## Lecture 14: February 10

*Lecturer: Jonathan Taylor*                                      *Scribes: Kenneth Tay*

## 14.1 Gibbs Sampling for Probit Model

We work through the Gibbs sampler for the probit model to show how data augmentation can make Gibbs sampling easier.

We first specify the model:

- Prior $g(\beta) \propto \exp\left(-\beta^T \Sigma^{-1} \beta / 2\right)$ (i.e. normal),

- Likelihood $P(Y_i = 1 \mid X_i, \beta) = \Phi(X_i^T \beta)$. We can also think of $Y_i$ in terms of a latent variable $\varepsilon_i$:

$$Y_i = \begin{cases} 1 & \text{if } \varepsilon_i \leq X_i^T \beta, \\ 0 & \text{otherwise.} \end{cases}$$

  Here, $\varepsilon_i \overset{iid}{\sim} \mathcal{N}(0,1)$.

With these two pieces, we can compute the posterior density for $\beta$:

$$h(\beta \mid Y) = h(\beta \mid \varepsilon \leq X\beta)$$

$$\propto \left[ \int_{\{\varepsilon \leq X\beta\}} e^{-\|\varepsilon\|^2/2} d\varepsilon \right] \cdot g(\beta),$$

$$h(\varepsilon, \beta \mid \varepsilon \leq X\beta) \propto e^{-\|\varepsilon\|^2/2} \cdot e^{-\beta^T \Sigma^{-1} \beta / 2} 1_{\{\varepsilon \leq X\beta\}}.$$

Note that if we ignore the indicator variable, the remainder corresponds to the distribution $\mathbb{N}\left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} I & 0 \\ 0 & \Sigma \end{pmatrix} \right)$.

**Goal: Sample** $(\varepsilon, \beta) \mid \varepsilon \leq X\beta$.

Recall the Gibbs sampler from the previous lecture. WLOG, assume that $\Sigma = I$. (If not, we can write $\beta = \Sigma^{-1/2}\gamma$ with $\gamma \sim \mathcal{N}(0, I)$.)

Consider what happens to each step of choosing a new $\varepsilon_i$ or $\beta_i$.

- When moving a particular $\varepsilon_i$, fix $\varepsilon_{-i}$ and $\beta$. Then the density is $\propto e^{-\varepsilon_i^2/2} 1_{\{\varepsilon_i \leq X_i^T \beta\}}$. We can draw $\varepsilon_i^{new}$ from this density, i.e. $\sim \mathcal{N}(0, 1) \mid (-\infty, X_i^T \beta)$.

  Note that all the $\varepsilon_i^{new}$'s can be drawn simultaneously (no need to be sequential as in the usual Gibbs sampler)!

- When moving a particular $\beta_i$, fix $\varepsilon$ and $\beta_{-i}$. The constraints from the indicator variable become

$$\varepsilon_j \leq X_j^T \beta \qquad\qquad\qquad \text{for } 1 \leq j \leq n,$$
$$\varepsilon_j \leq \sum_{l \neq i} X_{jl}\beta_l + X_{ji}\beta_i \qquad\qquad \text{for } 1 \leq j \leq n,$$
$$X_{ji}\beta_i \geq \varepsilon_j - \sum_{l \neq i} X_{jl}\beta_l. \qquad\qquad \text{for } 1 \leq j \leq n,$$

These constraints bound $\beta_i$ on either side, depending on the sign of $X_{ji}$:

$$\max_{j:X_{ji}>0} \frac{1}{X_{ji}} \left[\varepsilon_j - \sum_{l \neq i} X_{jl}\beta_l\right] \leq \beta_i \leq \min_{j:X_{ji}<0} \frac{1}{X_{ji}} \left[\varepsilon_j - \sum_{l \neq i} X_{jl}\beta_l\right],$$
$$L(\beta_{-i}, \varepsilon) \leq \beta_i \leq U(\beta_{-i}, \varepsilon).$$

We draw $\beta_i^{new}$ from $\mathcal{N}(0,1) \mid [L, U]$.


## 14.2   (Agresti 8) Multinomial Regression

Recall the multinomial distribution: If $Y \sim \text{Multinom}(N, \pi)$, where $\pi \in \mathbb{R}_0^k$, $\sum \pi_i = 1$, then it has mass function

$$f(y_1, \ldots, y_k \mid \pi) = \binom{N}{y_1, \ldots, y_k} \prod_{i=1}^{k} \pi_i^{y_i},$$

supported on $y_1, \ldots, y_k \in \mathbb{Z}_0$, $\sum y_i = N$. Note that there are actually only $k-1$ free parameters in this model. If we take the $k^{th}$ category as the baseline category, we can write the mass function as

$$f(y_1, \ldots, y_{k-1} \mid \pi) = \binom{N}{y_1, \ldots, y_k} \prod_{i=1}^{k} \pi_i^{y_i},$$

where $y_k = N - \sum_{i=1}^{k-1} y_i$, $\pi_k = 1 - \sum_{i=1}^{k-1} \pi_i$. This mass function is supported on $y_1, \ldots, y_{j-1} \in \mathbb{Z}_0$, $\sum_{i=1}^{k-1} y_i \leq N$.

We can compute the log likelihood for this model:

$$\log L(\pi_1, \ldots, \pi_{k-1} \mid Y) = C + \sum_{i=1}^{k-1} y_i \log \pi_i + \left(N - \sum_{i=1}^{k-1} y_i\right) \log \left(1 - \sum_{i=1}^{k-1} \pi_i\right)$$
$$= C + \sum_{i=1}^{k-1} y_i \log \left(\frac{\pi_i}{1 - \sum_{i=1}^{k-1} \pi_i}\right) + N \log \left(1 - \sum_{j=1}^{k-1} \pi_j\right),$$

where $C$ is some constant that does not depend on $\pi$. The identity above shows that this model is an exponential family with sufficient statistic $(y_1, \ldots, y_{k-1})$ and natural parameters $\eta_j = \log \left(\frac{\pi_j}{1 - \sum_{l=1}^{k-1} \pi_l}\right)$.

We can invert the relationship between $\eta$ and $\pi$ to obtain $\pi_j = \dfrac{e^{\eta_j}}{1 + \sum_{l=1}^{k-1} e^{\eta_l}}$. This allows us to rewrite the log likelihood in terms of the natural parameters:

$$\log L(\eta_1, \ldots, \eta_{k-1} \mid Y) = \sum_{i=1}^{k-1} \eta_i y_i - N \log\left(1 + \sum_{i=1}^{k-1} e^{\eta_i}\right).$$

## 14.2.1 Baseline Multinomial Logit

In this model, we have $Y_i \overset{iid}{\sim} \mathrm{Multinom}(N_i, \pi_\beta(X_i))$, where for $1 \leq j \leq k-1$,

$$[\pi_\beta(X_i)]_j = \frac{\exp(X_i^T \beta_j)}{1 + \sum_{l=1}^{k-1} \exp(X_i^T \beta_l)}, \quad \text{or equivalently,} \quad \eta_i = X_i^T \beta.$$

Note here that $\beta$ is a $p \times (k-1)$ matrix. (In logistic regression, $\beta$ was just a $p \times 1$ matrix.) In the above expression for $[\pi_\beta(X_i)]_j$, $\beta_j$ refers to the $j^{th}$ column of $\beta$.

We can compute the log likelihood and its gradient for this model:

$$\log L(\beta \mid Y) = \sum_{i=1}^{n}\left[\sum_{l=1}^{k-1} Y_{il}(X_i^T \beta_l) - N_i \log\left(1 + \sum_{l=1}^{k-1} e^{X_i^T \beta_l}\right)\right],$$

$$\nabla \log L(\beta \mid Y) = X^T[Y - \mathbb{E}_\beta(Y)] \in \mathbb{R}^{p \times (k-1)}.$$

Note that the gradient of the log likelihood has the exact same form that we had for logistic and loglinear regression.

The Hessian $\nabla^2 \log L(\beta \mid Y)$ is some kind of tensor. It can also be thought of as $\mathrm{Cov}\,(\nabla \log L(\beta \mid Y))$. It is given by

$$[\nabla^2 \log L(\beta \mid Y)]_{ijkl} = \sum_{c=1}^{n} X_{ci} X_{ck} \mathrm{Cov}_\beta(Y_c)_{jl}.$$