# LARGE RANDOM INTERCEPT LOGISTIC REGRESSION MODELS

JAN DE LEEUW

## 1. Model

Suppose we have $m$ contexts, indexed by $j$, and in context $j$ there are $n_j$ individuals, indexed by $i$. Responses, denoted by $\underline{y}_{ij}$ are binary[1]. For each $(i,j)$ there is also a vector with $p$ fixed predictors $x_{ij}$.

The model we consider is a *random intercept logistic regression model*. For each context, there is a random intercept $\underline{u}_j$. Thus

$$\text{(1)} \qquad \text{prob}(\underline{y}_{ij} = 1 | \underline{u}_j = u_j) = f(u_j + x'_{ij}\beta),$$

where $f$ is the logistic function, i.e.

$$\text{(2)} \qquad f(s) = \frac{\exp(s)}{1 + \exp(s)}.$$

We also assume that the $m$ random intercepts are i.i.d. And that, given the $\underline{u}_j$, all $\underline{y}_{ij}$ are independent.

For the $\underline{u}_j$ we also make the assumption that they are discrete, and they take the values $u_1, \cdots, u_t$ with probabilities $p_1, \cdots, p_t$. We will assume, in the sequel, that the $p_s$ are known, and that the $u_s$ and known up to a multiplicative constant $\sigma$. In the application we have in mind the $u_s$ are the knots and the $p_s$ are the weights associasted with a $t$-point Gauss-Hermite integration formula, i.e. we want to approximate a normal random intercept. In other application we can imagine actually optimizing over both the knots and the weights (and the number of knots), to obtain a semi-parametric logistic random intercept model (see Section 4).

The paper is concerned with the situation in which there are so many predictors that second order techniques are not really possible. We assume that it is impractical to compute the usual $p \times p$ matrices approximating the second derivatives that are used in the Newton or Gauss-Newton methods, i.e. in the methods of scoring and iterative weighted least squares.

---

[1]Random variables are underlined.

## 2. Likelihood

We now derive an expression for the log-likelihood from the assumptions in the previous section. Define

$$(3) \qquad \pi_{ij|s} = f_{ijs}^{y_{ij}}(1 - f_{ijs})^{1-y_{ij}},$$

$$(4) \qquad \pi_{j|s} = \prod_{i=1}^{n_j} \pi_{ij|s},$$

$$(5) \qquad \pi_j = \sum_{s=1}^{t} p_s \pi_{j|s},$$

where

$$(6) \qquad f_{ijs} = \frac{\exp(g_{ijs})}{1 + \exp(g_{ijs})},$$

and

$$(7) \qquad g_{ijs} = x'_{ij}\beta + \sigma u_s.$$

In fact, there is a more convenient way to write (7). If we define the $p + 1$-element vectors $z_{ijs}$ by

$$(8) \qquad z_{ijsr} = \begin{cases} x_{ijr} & \text{if } r \leq p \\ u_s & \text{if } r = p+1 \end{cases},$$

and the $p + 1$-element vector $\gamma$ by

$$(9) \qquad \gamma_r = \begin{cases} \beta_r & \text{if } r \leq p \\ \sigma & \text{if } r = p+1 \end{cases},$$

then

$$(10) \qquad g_{ijs} = z'_{ijs}\gamma.$$

Now

$$(11) \qquad \mathcal{L} = \sum_{j=1}^{m} \log \pi_j.$$

and some calculation gives

$$(12) \qquad \frac{\partial \mathcal{L}}{\partial \gamma_r} = \sum_{j=1}^{m} \frac{1}{\pi_j} \sum_{s=1}^{t} p_s \pi_{j|s} \sum_{i=1}^{n_j} (y_{ij} - f_{ijs}) z_{ijsr}.$$

From the computational point of view, observe that we never compute the $\pi_{ij|s}$. We simply cumulate the $\pi_{j|s}$ by multiplying the appropriate terms. Also, we never store the $z_{ijs}$, which have a lot of redundancy. It suffices to store the $x_{ij}$ and $u_s$. Formulas (11) and (12) are enough to

be able to apply conjugate gradient or other low-storage optimization methods, such as the ones in SAS/NLP[2].

## 3. Regularization by Majorization

The likelihood function discussed in the previous section is pretty complicated. Mixture models tend to give flat likelihoods, often with multiple local maxima. In this section we discuss an alternative approach, which combines conjugate gradient methods with use of the EM algorithm to majorize the likelihood locally by a concave function.

In order to majorize, we use the concavity of the logarithm (in this context also known as Jensen's inequality). Suppose $\tilde{\gamma}$ is our current best estimate of the parameters. Then

$$(13) \quad \sum_{j=1}^{m} \log \frac{\pi_j(\gamma)}{\pi_j(\tilde{\gamma})} = \sum_{j=1}^{m} \log \frac{\sum_{s=1}^{t} p_s \frac{\pi_{j|s}(\gamma)}{\pi_{j|s}(\tilde{\gamma})} \pi_{j|s}(\tilde{\gamma})}{\sum_{s=1}^{t} p_s \pi_{j|s}(\tilde{\gamma})} \geq$$

$$\sum_{j=1}^{m} \sum_{s=1}^{t} \frac{p_s \pi_{j|s}(\tilde{\gamma})}{\pi_j(\tilde{\gamma})} \log \frac{\pi_{j|s}(\gamma)}{\pi_{j|s}(\tilde{\gamma})} = \sum_{j=1}^{m} \sum_{s=1}^{t} \pi_{s|j}(\tilde{\gamma}) \log \frac{\pi_{j|s}(\gamma)}{\pi_{j|s}(\tilde{\gamma})}.$$

This can be written as

$$(14) \qquad \mathcal{L}(\gamma) \geq \mathcal{L}(\tilde{\gamma}) + \mathcal{K}(\gamma, \tilde{\gamma}) - \mathcal{K}(\tilde{\gamma}, \tilde{\gamma}),$$

with equality if and only if $\gamma = \tilde{\gamma}$. The only part depending on $\gamma$ is

$$(15) \quad \mathcal{K}(\gamma, \tilde{\gamma}) = \sum_{j=1}^{m} \sum_{s=1}^{t} \pi_{s|j}(\tilde{\gamma}) \log \pi_{j|s}(\gamma) =$$

$$\sum_{j=1}^{m} \sum_{s=1}^{t} \pi_{s|j}(\tilde{\gamma}) \sum_{i=1}^{n_j} [y_{ij} z'_{ijs} \gamma + \log(1 - f_{ijs})].$$

$\mathcal{K}(\gamma, \tilde{\gamma})$ is a concave function in $\gamma$, which must be maximized over $\gamma$. Suppose we do not necessarily maximize, but merely choose $\hat{\gamma}$ such that

$$(16) \qquad \mathcal{K}(\hat{\gamma}, \tilde{\gamma}) > \mathcal{K}(\tilde{\gamma}, \tilde{\gamma})$$

Then (14) says that

$$(17) \qquad \mathcal{L}(\hat{\gamma}) \geq \mathcal{L}(\tilde{\gamma}) + \mathcal{K}(\hat{\gamma}, \tilde{\gamma}) - \mathcal{K}(\tilde{\gamma}, \tilde{\gamma}),$$

and (16) then implies

$$(18) \qquad \mathcal{L}(\hat{\gamma}) > \mathcal{L}(\tilde{\gamma}).$$

---

[2]Observe there is no need to require that $\sigma \geq 0$. If our iterations converge to a $\sigma < 0$ we just change its sign and that of the $u_s$

Thus the likelihood is increased in each step, and since it is bounded above this least to a convergent algorithm (provided the function that updates $\gamma$ is chosen to be continuous).

Mazimizing $\mathcal{K}(\gamma, \tilde{\gamma})$ is actually very close to a standard logistic regression problem. In fact, the derivatives are

$$(19) \qquad \frac{\partial \mathcal{K}}{\partial \gamma_r} = \sum_{j=1}^{m} \sum_{s=1}^{t} \pi_{s|j}(\tilde{\gamma}) \sum_{i=1}^{n_j} z_{ijsr}(y_{ij} - f_{ijs}).$$

We now have various options. We can make one or more ascent steps to increase $\mathcal{K}(\gamma, \tilde{\gamma})$, then use the newly found $\gamma$ for $\tilde{\gamma}$, and continue with the next majorization.

## 4. Semiparametric Random Intercept Logistic Regression

If the knots $u_s$ and the weights $p_s$ are not known, we can optimize over them. This is easy to do for the $u_s$. Simply redefine (7) by writing

$$(20) \qquad g_{ijs} = x'_{ij}\beta + u_s,$$

where the $u_s$ are now additional unknowns[3]. This means that the vectors $z_{ijs}$ now have $p + t$ elements, defined by

$$(21) \qquad z_{ijsr} = \begin{cases} x_{ijr} & \text{if } r \leq p \\ 1 & \text{if } r = p + s \\ 0 & \text{otherwise.} \end{cases}$$

while $\gamma$ now has both the $p$ elements of $\beta$ and the $t$ elements of $u$.

Optimizing over $p$ can be done by utilizing *self-consistency*, which is basically just another example of majorization or EM. The iterative equation is

$$(22) \qquad p_s = \frac{1}{m} \sum_{j=1}^{m} \pi_{s|j}.$$

Again, in the actual implementation we can cycle through the unknowns $\gamma$ and $p$ in various ways. Generally, we alternate any number of $\gamma$ iterations with any number of $p$ iterations.

UCLA Program in Statistics, 405 Hilgard Avenue, Los Angeles, CA 90095-1554

---

[3]Again, there is no reason to require the $u_s$ to be ordered in any way. We can just order them after convergence, if we think this is appropriate