

### Exercise 15

a) Density of the Inverse Gaussian distribution can be transformed as:

$$\begin{aligned} f(y|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi y^3 \sigma}} \exp\left(-\frac{1}{2y} \left(\frac{y-\mu}{\mu\sigma}\right)^2\right) = \left[\frac{1}{\sqrt{2\pi y^3 \sigma}} =: k(y, \sigma)\right] = \\ &= k(y, \sigma) \exp\left(-\frac{y^2 - 2\mu y + \mu^2}{2y\mu^2\sigma^2}\right) = k(y, \sigma) \exp\left(-\frac{1}{2} \left(\frac{y}{\mu^2\sigma^2} - \frac{2}{\mu\sigma^2} + \frac{1}{y\sigma^2}\right)\right) = \\ &= \left[c(y, \sigma) := k(y, \sigma) \exp\left(-\frac{1}{2} \left(\frac{1}{y\sigma^2}\right)\right)\right] = c(y, \sigma) \exp\left(\frac{-\frac{y}{2\mu^2} + \frac{1}{\mu}}{\sigma^2}\right) = \\ &= \left[\phi = \sigma^2, \theta = -\frac{1}{2\mu^2}, a(\theta) = -\sqrt{-2\theta} = \frac{1}{\mu}\right] = c(y, \phi) \exp\left(\frac{\theta y - a(\theta)}{\phi}\right) \end{aligned}$$

Thus Inverse Gaussian distribution with  $[\mu, \sigma^2]$  parametrization belongs to the exponential family with  $\phi = \sigma^2, \theta = -\frac{1}{2\mu^2}, a(\theta) = -\sqrt{-2\theta}$ .

b) We remember that for the Exponential Family distributions mean is  $E[y] = \mu = \dot{a}(\theta)$  and variance function is  $V[\mu] = \ddot{a}(\theta)$ , which are thus in our case equal to  $\mu = \dot{a}(\theta) = \frac{2}{2\sqrt{-2\theta}} = \frac{1}{\sqrt{-2\theta}}$  and  $V[\mu] = \ddot{a}(\theta) = \frac{-0.5 \times (-2)}{\sqrt{(-2\theta)^3}} = \frac{1}{\sqrt{(-2\theta)^3}} = \mu^3$  correspondingly.

c) The canonical link has the property that  $g(\mu) = \theta = \eta = x^T \beta$  it is invariant to being multiplied by the constant in terms of remaining canonical, i.e.  $g^*(\mu) = kg(\mu) = kx^T \beta = x^T \gamma$  is also canonical. Since in our case we have shown that  $\theta = -\frac{1}{2\mu^2}$  then  $g(\mu) = -\frac{1}{2\mu^2}$  or alternatively  $g^*(\mu) = \frac{1}{\mu^2}$  is also a canonical link for the Inverse Gaussian responses.

### Exercise 14

As a reminder to the lectures and addition to the discussion we had on plenaries.

We will in this exercise look at the connection between *weighted least squares* and the *Fisher-scoring algorithm* for GLMs (the iterative algorithm used to estimate the MLE of  $\beta$ ).

First look at the *weighted least squares*:

a) Consider the linear model (identity link):

$$y_i = \beta_0 x_{i0} + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n$$

The same regression model written on matrix form is:

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times (p+1)} \boldsymbol{\beta}_{(p+1) \times 1} + \boldsymbol{\epsilon}_{n \times 1}$$

Here  $x_{i0} = 1$  and  $\epsilon_i \sim N(0, \sigma^2/w_i)$ , where  $w_i$  signals the precision in observation  $i$ .

Let

$$y_i^* = \sqrt{w_i} y_i, \quad x_{ij}^* = \sqrt{w_i} x_{ij}, \quad \epsilon_i^* = \sqrt{w_i} \epsilon_i, \quad i = 1, \dots, n, \quad j = 0, \dots, p$$

Show that we now can write a regression model with  $y_i^*$  as response and  $x_{ij}^*$  as covariates where

the noise terms have constant variance:

$$\begin{aligned}\sqrt{w_i} \cdot y_i &= \sqrt{w_i} \cdot \left( \beta_0 x_{i0} + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i \right) \\ y_i^* &= \beta_0 x_{ij}^* + \dots + \beta_p x_{ip}^* + \epsilon_i^*\end{aligned}$$

where

$$E[\epsilon_i^*] = \sqrt{w_i} \cdot E[\epsilon_i] = 0 \quad \text{Var}[\epsilon_i^*] = \text{Var}[\sqrt{w_i} \epsilon_i] = w_i \text{Var}[\epsilon_i] = \sigma^2.$$

So that  $\mathbf{Y}^* \sim N(\mathbf{X}^* \boldsymbol{\beta}, \sigma^2 \mathbf{I})$ .

b) Calculate the least squares estimate  $\hat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$ :

Begin with the formulation using  $\mathbf{Y}^*$

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^{*T} \mathbf{X}^*)^{-1} \mathbf{X}^{*T} \mathbf{Y}^* \\ &= ((\sqrt{\mathbf{W}} \mathbf{X})^T \sqrt{\mathbf{W}} \mathbf{X})^{-1} (\sqrt{\mathbf{W}} \mathbf{X})^T \sqrt{\mathbf{W}} \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y}\end{aligned}$$

where  $\mathbf{W} = \text{diag}\{w_i\}$ .

We will now turn to GLMs:

c) We remember from class (page 67 in J and H) that the Score function and Expected Fisher info. for GLM are:

$$s_j(\boldsymbol{\beta}) = \frac{1}{\phi} \sum_{i=1}^n x_{ij} \frac{y_i - \mu_i}{g'(\mu_i) V(\mu_i)}, \quad I_{j,k}(\boldsymbol{\beta}) = \frac{1}{\phi} \sum_{i=1}^n \frac{x_{ij} x_{ik}}{g'(\mu_i)^2 V(\mu_i)}$$

where the  $\mu_i$ 's are indirectly specified through  $\boldsymbol{\beta}$ .  $V = \frac{d^2 a(\theta_i)}{(d\theta_i)^2}$

Writing it on matrix form:  $\boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\beta}) = (\mu_1, \dots, \mu_p)^T$ ,  $\mathbf{G}(\boldsymbol{\beta}) = \text{diag}\{g'(\mu_i)\}$ ,  $\mathbf{W}(\boldsymbol{\beta}) = \text{diag}\left\{\frac{1}{g'(\mu_i)^2 V(\mu_i)}\right\}$

Show that we now can express  $s$  and  $I$  in terms of:

$$\begin{aligned}s(\boldsymbol{\beta}) &= \frac{1}{\phi} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}) \mathbf{G}(\boldsymbol{\beta}) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \\ &= \frac{1}{\phi} \mathbf{X}^T \text{diag}\left(\frac{1}{g'(\mu_i)^2 V(\mu_i)}\right) \text{diag}(g'(\mu_i)) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \\ &= \frac{1}{\phi} \mathbf{X}^T \text{diag}\left(\frac{1}{g'(\mu_i) V(\mu_i)}\right) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \\ &= \frac{1}{\phi} \begin{bmatrix} x_{10} \frac{1}{g'(\mu_1) V(\mu_1)} & \dots & x_{n0} \frac{1}{g'(\mu_n) V(\mu_n)} \\ \vdots & & \vdots \\ x_{1p} \frac{1}{g'(\mu_1) V(\mu_1)} & \dots & x_{np} \frac{1}{g'(\mu_n) V(\mu_n)} \end{bmatrix} \begin{bmatrix} y_i - \mu_i \\ \vdots \\ y_n - \mu_n \end{bmatrix} \\ &= \frac{1}{\phi} \begin{bmatrix} \sum_{i=1}^n x_{i0} \frac{y_i - \mu_i}{g'(\mu_i) V(\mu_i)} \\ \vdots \\ \sum_{i=1}^n x_{ip} \frac{y_i - \mu_i}{g'(\mu_i) V(\mu_i)} \end{bmatrix}\end{aligned}$$

Hence,

$$s_j(\boldsymbol{\beta}) = \frac{1}{\phi} \sum_{i=1}^n x_{ij} \frac{y_i - \mu_i}{g'(\mu_i)V(\mu_i)} \quad j = 0, \dots, p.$$

$$\mathbf{I}(\boldsymbol{\beta}) = \frac{1}{\phi} \mathbf{X}^T \mathbf{W} \mathbf{X} = \frac{1}{\phi} \mathbf{X}^T \begin{bmatrix} \frac{x_{1,0}}{g'(\mu_1)^2 V(\mu_1)} & \cdots & \frac{x_{1,(p+1)}}{g'(\mu_1)^2 V(\mu_1)} \\ \vdots & & \vdots \\ \frac{x_{n,0}}{g'(\mu_n)^2 V(\mu_n)} & \cdots & \frac{x_{n,(p+1)}}{g'(\mu_n)^2 V(\mu_n)} \end{bmatrix}$$

$$= \frac{1}{\phi} \begin{bmatrix} \sum_{i=1}^n \frac{x_{i,0}x_{i,0}}{g'(\mu_i)^2 V(\mu_i)} & \cdots & \sum_{i=1}^n \frac{x_{i,0}x_{i,(p+1)}}{g'(\mu_i)^2 V(\mu_i)} \\ \vdots & \sum_{i=1}^n \frac{x_{ij}x_{ik}}{g'(\mu_i)^2 V(\mu_i)} & \vdots \\ \sum_{i=1}^n \frac{x_{i,(p+1)}x_{i,0}}{g'(\mu_i)^2 V(\mu_i)} & \cdots & \sum_{i=1}^n \frac{x_{i,(p+1)}x_{i,(p+1)}}{g'(\mu_i)^2 V(\mu_i)} \end{bmatrix}$$

$$I_{jk} = \frac{1}{\phi} \left\{ \sum_{i=1}^n \frac{x_{ij}x_{ik}}{g'(\mu_i)^2 V(\mu_i)} \right\}$$

(It says  $J$  in the exercise, but it should be  $I$ .  $J$  is often used as the observed fisher info)

d) Now we want to find an estimate of  $\boldsymbol{\beta}$ .

Show that the Fisher scoring algorithm can be written as:

$$\boldsymbol{\beta}^{(k+1)} = (\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{z}^{(k)}$$

where

$$\mathbf{z}^{(k)} = \mathbf{X} \boldsymbol{\beta}^{(k)} + \mathbf{G}(\boldsymbol{\beta}^{(k)}) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(k)}))$$

Begin with the Newton-Rahpson iteration: (page 69)

$$\begin{aligned} \boldsymbol{\beta}^{(k+1)} &= \boldsymbol{\beta}^{(k)} - l''(\boldsymbol{\beta}^{(k)})^{-1} l'(\boldsymbol{\beta}^{(k)}) \\ &= \boldsymbol{\beta}^{(k)} + \mathbf{J}(\boldsymbol{\beta}^{(k)})^{-1} \mathbf{s}(\boldsymbol{\beta}^{(k)}) \end{aligned}$$

Fisher suggested to replace  $\mathbf{J}$  by its expectation, thus  $\mathbf{I}$ . Using the previous results for  $\mathbf{s}$  and  $\mathbf{I}$  we get:

$$\begin{aligned} \boldsymbol{\beta}^{(k+1)} &= \boldsymbol{\beta}^{(k)} + \left( \frac{1}{\phi} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X} \right)^{-1} \frac{1}{\phi} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{G}(\boldsymbol{\beta}^{(k)}) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(k)})) \\ &= \boldsymbol{\beta}^{(k)} + \left( \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) [\mathbf{G}(\boldsymbol{\beta}^{(k)}) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(k)})) + \mathbf{X} \boldsymbol{\beta}^{(k)} - \mathbf{X} \boldsymbol{\beta}^{(k)}] \\ &= \boldsymbol{\beta}^{(k)} + \left( \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) [\mathbf{z}^{(k)} - \mathbf{X} \boldsymbol{\beta}^{(k)}] \end{aligned}$$

Let  $\mathbf{A} = \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X}$

$$\begin{aligned} \boldsymbol{\beta}^{(k+1)} &= \boldsymbol{\beta}^{(k)} + \mathbf{A}^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) [\mathbf{z}^{(k)} - \mathbf{X} \boldsymbol{\beta}^{(k)}] \\ &= \boldsymbol{\beta}^{(k)} + \mathbf{A}^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{z}^{(k)} - \mathbf{A}^{-1} \mathbf{A} \boldsymbol{\beta}^{(k)} \\ &= \boldsymbol{\beta}^{(k)} + \mathbf{A}^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{z}^{(k)} - \boldsymbol{\beta}^{(k)} \\ &= \mathbf{A}^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{z}^{(k)} \\ &= (\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{z}^{(k)} \end{aligned}$$

To find the estimate of  $\boldsymbol{\beta}$  in a GLM setting, we start with a guess  $\boldsymbol{\beta}^{(0)}$ , calculate  $\mathbf{W}(\boldsymbol{\beta}^{(0)})$  and  $\mathbf{z}^{(0)}$ , calculate a new guess  $\boldsymbol{\beta}^{(1)}$  and so on ... until we only see a small change in the update of  $\boldsymbol{\beta}$ .

The expression from the Fisher algorithm is similar to the weighted least square estimation. For simplicity we show the following results for one observation and for  $\beta$  being a scalar.

Looking at a link function  $g(\mu) = x\beta$  we find, by Taylor expansion:

$$\begin{aligned} g(y) &\approx g(\mu) + g'(\mu)(y - \mu) \\ g'(\mu) &\approx \frac{g(y) - g(\mu)}{y - \mu} \\ z &= x\beta + g'(\mu)(y - \mu) \approx g(\mu) + \frac{g(y) - g(\mu)}{y - \mu}(y - \mu) = g(y) \end{aligned}$$

So  $\boldsymbol{\beta}^{(k+1)} \approx (\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(k)}) g(\mathbf{Y})$ .

For the identity link:  $g(y) = y$  and  $g'(\mu) = 1$  we get the weighted least square formula.

$$z = x\beta + 1(y - \mu) = \mu + y - \mu = y$$

e) Looking at an example:

Assume now that  $Y_i \sim N(\mu_i, \sigma^2)$  with  $\mu_i = \sum_{j=0}^p \beta_j x_{ij}$ .

Show that the Fisher scoring algorithm converges in the first iteration to the least squares estimator.

The expectation  $E(y_i) = \mu_i = \sum_{j=0}^p \beta_j x_{ij} = \eta_i$ , thus we have the identity link:  $g(\mu_i) = \mu_i$ .

From the normal distribution (we have showed earlier that)

$$\theta_i = \mu_i, \quad a(\theta_i) = \frac{1}{2} \mu_i^2 = \frac{1}{2} \theta_i^2, \quad a'(\theta_i) = \theta_i, \quad a''(\theta) = 1,$$

and

$$V(\mu) = a''(\theta) = 1,$$

such that

$$\mathbf{W} = \text{diag}\left\{\frac{1}{g'(\mu_i)^2 V(\mu_i)}\right\} = \text{diag}(1)$$

---

Start with  $\boldsymbol{\beta}^{(0)}$ ,  $\boldsymbol{\mu}(\boldsymbol{\beta}^{(0)}) = \mathbf{X}\boldsymbol{\beta}^{(0)}$

$$\begin{aligned}\boldsymbol{\beta}^{(1)} &= (\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(0)}) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(0)}) \mathbf{z}(\boldsymbol{\beta}^{(0)}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{z}(\boldsymbol{\beta}^{(0)}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \left( \mathbf{X}\boldsymbol{\beta}^{(0)} + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(0)}) \right) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

which is the least square solution!