

CSC2515 – Assignment #3

Due: Nov.25, 2pm at the **START** of class

Worth: 18%

Late assignments not accepted.

1 Mixture of Base Rates (7%)

In this question you will derive the EM algorithm for a very simple model called “mixture of base rates”, which is the unsupervised equivalent of Naive Bayes classification. We model some discrete data \mathbf{x} by assuming that there is an *unobserved* cluster variable k and that given this cluster variable, all the components x_i of \mathbf{x} are conditionally independent. The equation is:

$$p(\mathbf{x}) = \sum_k p(k) \prod_i p(x_i|k) = \sum_k \alpha_k \prod_i \prod_{j \in V_i} a_{ijk}^{[x_i=j]}$$

where α_k is the prior probability of cluster k , a_{ijk} are the model parameters ($\sum_j a_{ijk} = 1 \forall i, k$), V_i is the set of possible values for the discrete variable x_i (e.g. if x_i is binary then $V_i = \{0, 1\}$), and the indicator $[x_i = j]$ is 1 if x_i takes on value j and 0 otherwise.

- Write down the complete (joint) log likelihood $p(\mathbf{x}, k)$.
- Derive the E-step of the EM algorithm by writing down the expression for the posterior probability $p(k|\mathbf{x}^n)$ of cluster k given a particular datapoint \mathbf{x}^n . Hint: use Bayes’ rule.
- Write down the expected complete log likelihood using the posterior probability you just derived: $\sum_n \sum_k p(k|\mathbf{x}^n) \log p(\mathbf{x}^n, k)$
- Derive the M-step of the EM algorithm by taking the derivative of this expected complete log likelihood with respect to the parameters a_{ijk} . Don’t forget to enforce the normalization constraint that $\sum_j a_{ijk} = 1 \forall i, k$.
- Also find the updates for the priors α_k . Don’t forget the normalization constraint.
- Assuming you have computed $q_k = \log \alpha_k + \sum_i \sum_{j \in V_i} [x_i = j] \log a_{ijk}$, how would you compute $p(k|x)$ in a numerically stable way using `logsum`?
- Fill in the ??:
For continuous vector valued data \mathbf{x} , the analog of this model is a “?? of ?? with diagonal ??”.

2 Learning Mixtures of Gaussians (11%)

In this question you will implement the EM algorithm for training a mixture of full covariance Gaussians:

$$p(\mathbf{x}) = \sum_k p(z = k)p(\mathbf{x}|z = k) = \sum_k \alpha_k |2\pi\Sigma_k|^{-1/2} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^\top \Sigma_k^{-1}(\mathbf{x} - \mu_k) \right\}$$

2.1 What to do

- Using the data provided in `a3geyser.mat`, train mixtures with 1,2,3 and 6 components, starting at 1000 random initializations and training until convergence.

In case you are interested, the data represents intervals between eruptions and durations of eruptions of the Old Faithful Geyser at Yellowstone National Park, USA.

2.2 What to hand in

- The equation you are using to update the covariance matrices in the M-step (see below).
- A plot of the best fit parameters (highest likelihood) that you found for each number of components. The plot should show the training data as points, the means of the Gaussians as larger symbols of some kind, and the one-standard deviation ellipses of the Gaussians. (4 plots total, on a single page.)
- A histogram, for each number of components, of the final log likelihoods of the training data and of the number of iterations of EM it took for training to converge. Don't use too few or too many bins in your histograms. (8 histograms total, arranged in two columns of 4, on one page.)

2.3 Some hints

- You might find the the function `plotGauss.m` helpful in making the one-sigma Gaussian ellipse plots.
- You can assess convergence of EM by monitoring the log likelihood and stopping when the fractional change $(\text{newlik} - \text{oldlik}) / \text{abs}(\text{newlik})$ is less than one part in a million. Your log likelihood should never decrease!
- Try to get the EM algorithm working *without* the covariance updates at first (ie just update the means). Then when that is working, and your likelihood is always going up, try adding in the covariance updates.
- Initialize the model parameters by setting the means randomly in some sensible range, for example the covariance of the entire data set. Set the covariances to some small but not tiny values. Make sure there is enough randomness in your initialization that you fall into different local maxima.
- You might find it interesting to plot the parameters after every iteration of EM, so you can see how things are going. If your code issues a `drawnow` command in MATLAB after plotting the data and all the Gaussian ellipses this will update the plots on the screen. (If these ellipses don't fit the data at all, you may have an error somewhere in your code or in your covariance update equation!)
- Every once in a while, if you are unlucky, a Gaussian cluster will get trapped on a single datapoint and its variance will get very small. This can cause numerical problems, so watch out for it. If this happens, you can just restart at a different random initialization.

2.4 Reminders

Amongst other things, your code will need to:

- Compute the log probability of every datapoint under each mixture component:

$$\ell_{nk} = \log p(\mathbf{x}^n, z = k | \mu_k, \Sigma_k)$$

- Compute the log probability of every datapoint:

$$\ell_n = \log p(\mathbf{x}^n) = \log \sum_k p(\mathbf{x}^n, z = k | \mu_k, \Sigma_k)$$

NB: Make sure to do this in a numerically stable way. In particular, *do not* use $\ell_n = \log \sum_k \exp(\ell_{nk})$.
Hint: You may find the function `logsum.m` helpful.

- Compute the E-step “responsibilities”:

$$r_{nk} = p(z = k | \mathbf{x}^n) = \frac{p(\mathbf{x}^n, z = k)}{p(\mathbf{x}^n)} = \exp(\ell_{nk} - \ell_n)$$

(You probably don’t want to use $\exp(\ell_{nk}) / \exp(\ell_n)$. Think about why.)

- Update the mixing proportions in the M-step:

$$\alpha_k^{new} = \frac{\sum_n r_{nk}}{N}$$

- Update the means in the M-step:

$$\mu_k^{new} = \frac{\sum_n r_{nk} \mathbf{x}^n}{\sum_n r_{nk}}$$

- Update the covariance matrices Σ_k in the M-step.
(I’ve left this equation for you to write down from the class notes and the book.)
Use full (not diagonal or spherical) covariances.

$$\Sigma_k^{new} = ???$$

- At each iteration, you should also compute the total log probability of the training set:

$$\ell = \sum_n \log p(\mathbf{x}^n) = \sum_n \ell_n$$

and verify that this never goes down.