

# ST758, Homework 7

Due Nov 13, 2013

Consider again maximizing the Dirichlet-multinomial log-likelihood

$$\begin{aligned} L(\boldsymbol{\alpha}) &= \sum_{i=1}^n \ln \binom{|\mathbf{x}_i|}{\mathbf{x}_i} + \sum_{i=1}^n \sum_{j=1}^d \sum_{k=0}^{x_{ij}-1} \ln(\alpha_j + k) - \sum_{i=1}^n \sum_{k=0}^{|\mathbf{x}_i|-1} \ln(|\boldsymbol{\alpha}| + k) \\ &= \sum_{i=1}^n \ln \binom{|\mathbf{x}_i|}{\mathbf{x}_i} + \sum_{i=1}^n \sum_{j=1}^d [\ln \Gamma(\alpha_j + x_{ij}) - \ln \Gamma(\alpha_j)] - \sum_{i=1}^n [\ln \Gamma(|\boldsymbol{\alpha}| + |\mathbf{x}_i|) - \ln \Gamma(|\boldsymbol{\alpha}|)] \end{aligned}$$

given iid observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . In HW6, we worked out a Newton's method. In this homework, we explore the EM and MM approach.

1. Suppose  $(P_1, \dots, P_d) \in \Delta_d$  follows a Dirichlet distribution with parameter  $(\alpha_1, \dots, \alpha_d)$ . Show that

$$\mathbf{E}(\ln P_j) = \Psi(\alpha_j) - \Psi(|\boldsymbol{\alpha}|),$$

where  $\Psi(z) = \Gamma'(z)/\Gamma(z)$  is the digamma function and  $|\boldsymbol{\alpha}| = \sum_{j=1}^d \alpha_j$ . (Hint: Differentiate the identity  $1 = \int_{\Delta_d} \frac{\Gamma(|\boldsymbol{\alpha}|)}{\prod_{j=1}^d \Gamma(\alpha_j)} \prod_{j=1}^d p_j^{\alpha_j-1} d\mathbf{p}$ .)

2. The admixture representation of the Dirichlet-multinomial distribution suggests that we can treat the unobserved multinomial parameters  $\mathbf{p}_1, \dots, \mathbf{p}_n$  as missing data and derive an EM algorithm. Show that the Q function is

$$Q(\boldsymbol{\alpha}|\boldsymbol{\alpha}^{(t)}) = \sum_{j=1}^d \sum_{i=1}^n \alpha_j \left[ \Psi(x_{ij} + \alpha_j^{(t)}) - \Psi(|\mathbf{x}_i| + |\boldsymbol{\alpha}^{(t)}|) \right] - n \sum_{j=1}^d \ln \Gamma(\alpha_j) + n \ln \Gamma(|\boldsymbol{\alpha}|) + c^{(t)},$$

where  $c^{(t)}$  is a constant irrelevant to optimization. Comment on why it is not easy to maximize the Q function.

3. We derive an MM algorithm for maximizing  $L$ . Consider the first formulation of the log-likelihood that contains terms  $\ln(\alpha_j + k)$  and  $-\ln(|\boldsymbol{\alpha}| + k)$ . Applying Jensen's inequality to the concave term  $\ln(\alpha_j + k)$  and supporting hyperplane inequality to the convex term  $-\ln(|\boldsymbol{\alpha}| + k)$ , show that a minorizing function to  $L(\boldsymbol{\alpha})$  is

$$g(\boldsymbol{\alpha}|\boldsymbol{\alpha}^{(t)}) = - \sum_{k=0}^{\max_i |\mathbf{x}_i|-1} \frac{r_k}{|\boldsymbol{\alpha}^{(t)}| + k} |\boldsymbol{\alpha}| + \sum_{j=1}^d \sum_{k=0}^{\max_i x_{ij}-1} \frac{s_{jk} \alpha_j^{(t)}}{\alpha_j^{(t)} + k} \ln \alpha_j + c^{(t)},$$

where  $s_{jk} = \sum_{i=1}^n 1_{\{x_{ij} > k\}}$ ,  $r_k = \sum_{i=1}^n 1_{\{|\mathbf{x}_i| > k\}}$ , and  $c^{(t)}$  is a constant irrelevant to optimization. Maximizing the surrogate function  $g(\boldsymbol{\alpha}|\boldsymbol{\alpha}^{(t)})$  is trivial since  $\alpha_j$  are separated. Show that the MM updates are

$$\alpha_j^{(t+1)} = \frac{\sum_{k=0}^{\max_i x_{ij}-1} \frac{s_{jk}}{\alpha_j^{(t)} + k}}{\sum_{k=0}^{\max_i |\mathbf{x}_i|-1} \frac{r_k}{|\boldsymbol{\alpha}^{(t)}| + k}} \alpha_j^{(t)}, \quad j = 1, \dots, d.$$

The quantities  $s_{jk}$ ,  $r_k$ ,  $\max_i x_{ij}$  and  $\max_i |\mathbf{x}_i|$  only depend on data and can be pre-computed. Comment on whether the MM updates respect the parameter constraint  $\alpha_j > 0$ .

4. Write a function for finding MLE of Dirichlet-multinomial distribution given iid observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , using MM algorithm. The interface should be `dirmultfit.MM(x, alpha0 = NULL, maxiters = 1000, tolfun = 1e-6)`. The arguments are: `x` a  $n$ -by- $d$  matrix of counts, `alpha0` a  $d$  vector of starting point (optional), `maxiters` the maximum allowable MM iterations (default 1000), `tolfun` the tolerance for relative change in objective values (default 1e-6). The return value should be a list containing: `maximum` the log-likelihood at MLE, `estimate` the MLE, `gradient` the gradient at MLE, `hessian` the Hessian at MLE, `se` a  $d$  vector of standard errors, `iterations` the number of iterations performed.
5. Re-do HW6 Q10 using your new `dirmultfit.MM` function. Compare the number of iterations and run time by MM algorithm to those by Newton's method. Comment on the efficiency of Newton's algorithm vs MM algorithm for this problem.
6. Finally let us re-consider the EM algorithm. The difficulty with the M step in EM algorithm can be remedied. Discuss how we can further minorize the  $\ln \Gamma(|\boldsymbol{\alpha}|)$  term in the  $Q$  function to produce a minorizing function with all  $\alpha_j$  separated. For this homework, you do *not* need to implement this EM-MM hybrid algorithm. (Hint:  $z \mapsto \ln \Gamma(z)$  is a convex function.)