table = []
228     heatmap = []
229
230     for i in range(q_hat_3D.shape[0]):
231         row = []
232         heatmap_row = []
233         for j in range(q_hat_3D.shape[1]):
234
235             heatmap_row.append(np.max(q_hat_3D[i,j,:]))
236
237             for goal_state in goal_states:
238                 if (goal_state[0]==i) and (goal_state[1]==j):
239                     row.append('G')
240
241             if (start_state[0]==i) and (start_state[1]==j):
242                 row.append('S')
243             else:
244                 if np.max(q_hat_3D[i,j,:]) == 0:
245                     row.append(-1) # All zeros
246                 else:
247                     row.append(np.argmax(q_hat_3D[i,j,:]))
248         policy_table.append(row)
249         heatmap.append(heatmap_row)
250
251     return policy_table, heatmap
252
253 def plot_policy_from_q(q_hat, env):
254     q_hat_3D = np.reshape(q_hat, (env.m_size, env.m_size, env.num_actions))
255     max_val = q_hat_3D.max()
256     min_val = q_hat_3D.min()
257     start_state = env.get_coords_from_state(env._get_start_state)
258     goal_states = env._get_goal_state
259     goal_states = [env.get_coords_from_state(goal_state) for goal_state in goal_states]
260     policy_table, heatmap = get_policy_table(q_hat_3D, start_state, goal_states)
261     x,y = get_next_cell(start_state[0],start_state[1],heatmap,policy_table)
262     x_coords, y_coords = get_path(x,y,policy_table)
263     plot_table(env, policy_table, heatmap, goal_states, start_state,max_val,min_val, x_coords, y_coords)
264
265     return
```

```python
 1  import numpy as np
 2  import copy
 3  import math
 4  import random
 5
 6  ACTION_MEANING = {
 7      0: "UP",
 8      1: "RIGHT",
 9      2: "LEFT",
10      3: "DOWN",
11  }
12
13  SPACE_MEANING = {
14      1: "ROAD",
15      0: "BARRIER",
16      -1: "GOAL",
17  }
18
19
20  class MazeEnv:
21
22      def __init__(self, start=[6,3], goals=[[1, 8]]):
23          """Deterministic Maze Environment"""
24
25          self.m_size = 10
26          self.reward = 10
27          self.num_actions = 4
28          self.num_states = self.m_size * self.m_size
29
30          self.map = np.ones((self.m_size, self.m_size))
31          self.map[3, 4:9] = 0
32          self.map[4:8, 4] = 0
33          self.map[5, 2:4] = 0
34
35          for goal in goals:
36              self.map[goal[0], goal[1]] = -1
37
38          self.start = start
39          self.goals = goals
40          self.obs = self.start
41
42      def step(self, a):
43          """ Perform a action on the environment
44
45              Args:
46                  a (int): action integer
47
48              Returns:
49                  obs (list): observation list
50                  reward (int): reward for such action
51                  done (int): whether the goal is reached
52          """
53          done, reward = False, 0.0
54          next_obs = copy.copy(self.obs)
55
56          if a == 0:
57              next_obs[0] = next_obs[0] - 1
58          elif a == 1:
59              next_obs[1] = next_obs[1] + 1
60          elif a == 2:
61              next_obs[1] = next_obs[1] - 1
62          elif a == 3:
63              next_obs[0] = next_obs[0] + 1
64          else:
65              raise Exception("Action is Not Valid")
66
67          if self.is_valid_obs(next_obs):
68              self.obs = next_obs
69
70          if self.map[self.obs[0], self.obs[1]] == -1:
71              reward = self.reward
72              done = True
73
74          state = self.get_state_from_coords(self.obs[0], self.obs[1])
75
76          return state, reward, done
77
78      def is_valid_obs(self, obs):
79          """ Check whether the observation is valid
80
81              Args:
82                  obs (list): observation [x, y]
83
84              Returns:
85                  is_valid (bool)
```

```
 86          """
 87
 88          if obs[0] >= self.m_size or obs[0] < 0:
 89              return False
 90
 91          if obs[1] >= self.m_size or obs[1] < 0:
 92              return False
 93
 94          if self.map[obs[0], obs[1]] == 0:
 95              return False
 96
 97          return True
 98
 99      @property
100      def _get_obs(self):
101          """ Get current observation
102          """
103          return self.obs
104
105      @property
106      def _get_state(self):
107          """ Get current observation
108          """
109          return self.get_state_from_coords(self.obs[0], self.obs[1])
110
111      @property
112      def _get_start_state(self):
113          """ Get the start state
114          """
115          return self.get_state_from_coords(self.start[0], self.start[1])
116
117      @property
118      def _get_goal_state(self):
119          """ Get the start state
120          """
121          goals = []
122          for goal in self.goals:
123              goals.append(self.get_state_from_coords(goal[0], goal[1]))
124          return goals
125
126      def reset(self):
127          """ Reset the observation into starting point
128          """
129          self.obs = self.start
130          state = self.get_state_from_coords(self.obs[0], self.obs[1])
131          return state
132
133      def get_state_from_coords(self, row, col):
134          state = row * self.m_size + col
135          return state
136
137      def get_coords_from_state(self, state):
138          row = int(math.floor(state/self.m_size))
139          col = int(state % self.m_size)
140          return row, col
141
142
143 class ProbabilisticMazeEnv(MazeEnv):
144      """ (Q2.3) Hints: you can refer the implementation in MazeEnv
145      """
146
147      def __init__(self, goals=[[2, 8]], p_random=0.05):
148          """ Probabilistic Maze Environment
149
150              Args:
151                  goals (list): list of goals coordinates
152                  p_random (float): random action rate
153          """
154          self.m_size = 10
155          self.reward = 10
156          self.num_actions = 4
157          self.num_states = self.m_size * self.m_size
158          self.map = np.ones((self.m_size, self.m_size))
159          self.map[3, 4:9] = 0
160          self.map[4:8, 4] = 0
161          self.map[5, 2:4] = 0
162          for goal in goals:
163              self.map[goal[0], goal[1]] = -1
164          self.goals = goals
165          self.start = [6,3]
166          self.obs = self.start
167          self.p_random = p_random
168
169      def step(self, a):
170          done, reward = False, 0.0
171          next_obs = copy.copy(self.obs)
```

```
172
173          rand  = random.uniform(0, 1)
174          if rand <= self.p_random:
175              a = np.random.choice(4)
176
177          if a == 0:
178              next_obs[0] = next_obs[0] - 1
179          elif a == 1:
180              next_obs[1] = next_obs[1] + 1
181          elif a == 2:
182              next_obs[1] = next_obs[1] - 1
183          elif a == 3:
184              next_obs[0] = next_obs[0] + 1
185          else:
186              raise Exception("Action is Not Valid")
187
188          if self.is_valid_obs(next_obs):
189              self.obs = next_obs
190
191          if self.map[self.obs[0], self.obs[1]] == -1:
192              reward = self.reward
193              done = True
194
195          state = self.get_state_from_coords(self.obs[0], self.obs[1])
196
197          return state, reward, done
198
199
```

```
 1 import numpy as np
 2 import math
 3 import copy
 4
 5 def qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy, init_beta=None, k_exp_sched=None):
 6     """ Runs tabular Q learning algorithm for stochastic environment.
 7
 8     Args:
 9         env: instance of environment object
10         num_iters (int): Number of episodes to run Q-learning algorithm
11         alpha (float): The learning rate between [0,1]
12         gamma (float): Discount factor, between [0,1]
13         epsilon (float): Probability in [0,1] that the agent selects a random move instead of
14             selecting greedily from Q value
15         max_steps (int): Maximum number of steps in the environment per episode
16         use_softmax_policy (bool): Whether to use softmax policy (True) or Epsilon-Greedy (False)
17         init_beta (float): If using stochastic policy, sets the initial beta as the parameter for the softmax
18         k_exp_sched (float): If using stochastic policy, sets hyperparameter for exponential schedule
19             on beta
20
21     Returns:
22         q_hat: A Q-value table shaped [num_states, num_actions] for environment with with num_states
23             number of states (e.g. num rows * num columns for grid) and num_actions number of possible
24             actions (e.g. 4 actions up/down/left/right)
25         steps_vs_iters: An array of size num_iters. Each element denotes the number
26             of steps in the environment that the agent took to get to the goal
27             (capped to max_steps)
28     """
29     action_space_size = env.num_actions
30     state_space_size = env.num_states
31     q_hat = np.zeros(shape=(state_space_size, action_space_size))
32     steps_vs_iters = np.zeros(num_iters)
33
34     for i in range(num_iters):
35         # TODO: Initialize current state by resetting the environment
36         # curr_state = ...
37         curr_state = env.reset()
38         num_steps = 0
39         done = False
40
41         # TODO: Keep looping while environment isn't done and less than maximum steps
42         # while ...:
43         while (done != True) 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
52                 beta = beta_exp_schedule(init_beta, i, k_exp_sched)
53                 action = softmax_policy(q_hat, beta, curr_state, action_space_size)
54             else:
55                 # TODO: Epsilon-greedy
56                 action = epsilon_greedy(q_hat, epsilon, curr_state, action_space_size)
57             next_state, reward, done = env.step(action)
```

```
 58
 59                     # TODO: Update Q_value
 60                     if next_state != curr_state:
 61                         new_value = alpha * (reward + gamma * np.max(q_hat[next_state, :]) - q_hat[curr_state, action])
 62                         q_hat[curr_state, action] = q_hat[curr_state, action] + new_value
 63                         curr_state = next_state
 64             steps_vs_iters[i] = num_steps
 65
 66     return q_hat, steps_vs_iters
 67
 68
 69 def epsilon_greedy(q_hat, epsilon, state, action_space_size):
 70     """ Chooses a random action with p_rand_move probability,
 71     otherwise choose the action with highest Q value for
 72     current observation
 73
 74     Args:
 75         q_hat_3D: A Q-value table shaped [num_rows, num_col, num_actions] for
 76             grid environment with num_rows rows and num_col columns and num_actions
 77             number of possible actions
 78         epsilon (float): Probability in [0,1] that the agent selects a random
 79             move instead of selecting greedily from Q value
 80         obs: A 2-element array with integer element denoting the row and column
 81             that the agent is in
 82         action_space_size (int): number of possible actions
 83
 84     Returns:
 85         action (int): A number in the range [0, action_space_size-1]
 86             denoting the action the agent will take
 87     """
 88     # TODO: Implement your code here
 89     # Hint: Sample from a uniform distribution and check if the sample is below
 90     # a certain threshold
 91     prob = np.random.uniform()
 92     if np.all(q_hat[state, :] == 0) == True:
 93         action = np.random.choice(action_space_size)
 94     elif (prob <= epsilon):
 95         action = np.random.choice(action_space_size)
 96     else:
 97         action = q_hat[state,:].argmax()
 98     return action
 99
100
101 def softmax_policy(q_hat, beta, state, action_space_size):
102     """ Choose action using policy derived from Q, using
103     softmax of the Q values divided by the temperature.
104
105     Args:
106         q_hat: A Q-value table shaped [num_rows, num_col, num_actions] for
107             grid environment with num_rows rows and num_col columns
108         beta (float): Parameter for controlling the stochasticity of the action
109         obs: A 2-element array with integer element denoting the row and column
110             that the agent is in
111
112     Returns:
113         action (int): A number in the range [0, action_space_size-1]
114             denoting the action the agent will take
115     """
116     prob = stable_softmax(beta * q_hat)
117     if all(ele == 0 for ele in q_hat[state]):
118         action = np.random.choice(action_space_size)
119     else:
120         action = np.random.choice(action_space_size, p=prob[state,:])
121     return action
122
123
124 def beta_exp_schedule(init_beta, iteration, k=0.1):
125     beta = init_beta * np.exp(k * iteration)
126     return beta
127
128 def stable_softmax(x, axis=1):
129     """ Numerically stable softmax:
130     softmax(x) = e^x /(sum(e^x))
131                = e^x / (e^max(x) * sum(e^x/e^max(x)))
132
133     Args:
134         x: An N-dimensional array of floats
135         axis: The axis for normalizing over.
136
137     Returns:
138         output: softmax(x) along the specified dimension
139     """
140     max_x = np.max(x, axis, keepdims=True)
141     z = np.exp(x - max_x)
142     output = z / np.sum(z, axis, keepdims=True)
143
```

```
144     return output
```

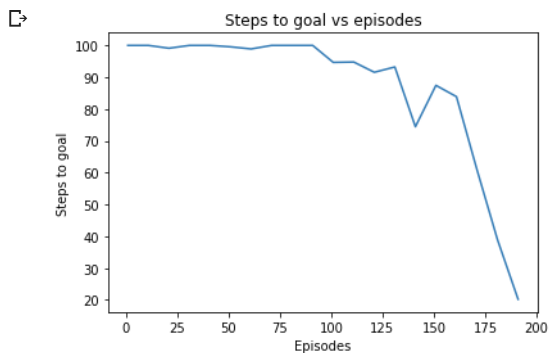# ▾ 1. Basic Q Learning experiments

(a) Run your algorithm several times on the given environment. Use the following hyperparameters:

1. Number of episodes = 200
2. Alpha ($\alpha$) learning rate = 1.0
3. Maximum number of steps per episode = 100. An episode ends when the agent reaches a goal state, or uses the maximum number of steps per episode
4. Gamma ($\gamma$) discount factor = 0.9
5. Epsilon ($\epsilon$) for $\epsilon$-greedy = 0.1 (10% of the time)

```
1 # TODO: Fill this in
2 num_iters = 200
3 alpha = 1.0
4 gamma = 0.9
5 epsilon = 0.1
6 max_steps = 100
7 use_softmax_policy = False
8
9 # TODO: Instantiate the MazeEnv environment with default arguments
10 env = MazeEnv()
11
12 # TODO: Run Q-learning:
13 q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy)
```
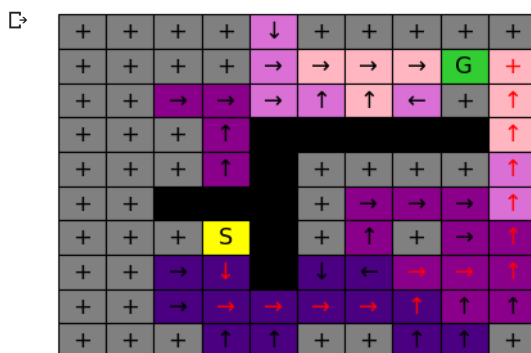
Plot the steps to goal vs training iterations (episodes):

```
1 # TODO: Plot the steps vs iterations
2 plot_steps_vs_iters(steps_vs_iters)
```



Visualize the learned greedy policy from the Q values:

```
1 # TODO: plot the policy from the Q value
2 plot_policy_from_q(q_hat, env)
```



```
<Figure size 720x720 with 0 Axes>
```

(b) Run your algorithm by passing in a list of 2 goal locations: (1,8) and (5,6). Note: we are using 0-indexing, where (0,0) is top left corner. Report on the results.

```
1 # TODO: Fill this in (same as before)
2 num_iters = 200
3 alpha = 1.0
4 gamma = 0.9
5 epsilon = 0.1
6 max_steps = 100
7 use_softmax_policy = False
```

8

```
 9 # TODO: Set the goal
10 goal_locs = [[1, 8], [5, 6]]
11 env = MazeEnv(goals= goal_locs)
12
13 # TODO: Run Q-learning:
14 q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy)
```
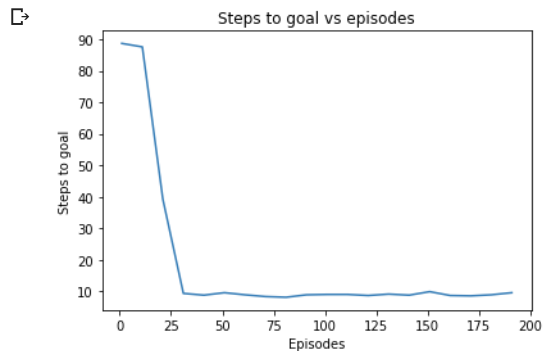
Plot the steps to goal vs training iterations (episodes):

```
1 # TODO: Plot the steps vs iterations
2 plot_steps_vs_iters(steps_vs_iters)
```
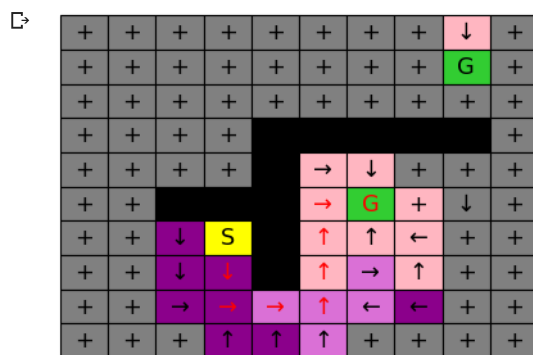


Plot the steps to goal vs training iterations (episodes):

```
1 # TODO: plot the policy from the Q values
2 plot_policy_from_q(q_hat, env)
```



```
<Figure size 720x720 with 0 Axes>
```

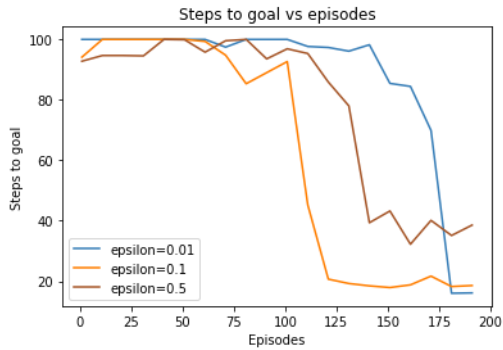## ▼ 2. Experiment with the exploration strategy, in the original environment

(a) Try different $\epsilon$ values in $\epsilon$-greedy exploration: We asked you to use a rate of $\epsilon$=10%, but try also 50% and 1%. Graph the results (for 3 epsilon values) and discuss the costs and benefits of higher and lower exploration rates.

```
 1 # TODO: Fill this in (same as before)
 2 num_iters = 200
 3 alpha = 1.0
 4 gamma = 0.9
 5 max_steps = 100
 6 use_softmax_policy = False
 7
 8 # TODO: set the epsilon lists in increasing order:
 9 epsilon_list = [0.01, 0.1, 0.5]
10
11 env = MazeEnv()
12
13 steps_vs_iters_list = []
14 for epsilon in epsilon_list:
15     q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy)
16     steps_vs_iters_list.append(steps_vs_iters)
```

```
1 # TODO: Plot the results
2 label_list = ["epsilon={}".format(eps) for eps in epsilon_list]
3 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)
```
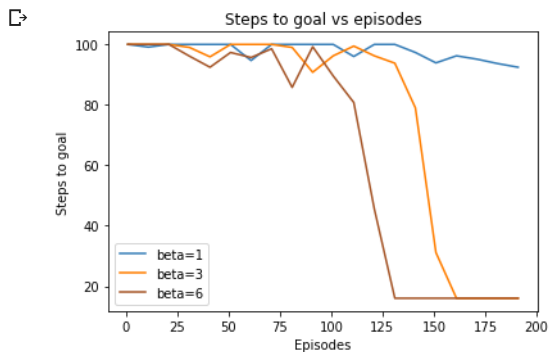
(b) Try exploring with policy derived from **softmax of Q-values** described in the Q learning lecture. Use the values of $\beta \in \{1, 3, 6\}$ for your experiment, keeping $\beta$ fixed throughout the training.

```
1 # TODO: Fill this in for Static Beta with softmax of Q-values
2 num_iters = 200
3 alpha = 1.0
4 gamma = 0.9
5 epsilon = 0.1
6 max_steps = 100
7
8 # TODO: Set the beta
9 beta_list = [1, 3, 6]
10 use_softmax_policy = True
11 k_exp_schedule = 0.0
12 env = MazeEnv()
13 steps_vs_iters_list = []
14 for beta in beta_list:
15     q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy, beta, k_exp_schedule
16     steps_vs_iters_list.append(steps_vs_iters)
```

```
1 label_list = ["beta={}".format(beta) for beta in beta_list]
2 # TODO:
3 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)
```



(c) Instead of fixing the $\beta = \beta_0$ to the initial value, we will increase the value of $\beta$ as the number of episodes $t$ increase:
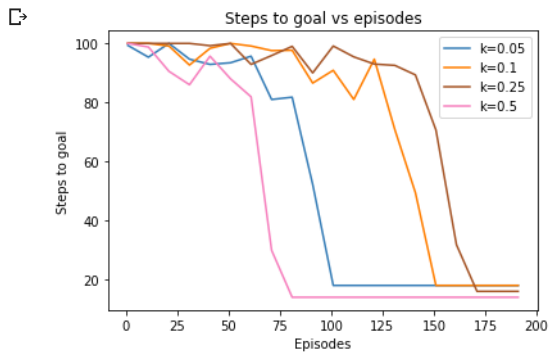
$$\beta(t) = \beta_0 e^{kt}$$

That is, the $\beta$ value is fixed for a particular episode. Run the training again for different values of $k \in \{0.05, 0.1, 0.25, 0.5\}$, keeping $\beta_0 = 1.0$. Compare the results obtained with this approach to those obtained with a static $\beta$ value.

```
1 # TODO: Fill this in for Dynamic Beta
2 num_iters = 200
3 alpha = 1.0
4 gamma = 0.9
5 epsilon = 0.1
6 max_steps = 100
7
8 # TODO: Set the beta
9 beta = 1.0
10 use_softmax_policy = True
11 k_exp_schedule_list = [0.05, 0.1, 0.25, 0.5]
12 env = MazeEnv()
13
14 steps_vs_iters_list = []
15 for k_exp_schedule in k_exp_schedule_list:
16     q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy, beta, k_exp_schedule
17     steps_vs_iters_list.append(steps_vs_iters)
```

```
1 # TODO: Plot the steps vs iterations
```

```
2 label_list = ["k={}".format(k_exp_schedule) for k_exp_schedule in k_exp_schedule_list]
3 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)
```
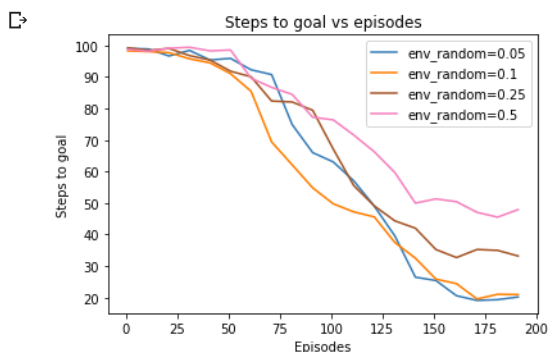


## 3. Stochastic Environments

(a) Make the environment stochastic (uncertain), such that the agent only has a 95% chance of moving in the chosen direction, and has a 5% chance of moving in some random direction.

```
1 # TODO: Implement ProbabilisticMazeEnv in maze.py
```

(b) Change the learning rule to handle the non-determinism, and experiment with different probability of environment performing random action $p_{rand} \in \{0.05, 0.1, 0.25, 0.5\}$ in this new rule. How does performance vary as the environment becomes more stochastic?

Use the same parameters as in first part, except change the alpha ($\alpha$) value to be **less than 1**, e.g. 0.5.

```
 1 # TODO: Fill this in for Dynamic Beta
 2 num_iters = 200
 3 alpha = 0.5
 4 gamma = 0.9
 5 epsilon = 0.1
 6 max_steps = 100
 7 use_softmax_policy = False
 8 beta = 1.0
 9
10 # Set the environment probability of random
11 env_p_rand_list = [0.05, 0.1, 0.25, 0.5]
12
13 steps_vs_iters_list = []
14 for env_p_rand in env_p_rand_list:
15     # Instantiate with ProbabilisticMazeEnv
16     env = ProbabilisticMazeEnv(p_random=env_p_rand)
17
18     # Note: We will repeat for several runs of the algorithm to make the result less noisy
19     avg_steps_vs_iters = np.zeros(num_iters)
20     for i in range(10):
21         # q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy, init_beta=6, k
22         q_hat, steps_vs_iters = qlearn(env, num_iters, alpha, gamma, epsilon, max_steps, use_softmax_policy)
23         avg_steps_vs_iters += steps_vs_iters
24     avg_steps_vs_iters /= 10
25     steps_vs_iters_list.append(avg_steps_vs_iters)
```

```
1 label_list = ["env_random={}".format(env_p_rand) for env_p_rand in env_p_rand_list]
2 plot_several_steps_vs_iters(steps_vs_iters_list, label_list)
```



## 3. Did you complete the course evaluation?

**ABSOLUTELY YES** Thank you so much my dearest professors and TAs.