

Implementing EM for Gaussian mixtures

In this assignment you will

- implement the EM algorithm for a Gaussian mixture model
- apply your implementation to cluster images
- explore clustering results and interpret the output of the EM algorithm

If you are using GraphLab Create

An IPython Notebook has been provided below to you for this assignment. This notebook contains the instructions, quiz questions and partially-completed code for you to use as well as some cells to test your code.

- Download the image dataset in SFrame format: [images.sf.zip](#)
- Download the companion IPython notebook: [3_em-for-gmm_blank.ipynb](#)
- Download the sample picture: [chosen_images.png](#)
- Save all the files in the same directory (where you are calling IPython notebook from) and unzip the data file.

Open the companion IPython notebook and follow the instructions in the notebook. The instructions below do not apply to users of GraphLab Create.

If you are not using GraphLab Create

It is possible to complete this assignment without using GraphLab Create. The instructions below are geared towards Python users, but you are free to adapt them to your specific environment.

Disclaimer. We have tested the assessment using the standard Python installation (with access to scikit-learn). However, the assessment may not be compatible with other tools (e.g. Matlab, R).

Download the dataset

- Download the image dataset in SFrame format: [images.sf.zip](#) (Those experimenting with other tools, get [images.zip](#) instead.)

Import packages

```
1 import sframe  
    instruction  
2 import numpy as np  
3 import matplotlib.pyplot as plt  
4 from scipy.stats import multivariate_normal  
    distribution  
5 import copy  
6  
7 %matplotlib inline
```

see below for install
dense matrices
plotting
multivariate Gaussian
deep copies

About SFrame. SFrame is a dataframe library Dato has released free-of-charge. Its source code is available [here](#). You may install SFrame via pip:

```
1 pip install --upgrade sframe
```

Implementing the EM algorithm for Gaussian mixture models

In this section, you will implement the EM algorithm. We will take the following steps:

- Create some synthetic data.
- Provide a log likelihood function for this model.
- Implement the EM algorithm.
- Visualize the progress of the parameters during the course of running EM.
- Visualize the convergence of the model.

Synthetic dataset. To help us develop and test our implementation, we will generate some observations from a mixture of Gaussians and then run our EM algorithm to discover the mixture components. We'll begin with a function to generate the data, and a quick plot to visualize its output for a 2-dimensional mixture of three Gaussians.

Now we will create a function to generate data from a mixture of Gaussians model.

```
1 def generate_MoG_data(num_data, means, covariances, weights):
2     """ Creates a list of data points """
3     num_clusters = len(weights)
4     data = []
5     for i in range(num_data):
6         # Use np.random.choice and weights to pick a cluster id greater than or equal to 0 and
7         # less than num_clusters.
8         k = np.random.choice(len(weights), 1, p=weights)[0]
9         # Use np.random.multivariate_normal to create data from this cluster
10        x = np.random.multivariate_normal(means[k], covariances[k])
11
12        data.append(x)
13    return data
```

After specifying a particular set of clusters (so that the results are reproducible across assignments), we use the above function to generate a dataset.

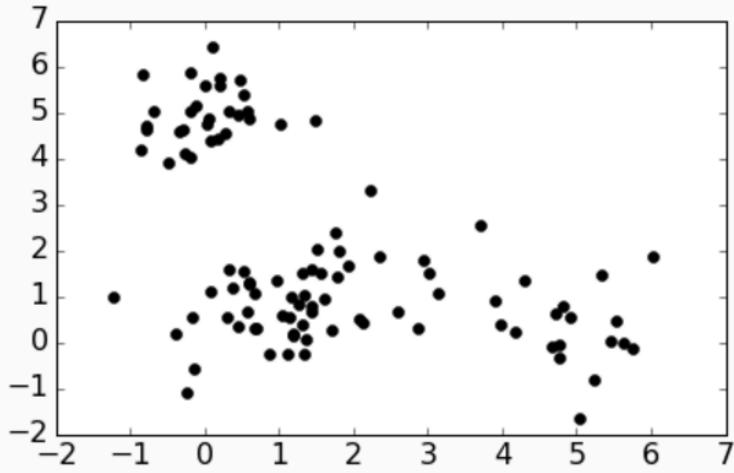
```
1 # Model parameters
2 init_means = [
3     [5, 0], # mean of cluster 1
4     [1, 1], # mean of cluster 2
5     [0, 5] # mean of cluster 3
6 ]
7 init_covariances = [
8     [[.5, 0.], [0, .5]], # covariance of cluster 1
9     [[1., .7], [0, .7]], # covariance of cluster 2
10    [[.5, 0.], [0, .5]] # covariance of cluster 3
11 ]
12 init_weights = [1/4., 1/2., 1/4.] # weights of each cluster
13
14 # Generate data
15 np.random.seed(4)
16 data = generate_MoG_data(100, init_means, init_covariances, init_weights)
```

Checkpoint: To verify your implementation above, make sure the following code does not return an error.

```
1 assert len(data) == 100
2 assert len(data[0]) == 2
3 print 'Checkpoint passed!'
```

Now plot the data you created above. The plot should be a scatterplot with 100 points that appear to roughly fall into three clusters.

```
1 plt.figure()
2 d = np.vstack(data)
3 plt.plot(d[:,0], d[:,1], 'ko')
4 plt.rcParams.update({'font.size':16})
5 plt.tight_layout()
```



Log likelihood. We provide a function to calculate log likelihood for mixture of Gaussians. The log likelihood quantifies the probability of observing a given set of data under a particular setting of the parameters in our model. We will use this to assess convergence of our EM algorithm; specifically, we will keep looping through EM update steps until the log likelihood ceases to increase at a certain rate.

```
1 def log_sum_exp(Z):
2     """ Compute log(\sum_i exp(Z_i)) for some array Z."""
3     return np.max(Z) + np.log(np.sum(np.exp(Z - np.max(Z))))
4
5 def loglikelihood(data, weights, means, covs):
6     """ Compute the loglikelihood of the data for a Gaussian mixture model with the given
7         parameters. """
8     num_clusters = len(means)
9     num_dim = len(data[0])
10
11    ll = 0
12    for d in data:
13
14        Z = np.zeros(num_clusters)
15        for k in range(num_clusters):
16
17            # Compute (x-mu)^T * Sigma^{-1} * (x-mu)
18            delta = np.array(d) - means[k]
19            exponent_term = np.dot(delta.T, np.dot(np.linalg.inv(covs[k]), delta))
20
21            # Compute loglikelihood contribution for this data point and this cluster
22            Z[k] += np.log(weights[k])
23            Z[k] -= 1/2. * (num_dim * np.log(2*np.pi) + np.log(np.linalg.det(covs[k])) +
24                            exponent_term)
25
26        # Increment loglikelihood contribution of this data point across all clusters
27        ll += log_sum_exp(Z)
28
29    return ll
```

Implementation of EM. You will now complete an implementation that can run EM on the data you just created. It uses the loglikelihood function we provided above.

Fill in the places where you find ## YOUR CODE HERE. There are seven places in this function for you to fill in.

Hint: Some useful functions

- [multivariate_normal.pdf](#): lets you compute the likelihood of seeing a data point in a multivariate Gaussian distribution.
- [np.outer](#): comes in handy when estimating the covariance matrix from data.

```

1 def EM(data, init_means, init_covariances, init_weights, maxiter=1000, thresh=1e-4):
2
3     # Make copies of initial parameters, which we will update during each iteration
4     means = init_means[:]
5     covariances = init_covariances[:]
6     weights = init_weights[:]
7
8     # Infer dimensions of dataset and the number of clusters
9     num_data = len(data)
10    num_dim = len(data[0])
11    num_clusters = len(means)
12
13    # Initialize some useful variables
14    resp = np.zeros((num_data, num_clusters))
15    ll = loglikelihood(data, weights, means, covariances)
16    ll_trace = [ll]
17
18    for i in range(maxiter):
19        if i % 5 == 0:
20            print("Iteration %s" % i)
21
22        # E-step: compute responsibilities
23        # Update resp matrix so that resp[j, k] is the responsibility of cluster k for data
24        # point j.
25        # Hint: To compute likelihood of seeing data point j given cluster k, use
26        # multivariate_normal.pdf.
27        for j in range(num_data):
28            for k in range(num_clusters):
29                # YOUR CODE HERE
30                resp[j, k] = ...
31
32        row_sums = resp.sum(axis=1)[:, np.newaxis]
33        resp = resp / row_sums # normalize over all possible cluster assignments
34
35        # M-step
36        # Compute the total responsibility assigned to each cluster, which will be useful when
37        # implementing M-steps below. In the lectures this is called N^{soft}
38        counts = np.sum(resp, axis=0)

```

```

36
37     for k in range(num_clusters):
38
39         # Update the weight for cluster k using the M-step update rule for the cluster
40         # weight,  $\hat{\pi}_k$ .
41         # YOUR CODE HERE
42         weights[k] = ...
43
44         # Update means for cluster k using the M-step update rule for the mean variables.
45         # This will assign the variable means[k] to be our estimate for  $\hat{\mu}_k$ .
46         weighted_sum = 0
47         for j in range(num_data):
48             # YOUR CODE HERE
49             weighted_sum += ...
50         # YOUR CODE HERE
51         means[k] = ...
52
53         # Update covariances for cluster k using the M-step update rule for covariance
54         # variables.
55         # This will assign the variable covariances[k] to be the estimate for  $\hat{\Sigma}_k$ .
56         weighted_sum = np.zeros((num_dim, num_dim))
57         for j in range(num_data):
58             # YOUR CODE HERE (Hint: Use np.outer on the data[j] and this cluster's mean)
59             weighted_sum += ...
60         # YOUR CODE HERE
61         covariances[k] = ...
62
63         # Compute the loglikelihood at this iteration
64         # YOUR CODE HERE
65         ll_latest = ...
66         ll_trace.append(ll_latest)
67
68         # Check for convergence in log-likelihood and store
69         if (ll_latest - ll) < thresh and ll_latest > -np.inf:
70             break
71         ll = ll_latest
72
73         if i % 5 != 0:
74             print("Iteration %s" % i)
75
76         out = {'weights': weights, 'means': means, 'covs': covariances, 'loglik': ll_trace, 'resp':
77             resp}

```

Testing the implementation on the simulated data. Now we'll fit a mixture of Gaussians to this data using our implementation of the EM algorithm. As with k-means, it is important to ask how we obtain an initial configuration of mixing weights and component parameters. In this simple case, we'll take three random points to be the initial cluster means, use the empirical covariance of the data to be the initial covariance in each cluster (a clear overestimate), and set the initial mixing weights to be uniform across clusters.

```
1 np.random.seed(4)
2
3 # Initialization of parameters
4 chosen = np.random.choice(len(data), 3, replace=False)
5 initial_means = [data[x] for x in chosen]
6 initial_covs = [np.cov(data, rowvar=0)] * 3
7 initial_weights = [1/3.] * 3
8
9 # Run EM
10 results = EM(data, initial_means, initial_covs, initial_weights)
```

Checkpoint. For this particular example, the EM algorithm is expected to terminate in 30 iterations. That is, the last line of the log should say "Iteration 29". If your function stopped too early or too late, you should re-visit your code.

Our algorithm returns a dictionary with five elements:

- 'loglik': a record of the log likelihood at each iteration
- 'resp': the final responsibility matrix
- 'means': a list of K means
- 'covs': a list of K covariance matrices
- 'weights': the weights corresponding to each model component

Quiz Question: What is the weight that EM assigns to the first component after running the above code block?

Quiz Question: Using the same set of results, obtain the mean that EM assigns the second component. What is the mean in the first dimension?

Quiz Question: Using the same set of results, obtain the mean that EM assigns the second component. What is the mean in the first dimension?

Plot progress of parameters. One useful feature of testing our implementation on low-dimensional simulated data is that we can easily visualize the results.

We will use the following plot_contours function to visualize the Gaussian components over the data at three different points in the algorithm's execution:

1. At initialization (using initial_mu, initial_cov, and initial_weights)
2. After running the algorithm to completion
3. After just 12 iterations (using parameters estimates returned when setting maxiter=12)

```
1 import matplotlib.mlab as mlab
2 def plot_contours(data, means, covs, title):
3     plt.figure()
4     plt.plot([x[0] for x in data], [y[1] for y in data], 'ko') # data
5
6     delta = 0.025
7     k = len(means)
8     x = np.arange(-2.0, 7.0, delta)
9     y = np.arange(-2.0, 7.0, delta)
10    X, Y = np.meshgrid(x, y)
11    col = ['green', 'red', 'indigo']
12    for i in range(k):
13        mean = means[i]
14        cov = covs[i]
15        sigmax = np.sqrt(cov[0][0])
16        sigmay = np.sqrt(cov[1][1])
17        sigmaxy = cov[0][1]/(sigmax*sigmay)
18        Z = mlab.bivariate_normal(X, Y, sigmax, sigmay, mean[0], mean[1], sigmaxy)
19        plt.contour(X, Y, Z, colors = col[i])
20        plt.title(title)
21    plt.rcParams.update({'font.size':16})
22    plt.tight_layout()
```

If you are using other tools, look for library functions to visualize 2D Gaussian distributions.

Run the following code block to visualize the progress of EM at initialization, after 12 iterations, and after convergence.

```
1 # Parameters after initialization
2 plot_contours(data, initial_means, initial_covs, 'Initial clusters')
3
4 # Parameters after 12 iterations
5 results = ... # YOUR CODE HERE
6 plot_contours(data, results['means'], results['covs'], 'Clusters after 12 iterations')
7
8 # Parameters after running EM to convergence
9 results = EM(data, initial_means, initial_covs, initial_weights)
10 plot_contours(data, results['means'], results['covs'], 'Final clusters')
```

Quiz Question: Plot the loglikelihood that is observed at each iteration. Is the loglikelihood plot monotonically increasing, monotonically decreasing, or neither [multiple choice]? Complete the following code block to answer the question.

```
1 results = EM(data, initial_means, initial_covs, initial_weights)
2
3 # YOUR CODE HERE
4 loglikelihoods = ...
5
6 plt.plot(range(len(loglikelihoods)), loglikelihoods, linewidth=4)
7 plt.xlabel('Iteration')
8 plt.ylabel('Log-likelihood')
9 plt.rcParams.update({'font.size':16})
10 plt.tight_layout()
```

Fitting a Gaussian mixture model for image data

Now that we're confident in our implementation of the EM algorithm, we'll apply it to cluster some more interesting data. In particular, we have a set of images that come from four categories: sunsets, rivers, trees and forests, and cloudy skies. For each image we are given the average intensity of its red, green, and blue pixels, so we have a 3-dimensional representation of our data. Our goal is to find a good clustering of these images using our EM implementation; ideally our algorithm would find clusters that roughly correspond to the four image categories.

To begin with, we'll take a look at the data and get it in a form suitable for input to our algorithm. The data are provided in SFrame format:

```
1 images = sframe.SFrame('images.sf/')
2 images['rgb'] = images.pack_columns(['red', 'green', 'blue'])['X4']
```

(Optional) Generating dataset from original images. You are free to experiment with other tools to generate 3-dimensional representation of images yourself. For each images, you should compute the average pixel value of each color channel and then divide the average by 256:

```
1 # arr = 3D array of pixels, with dimensions (width, height, 3)
2
3 r, g, b = np.mean(arr[:, :, 0]/256.0), np.mean(arr[:, :, 1]/256.0), np.mean(arr[:, :, 2]/256.0)IniIni
```

Initialization. We need to come up with initial estimates for the mixture weights and component parameters. Let's take three images to be our initial cluster centers, and let's initialize the covariance matrix of each cluster to be diagonal with each element equal to the sample variance from the full data. As in our test on simulated data, we'll start by assuming each mixture component has equal weight.

This may take a few minutes to run.

```
1 np.random.seed(1)
2
3 # Initialize parameters
4 init_means = [images['rgb'][x] for x in np.random.choice(len(images), 4, replace=False)]
5 cov = np.diag([images['red'].var(), images['green'].var(), images['blue'].var()])
6 init_covariances = [cov, cov, cov, cov]
7 init_weights = [1/4., 1/4., 1/4., 1/4.]
8
9 # Convert rgb data to numpy arrays
10 img_data = [np.array(i) for i in images['rgb']]
11
12 # Run our EM algorithm on the image data using the above initializations.
13 # This should converge in about 125 iterations
14 out = EM(img_data, init_means, init_covariances, init_weights)
```

The following sections will evaluate the results by asking the following questions:

- **Convergence:** How did the log likelihood change across iterations? Did the algorithm achieve convergence?
- **Uncertainty:** How did cluster assignment and uncertainty evolve?
- **Interpretability:** Can we view some example images from each cluster? Do these clusters correspond to known image categories?

Evaluating convergence. Let's start by plotting the log likelihood at each iteration - we know that the EM algorithm guarantees that the log likelihood can only increase (or stay the same) after each iteration, so if our implementation is correct then we should see an increasing function.

```
1 ll = out['loglik']
2 plt.plot(range(len(ll)),ll,linewidth=4)
3 plt.xlabel('Iteration')
4 plt.ylabel('Log-likelihood')
5 plt.rcParams.update({'font.size':16})
6 plt.tight_layout()
```

The log likelihood increases so quickly on the first few iterations that we can barely see the plotted line. Let's plot the log likelihood after the first three iterations to get a clearer view of what's going on:

```
1 plt.figure()
2 plt.plot(range(3,len(ll)),ll[3:],linewidth=4)
3 plt.xlabel('Iteration')
4 plt.ylabel('Log-likelihood')
5 plt.rcParams.update({'font.size':16})
6 plt.tight_layout()
```

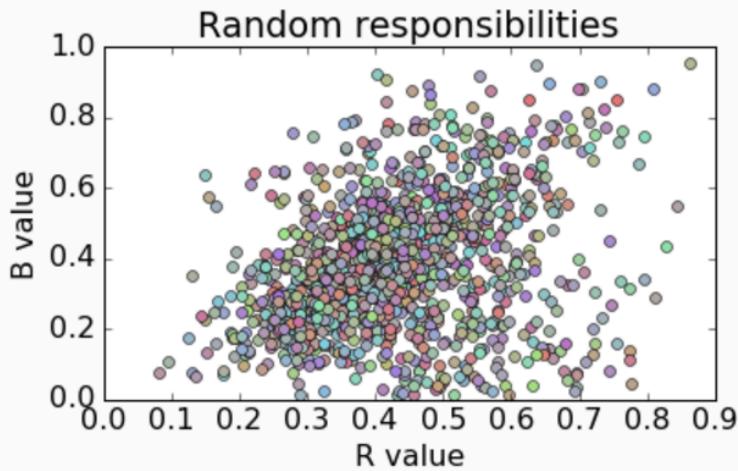
Evaluating uncertainty. Next we'll explore the evolution of cluster assignment and uncertainty. Remember that the EM algorithm represents uncertainty about the cluster assignment of each data point through the responsibility matrix. Rather than making a 'hard' assignment of each data point to a single cluster, the algorithm computes the responsibility of each cluster for each data point, where the responsibility corresponds to our certainty that the observation came from that cluster.

We can track the evolution of the responsibilities across iterations to see how these 'soft' cluster assignments change as the algorithm fits the Gaussian mixture model to the data; one good way to do this is to plot the data and color each point according to its cluster responsibilities. Our data are three-dimensional, which can make visualization difficult, so to make things easier we will plot the data using only two dimensions, taking just the [R G], [G B] or [R B] values instead of the full [R G B] measurement for each observation.

```
1 import colorsys
2 def plot_responsibilities_in_RB(img, resp, title):
3     N, K = resp.shape
4
5     HSV_tuples = [(x*1.0/K, 0.5, 0.9) for x in range(K)]
6     RGB_tuples = map(lambda x: colorsys.hsv_to_rgb(*x), HSV_tuples)
7
8     R = img['red']
9     B = img['blue']
10    resp_by_img_int = [[resp[n][k] for k in range(K)] for n in range(N)]
11    cols = [tuple(np.dot(resp_by_img_int[n], np.array(RGB_tuples))) for n in range(N)]
12
13    plt.figure()
14    for n in range(len(R)):
15        plt.plot(R[n], B[n], 'o', c=cols[n])
16    plt.title(title)
17    plt.xlabel('R value')
18    plt.ylabel('B value')
19    plt.rcParams.update({'font.size':16})
20    plt.tight_layout()
```

To begin, we will visualize what happens when each data has random responsibilities.

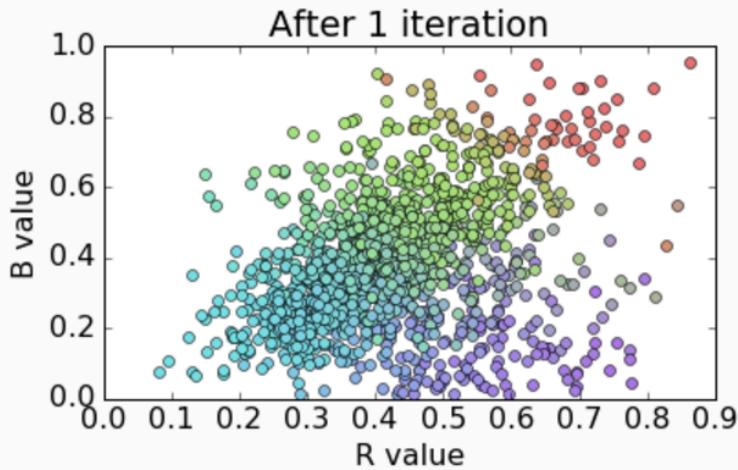
```
1 N, K = out['resp'].shape
2 random_resp = np.random.dirichlet(np.ones(K), N)
3 plot_responsibilities_in_RB(images, random_resp, 'Random responsibilities')
```



We now use the above plotting function to visualize the responsibilities after 1 iteration.

```

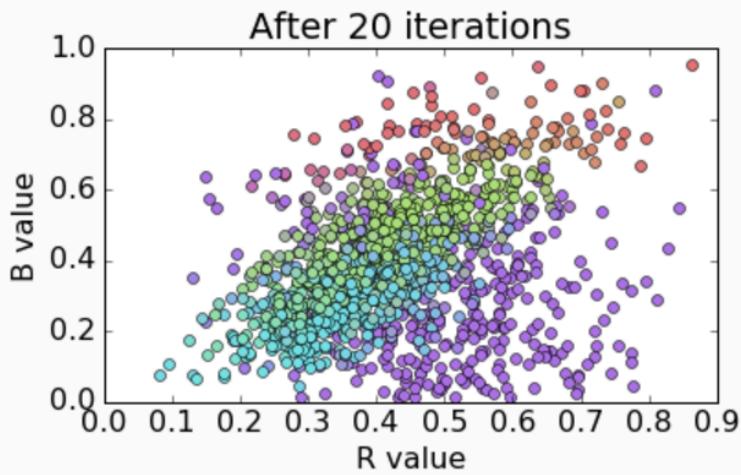
1 N, K = out['resp'].shape
2 random_resp = np.random.dirichlet(np.ones(K), N)
3 plot_responsibilities_in_RB(images, random_resp, 'Random responsibilities')
```



We now use the above plotting function to visualize the responsibilities after 1 iteration.

```

1 out = EM(img_data, init_means, init_covariances, init_weights, maxiter=1)
2 plot_responsibilities_in_RB(images, out['resp'], 'After 1 iteration')
```



Plotting the responsibilities over time in [R B] space shows a meaningful change in cluster assignments over the course of the algorithm's execution. While the clusters look significantly better organized at the end of the algorithm than they did at the start, it appears from our plot that they are still not very well separated. We note that this is due in part our decision to plot 3D data in a 2D space; everything that was separated along the G axis is now "squashed" down onto the flat [R B] plane. If we were to plot the data in full [R G B] space, then we would expect to see further separation of the final clusters. We'll explore the cluster interpretability more in the next section.

Interpreting each cluster

Let's dig into the clusters obtained from our EM implementation. Recall that our goal in this section is to cluster images based on their RGB values. We can evaluate the quality of our clustering by taking a look at a few images that 'belong' to each cluster. We hope to find that the clusters discovered by our EM algorithm correspond to different image categories - in this case, we know that our images came from four categories ('cloudy sky', 'rivers', 'sunsets', and 'trees and forests'), so we would expect to find that each component of our fitted mixture model roughly corresponds to one of these categories.

If we want to examine some example images from each cluster, we first need to consider how we can determine cluster assignments of the images from our algorithm output. This was easy with k-means - every data point had a 'hard' assignment to a single cluster, and all we had to do was find the cluster center closest to the data point of interest. Here, our clusters are described by probability distributions (specifically, Gaussians) rather than single points, and our model maintains some uncertainty about the cluster assignment of each observation.

One way to phrase the question of cluster assignment for mixture models is as follows: how do we calculate the distance of a point from a distribution? Note that simple Euclidean distance might not be appropriate since (non-scaled) Euclidean distance doesn't take direction into account. For example, if a Gaussian mixture component is very stretched in one direction but narrow in another, then a data point one unit away along the 'stretched' dimension has much higher probability (and so would be thought of as closer) than a data point one unit away along the 'narrow' dimension.

In fact, the correct distance metric to use in this case is known as Mahalanobis distance. For a Gaussian distribution, this distance is proportional to the square root of the negative log likelihood. This makes sense intuitively - reducing the Mahalanobis distance of an observation from a cluster is equivalent to increasing that observation's probability according to the Gaussian that is used to represent the cluster. This also means that we can find the cluster assignment of an observation by taking the Gaussian component for which that observation scores highest. We'll use this fact to find the top examples that are 'closest' to each cluster.

Quiz Question: Calculate the likelihood (score) of the first image in our data set (`images[0]`) under each Gaussian component through a call to `multivariate_normal.pdf`. Given these values, what cluster assignment should we make for this image?

Now we calculate cluster assignments for the entire image dataset using the result of running EM for 20 iterations above:

```
1 means = out['means']
2 covariances = out['covs']
3 rgb = images['rgb']
4 N = len(images)
5 K = len(means)
6
7 assignments = [0]*N
8 probs = [0]*N
9
10 for i in range(N):
11     # Compute the score of data point i under each Gaussian component:
12     p = np.zeros(K)
13     for k in range(K):
14         # YOUR CODE HERE (Hint: use multivariate_normal.pdf and rgb[i])
15         p[k] = ...
16
17     # Compute assignments of each data point to a given cluster based on the above scores:
18     # YOUR CODE HERE
19     assignments[i] = ...
20
21     # For data point i, store the corresponding score under this cluster assignment:
22     # YOUR CODE HERE
23     probs[i] = ...
24
25 assignments = sframe.SFrame({'assignments':assignments, 'probs':probs, 'image': images['image']})
```

We'll use the 'assignments' SFrame to find the top images from each cluster by sorting the datapoints within each cluster by their score under that cluster (stored in probs). We can plot the corresponding images in the original data using show().

Create a function that returns the top 5 images assigned to a given category in our data (HINT: use the SFrame function topk(column, k) to find the k top values according to specified column in an SFrame).

```
1 def get_top_images(assignments, cluster, k=5):
2     # YOUR CODE HERE
3     images_in_cluster = ...
4     top_images = images_in_cluster.topk('probs', k)
5     return top_images['image']
```

Here are some utility function to display the top images.

- IPython notebook users: use `display_images()`.
- Others: use `save_images()`. This will save the images instead of displaying them on the screen.

```
1 def display_images(images):
2     from IPython.display import display
3
4     for image in images:
5         display(Image.open(BytesIO(image._image_data)))
6
7 def save_images(images, prefix):
8     for i, image in enumerate(images):
9         Image.open(BytesIO(image._image_data)).save(prefix % i)
```

Run the following to display the top images.

```
1 for component_id in range(4):
2     print 'Component {0:d}'.format(component_id)
3     images = get_top_images(assignments, component_id)
4     display_images(images)
5     #save_images(images, 'component_{0:d}_%d.jpg'.format(component_id))
6     print '\n'
```

These look pretty good! Our algorithm seems to have done a good job overall at 'discovering' the four categories that from which our image data was drawn. It seems to have had the most difficulty in distinguishing between rivers and cloudy skies, probably due to the similar color profiles of images in these categories; if we wanted to achieve better performance on distinguishing between these categories, we might need a richer representation of our data than simply the average [R G B] values for each image.

Quiz Question: Which of the following images are *not* in the list of top 5 images in the first cluster?

Image 1



Image 2



Image 3



Image 4



Image 5



Image 6



Image 7

