**CSE 527: Computational Biology**                                    **Assignment #3**

Turn this one in on paper; handwritten is fine, I don't recommend trying to typeset it. Extra credit is for extra practice and glory; it is not a big component of your grade.

1. **Bayes Rule:** In a certain population, an obese person has a 30 percent chance of having high blood pressure and a non-obese person has a 10 percent chance of having high blood pressure. Twenty percent of the population is obese. What is the conditional probability that a person is obese, given that the person has high blood pressure?

2. **Maximum Likelihood:** Let $x_1, x_2, \ldots, x_n$ be $n$ samples of a normal random variable $X$ with mean $\theta_1$ and variance $\theta_2$. In class I showed that the maximum likelihood estimates of $\theta_1$ and $\theta_2$ when both are unknown give a biased estimate of $\theta_2$. What is the MLE of $\theta_2 = \sigma^2$ if $\theta_1 = \mu$ is assumed to be known?

   For example, suppose I draw a sample of 3 measuring 9, 10 , 11. The sample mean is 10, sample standard deviation is $\sqrt{((9-10)^2 + (10-10)^2 + (11-10)^2)/3} = \sqrt{2/3} \approx .8$.

   On the other hand suppose I told you the population mean was 0. Drawing a sample of 9, 10, 11 now seems much less likely, but is certainly still possible, and is made more probable by increasing our estimate of the population variance (the *sample* variance is unchanged; it's defined in terms of the sample mean, not the population mean). So the question is: now that I know the population mean, what estimate of population variance make the data I just observed most likely?

   Extra Credit: Is your estimate of $\theta_2$ biased, i.e., does the expected value of $\hat{\theta}_2$ differ from $\theta_2$?

3. **EM:** In class, I sketched the EM algorithm for the two-component Gaussian Mixture Model only in the special case when both subpopulations were assumed to share the same variance and the mixing proportions ($\tau_1/\tau_2$) were assumed to be 50/50. Carry out the analysis for the general case where $\sigma_1^2, \sigma_2^2$ and $0 \leq \tau_1 \leq 1$ ($\tau_2 = 1 - \tau_1$) are arbitrary.

4. **Maximum Likelihood:** Suppose $X$ is a discrete random variable with three possible outcomes, say $A_1, A_2$ and $A_3$. Let $\theta = (p_1, p_2, p_3)$ be the probabilities of outcomes $A_1, A_2, A_3$, resp., (where $p_1 + p_2 + p_3 = 1$, of course). Suppose you have collected $n$ independent random samples $x_1, x_2, \ldots, x_n$ drawn from this distribution. Using the same basic approach as in the coin-flipping example in the class notes (Lecute 6, slide 8), show that the maximum likelihood estimators for the parameters $\theta$ are $\hat{\theta} = (n_1/n, n_2/n, n_3/n)$, where $n_i$ is the number of occurrences of outcome $A_i$ among $x_1, x_2, \ldots, x_n$. Hint: The three variables are coupled, since $p_3 = 1 - p_1 - p_2$. The algebra is mildly easier if you happen to remember Lagrange multipliers, but it's absolutely not essential; just substitute for $p_3$ using the identity before you differentiate. (FYI, this result generalizes to arbitrary multinomial distributions, not just 2 or 3 outcomes; see the slick proof in Chapter 11.)

5. **EM:** Recall that an *allele* of a gene is one variant of its DNA or protein sequence. Individuals generally carry two (possibly identical) alleles of each gene, one inherited from mother, one from father (genes on the X/Y chromosomes being exceptions). The ABO blood type gene has three common alleles in the human population: A, B and O. The blood type of an individual depends as follows on the pair of alleles that he or she has: type A if the pair is A/A or A/O; type B if the pair is B/B or B/O; type AB if the pair is A/B; type O if the pair is O/O. Let $p(A)$ be the fraction of A alleles in the population, $p(B)$, the fraction of B alleles and $p(O)$, the fraction of O alleles. These fractions are nonnegative and sum to 1. Under the standard

assumption in genetics of independent assortment, the probability that an individual has a given pair of alleles is the same as the probability of obtaining that pair in two random draws from the set of all alleles in the population: for example, the probability of the pair A/B is $2p(A)p(B)$. In a sample of 20 individuals, the first 9 have blood type A, the next 2 have blood type B, the next has blood type AB and the last 8 have blood type O. Derive the appropriate formulas needed to use the EM algorithm to determine the values of $p(A)$, $p(B)$ and $p(O)$ most likely to have given rise to this data. Then run the algorithm for a few iterations on the given data. Try it with a couple of very different starting estimates for the parameters. You may write a program to do the iteration, do it by hand, or give a spreadsheet with the relevant formulas and "fill down" a few rows to iterate. If you use a spreadsheet, turn in a printout of the formulas as well as the numbers; I think CONTROL-backquote causes Excel to show all formulas. Hint: The parameters are $p(A), p(B)$ and $p(O)$, the observed data are the blood types of the individuals and the hidden data are the pairs of alleles possessed by the individuals. The solution to problem 4 will help. Depending on how you set up the likelihood function, you might (or might not) need the multinomial distribution from pg 300 of the text.

(If you'd like info on the genetics of the ABO blood group system, the 1930 Nobel prize in Physiology or Medicine, have a look at Wikipedia `http://en.wikipedia.org/wiki/Abo_blood_group` or OMIM `http://www.ncbi.nlm.nih.gov/entrez/dispomim.cgi?id=110300`. In a nutshell, they are 3 alleles of a single gene on the ninth chromosome (9q34), which encodes a *glycosyltransferase*—an enzyme that modifies the carbohydrate content of the red blood cell antigens. The A and B alleles perform slightly (but immunologically significantly) different modifications; the O allele has a 1 base deletion, hence an altered reading frame, producing a very different protein with no apparent function at all, a so-called "null" allele, more or less explaining why the O allele is "silent." Aside from issues with blood transfusions, people with O blood type are apparently more susceptible to cholera. And, no, the "independent assortment" assumption for this gene is *not* well justified in the human population; prevalence is strongly dependent on geography. But we'll ignore that for this problem...)

## Extra Credit Problems:

6. **Maximum Likelihood:** Suppose $X$ is a random variable uniformly distributed between 0 and $\theta > 0$ for some unknown $\theta$. Based on a sample $x_1, x_2, \ldots, x_n$ of $X$, what is the maximum likelihood estimator of $\theta$? Is it biased?

7. **EM:** Generalize the EM algorithm from problem 3 to allow a fixed but arbitrary number $k \geq 1$ of components in the mixture, preferably allowing a choice of either a common variance $\sigma^2$ shared by all clusters, or a separate variance per cluster. Implement it and experiment with simulated data to see how well it recovers the parameters you used to generate the data. How quickly does the iteration converge? Does it ever seem to be converging to a local, not global, max? How well does it work with sparse data? Well-separated clusters? Highly overlapping clusters?

8. **Motif Finding:** Pick 10–20 genes from one prokaryotic organism, say *E. coli*, and run one of the motif finding tools we've discussed (MEME, Align-ace, ...) on the 200 base-pair region upstream of each. Does it find anything interesting, perhaps a TATA box? Repeat with 10–20 human genes. For a somewhat more ambitious exercise, try Footprinter on 10–20 orthologous genes.