**CPSC-340    Machine Learning and Data Mining    2005**

**Homework # 6**
Due Tuesday, Nov 15 at 12 noon.

NAME:_____

Signature:_____

STD. NUM: _____

Email: _____

---

**General guidelines for homeworks:**

You are encouraged to discuss the problems with others in the class, but all write-ups are to be done on your own.
**Homework grades will be based not only on getting the "correct answer," but also on good writing style and clear presentation of your solution.** It is your responsibility to make sure that the graders can easily follow your line of reasoning.
Try every problem. Even if you can't solve the problem, you will receive partial credit for explaining why you got stuck on a promising line of attack. More importantly, you will get valuable feedback that will help you learn the material.
**Please acknowledge the people with whom you discussed the problems and what sources you used to help you solve the problem (e.g. websites, books from the library).** This won't affect your grade but is important as academic honesty.
**When dealing with Matlab exercises, please attach a printout with all your code and show your results clearly.**

1. **Question 1: Marginalization and Conditioning**: A band called Radiohead is inspired by an old band called The Beatles. 50% of music critics think the beatles was a great ($G$) band, 40% that it was moderate ($M$) and 10% that it was awful ($A$). These critics have also compiled the following table:

$$
\begin{array}{cc}
 & \begin{array}{ccc} & B_2 & \\ G & M & A \end{array} \\
\begin{array}{cc} & G \\ B_1 & M \\ & A \end{array} &
\left( \begin{array}{ccc}
0.8 & 0.1 & 0.1 \\
0.1 & 0.9 & 0 \\
0.2 & 0.3 & 0.5
\end{array} \right)
\end{array}
$$

The table says that the probability of a new band ($B_2$) being great given that the inspiring band ($B_1$) was great is $P(B_2 = G | B_1 = G) = 0.8$. Similarly, $P(B_2 = G | B_1 = M) = 0.1$, $P(B_2 = M | B_1 = A) = 0.3$, and so on. Note that the numbers in the rows add up to 1, so the table is a probability transition matrix.

(i) Given what the critics think of the Beatles and the fact that the Beatles inspired Radiohead, what is the probability that Radiohead is a great band?

(ii) What is $P(B_1 = G | B_2 = G)$?

2. **Question 2: K-means and EM**: The file `clusterData.dat` on the homework website contains 200 rows of $(x_1, x_2)$ pairs.

(i) Plot this data. The plots in the following parts should be plotted on top of this plot.

(ii) Implement the K-means algorithm and apply it to the data, using 3 clusters and iterating until the algorithm converges. For initialisation, use the cluster centroids $\mu_1 = (0.5, 2.3)$, $\mu_2 = (0.5, 0.5)$ and $\mu_3 = (5.2, -0.1)$. Plot the evolution of the cluster centroids as the algorithm runs (that is, plot $\mu_c^{(t)}$ for each $c$ and iteration $t$).

(iii) Implement the EM algorithm for fitting a mixture of Gaussians,

$$p(x|\theta) = \sum_{c=1}^{3} p(c) \mathcal{N}(x|\mu_c, \Sigma_c),$$

on the same dataset. For initialisation, use the same means as in (ii) and set $p(c) = (1/3, 1/3, 1/3)$ and $\Sigma_1 = \Sigma_2 = \Sigma_3 = I$. Plot the evolution of the means and (optional) superimpose the elipses for the variances around the respective means.

(iv) Comment on the relative rates at which the two algorithms converge.

3. **Question 3: Vector quantisation and image compression**: In this problem, we will apply the K-means algorithm to lossy image compression, by reducing the number of colors used in an image.

The data website contains a 512x512 image of a mandrill represented in 24-bit colour. This means that, for each 262144 pixels in the image, there are three 8-bit numbers (each ranging from 0 to 255) that represent the green, red and blue intensity values for that pixel. The straightforward representation of this image therefore takes about 262144x3 = 786432 bytes (a byte being 8 bits). To compress the image, we will use K-means to reduce the image to K=16 colours. More specifically, each pixel in the image is considered a point in the three-dimensional (r,g,b)-space. To compress the image, we will cluster these points in colour-space into 16 clusters, and replace each pixel with the closest cluster centroid.

(i) Load and plot the image `mandrill-large.tiff` using the following Matlab commands:

```
A = double(imread('mandrill-large.tiff'));
imshow(uint8(round(A)));
```

The first line creates a three dimensional matrix, such that `A(:,:,1),A(:,:,2)` and `A(:,:,3)` are 512x512 arrays that respectively contain the red, green and blue values for each pixel. Since this image is large, K-means would take too long. Instead, you should load and cluster the image `mandrill-small.tiff`. In particular, treat each pixel's (r,g,b) values as an element of $\mathbb{R}^3$ and run K-means with 16 clusters on this image. Iterate to (preferably) convergence, but in any case for less than 30 iterations. For initialisation, set each cluster centroid to the (r,g,b) values of a randomly chosen pixel in the image.

(ii) Take the matrix $A$ from `mandril-large.tiff`, and replace each pixel's (r,g,b) values with the value of the closest cluster centroid. Display the new image, and compare it visually to the original image. Hand in all your code and a printout of your compressed image next to the original image (printing on a black-and-white printer is fine).