

Advanced Regression Methods for Independent Data

STAT/BIOST 570, 2020

Regression Models with Weaker Assumptions II

Mauricio Sadinle

Department of Biostatistics

 UNIVERSITY *of* WASHINGTON

Nonlinear Models with Weaker Assumptions

We now consider semiparametric inference for *nonlinear models*:

- ▶ We assume

$$Y_i = \mu(\mathbf{x}_i, \beta) + \epsilon_i,$$

where:

- ▶ $\mu(\mathbf{x}_i, \beta) = E(Y_i | \mathbf{x}_i)$: the regression function is given by the functional form $\mu(\mathbf{x}_i, \beta)$, which is nonlinear in β
- ▶ $E(\epsilon_i | \mathbf{x}_i) = 0$
- ▶ This formulation is equivalent to

$$Y_i | \mathbf{x}_i \sim H_i[\mu(\mathbf{x}_i, \beta)],$$

for some unspecified distribution H_i with mean $\mu(\mathbf{x}_i, \beta)$

Nonlinear Models with Weaker Assumptions

We now consider semiparametric inference for *nonlinear models*:

- ▶ We assume

$$Y_i = \mu(\mathbf{x}_i, \beta) + \epsilon_i,$$

where:

- ▶ $\mu(\mathbf{x}_i, \beta) = E(Y_i \mid \mathbf{x}_i)$: the regression function is given by the functional form $\mu(\mathbf{x}_i, \beta)$, which is nonlinear in β
- ▶ $E(\epsilon_i \mid \mathbf{x}_i) = 0$
- ▶ This formulation is equivalent to

$$Y_i \mid \mathbf{x}_i \sim H_i[\mu(\mathbf{x}_i, \beta)],$$

for some unspecified distribution H_i with mean $\mu(\mathbf{x}_i, \beta)$

Nonlinear Models with Weaker Assumptions

Some notation:

- For simplicity, we shall write

$$\mu_i := \mu_i(\boldsymbol{\beta}) := \mu(\mathbf{x}_i, \boldsymbol{\beta})$$

- We use the vector notation

$$\boldsymbol{\mu}(\boldsymbol{\beta}) := \begin{pmatrix} \mu_1(\boldsymbol{\beta}) \\ \vdots \\ \mu_n(\boldsymbol{\beta}) \end{pmatrix}$$

Estimating Equations for Nonlinear Models

- Inferences on β will be obtained from estimating equations with the form:

$$\mathbf{G}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i^T (Y_i - \mu_i(\beta)) = \frac{1}{n} \mathbf{Z}^T (\mathbf{Y} - \boldsymbol{\mu}(\beta)) = \mathbf{0},$$

- \mathbf{z}_i : $k + 1$ row vector (same length as β) that might depend on β but not on Y_i
- \mathbf{Z} : matrix containing the \mathbf{z}_i 's as its rows
- Note that *under the assumption* $\mu(\mathbf{x}_i, \beta) = E(Y_i \mid \mathbf{x}_i)$ we obtain

$$E[\mathbf{G}(\beta, Y_i, \mathbf{x}_i)] = E[\mathbf{z}_i^T (Y_i - \mu_i(\beta))] = \mathbf{0}$$

which we use to simplify some of the sandwich formulae

- These are called *linear unbiased estimating equations*

Estimating Equations for Nonlinear Models

- ▶ Inferences on β will be obtained from estimating equations with the form:

$$\mathbf{G}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i^T (Y_i - \mu_i(\beta)) = \frac{1}{n} \mathbf{Z}^T (\mathbf{Y} - \boldsymbol{\mu}(\beta)) = \mathbf{0},$$

- ▶ \mathbf{z}_i : $k + 1$ row vector (same length as β) that might depend on β but not on Y_i
- ▶ \mathbf{Z} : matrix containing the \mathbf{z}_i 's as its rows
- ▶ Note that *under the assumption* $\mu(\mathbf{x}_i, \beta) = E(Y_i \mid \mathbf{x}_i)$ we obtain

$$E[\mathbf{G}(\beta, Y_i, \mathbf{x}_i)] = E[\mathbf{z}_i^T (Y_i - \mu_i(\beta))] = \mathbf{0}$$

which we use to simplify some of the sandwich formulae

- ▶ These are called *linear unbiased estimating equations*

Least Squares for Nonlinear Models

- Consider the least-squares objective

$$SS(\beta) = \sum_{i=1}^n (Y_i - \mu_i(\beta))^2,$$

again with $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_k)^T$

- The minimizer of $SS(\beta)$ solves

$$\mathbf{G}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \mu_i}{\partial \beta} (Y_i - \mu_i(\beta)) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \frac{\partial \mu_i}{\partial \beta_0} (Y_i - \mu_i(\beta)) \\ \frac{\partial \mu_i}{\partial \beta_1} (Y_i - \mu_i(\beta)) \\ \vdots \\ \frac{\partial \mu_i}{\partial \beta_k} (Y_i - \mu_i(\beta)) \end{pmatrix} = \mathbf{0}.$$

Least Squares for Nonlinear Models

- Consider the least-squares objective

$$SS(\beta) = \sum_{i=1}^n (Y_i - \mu_i(\beta))^2,$$

again with $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_k)^T$

- The minimizer of $SS(\beta)$ solves

$$\mathbf{G}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \mu_i}{\partial \beta} (Y_i - \mu_i(\beta)) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \frac{\partial \mu_i}{\partial \beta_0} (Y_i - \mu_i(\beta)) \\ \frac{\partial \mu_i}{\partial \beta_1} (Y_i - \mu_i(\beta)) \\ \vdots \\ \frac{\partial \mu_i}{\partial \beta_k} (Y_i - \mu_i(\beta)) \end{pmatrix} = \mathbf{0}.$$

Least Squares for Nonlinear Models

- ▶ Noting that $\mu : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^n$, define

$$\mathbf{D} = \frac{\partial \mu}{\partial \beta} = \left(\frac{\partial \mu}{\partial \beta_0}, \quad \dots, \quad \frac{\partial \mu}{\partial \beta_k} \right) = \begin{pmatrix} \frac{\partial \mu_1}{\partial \beta_0} & \dots & \frac{\partial \mu_1}{\partial \beta_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mu_n}{\partial \beta_0} & \dots & \frac{\partial \mu_n}{\partial \beta_k} \end{pmatrix}$$

and note that

$$\mu'_i(\beta) := \frac{\partial}{\partial \beta} \mu_i(\beta) = \mathbf{D}_{[i,]}^\top \quad (\text{with } \mathbf{D}_{[i,]} \text{ the } i\text{th row of } \mathbf{D}),$$

- ▶ We can then express the least squares estimating equation for a nonlinear model in matrix format as

$$\mathbf{G}_n(\beta) = \frac{1}{n} \mathbf{D}^\top (\mathbf{Y} - \mu(\beta)) = \mathbf{0}$$

- ▶ The linear model is a special case where $\mu(\beta) = \mathbf{X}\beta$ and $\mathbf{D} = \mathbf{X}$, so the estimating equation reduces to the known form

$$\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta) = \mathbf{0}$$

Least Squares for Nonlinear Models

- ▶ Noting that $\mu : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^n$, define

$$\mathbf{D} = \frac{\partial \mu}{\partial \beta} = \left(\frac{\partial \mu}{\partial \beta_0}, \quad \dots, \quad \frac{\partial \mu}{\partial \beta_k} \right) = \begin{pmatrix} \frac{\partial \mu_1}{\partial \beta_0} & \dots & \frac{\partial \mu_1}{\partial \beta_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mu_n}{\partial \beta_0} & \dots & \frac{\partial \mu_n}{\partial \beta_k} \end{pmatrix}$$

and note that

$$\mu'_i(\beta) := \frac{\partial}{\partial \beta} \mu_i(\beta) = \mathbf{D}_{[i,]}^\top \quad (\text{with } \mathbf{D}_{[i,]} \text{ the } i\text{th row of } \mathbf{D}),$$

- ▶ We can then express the least squares estimating equation for a nonlinear model in matrix format as

$$\mathbf{G}_n(\beta) = \frac{1}{n} \mathbf{D}^\top (\mathbf{Y} - \mu(\beta)) = \mathbf{0}$$

- ▶ The linear model is a special case where $\mu(\beta) = \mathbf{X}\beta$ and $\mathbf{D} = \mathbf{X}$, so the estimating equation reduces to the known form

$$\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta) = \mathbf{0}$$

Least Squares for Nonlinear Models

- ▶ Noting that $\boldsymbol{\mu} : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^n$, define

$$\mathbf{D} = \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} = \left(\frac{\partial \boldsymbol{\mu}}{\partial \beta_0}, \quad \dots, \quad \frac{\partial \boldsymbol{\mu}}{\partial \beta_k} \right) = \begin{pmatrix} \frac{\partial \mu_1}{\partial \beta_0} & \dots & \frac{\partial \mu_1}{\partial \beta_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mu_n}{\partial \beta_0} & \dots & \frac{\partial \mu_n}{\partial \beta_k} \end{pmatrix}$$

and note that

$$\mu'_i(\boldsymbol{\beta}) := \frac{\partial}{\partial \boldsymbol{\beta}} \mu_i(\boldsymbol{\beta}) = \mathbf{D}_{[i,]}^\top \quad (\text{with } \mathbf{D}_{[i,]} \text{ the } i\text{th row of } \mathbf{D}),$$

- ▶ We can then express the least squares estimating equation for a nonlinear model in matrix format as

$$\mathbf{G}_n(\boldsymbol{\beta}) = \frac{1}{n} \mathbf{D}^\top (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) = \mathbf{0}$$

- ▶ The linear model is a special case where $\boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$ and $\mathbf{D} = \mathbf{X}$, so the estimating equation reduces to the known form

$$\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$$

Least Squares for Nonlinear Models

Different authors provide conditions that guarantee asymptotic normality of the least-squares estimator $\hat{\beta}_n$ in the nonlinear model:

$$\sqrt{n}\mathbf{B}_n^{-1/2}\mathbf{A}_n(\hat{\beta}_n - \beta) \xrightarrow{d} N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$$

where \mathbf{B}_n and \mathbf{A}_n are defined analogously as before, taking

$$\mathbf{G}(\beta, Y_i, \mathbf{x}_i) = \mu'_i(\beta)(Y_i - \mu_i(\beta)),$$

see, e.g. Boos and Stefanski (2013, sec 7.5.3)

Least Squares for Nonlinear Models

- We obtain

$$\begin{aligned}\mathbf{A}_n &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{G}'(\boldsymbol{\beta}, Y_i, \mathbf{x}_i)] \\ &= \frac{1}{n} \sum_{i=1}^n \{ \mathbb{E}[Y_i - \mu_i(\boldsymbol{\beta})] \mu_i''(\boldsymbol{\beta}) - \mu_i'(\boldsymbol{\beta}) \mu_i'(\boldsymbol{\beta})^\top \}\end{aligned}$$

where

$$\mu_i''(\boldsymbol{\beta}) = \frac{\partial^2}{\partial \boldsymbol{\beta}^\top \partial \boldsymbol{\beta}} \mu_i(\boldsymbol{\beta})$$

- Under the model assumption $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = \mathbb{E}(Y_i \mid \mathbf{x}_i)$

$$\mathbf{A}_n = -\frac{1}{n} \sum_{i=1}^n \mu_i'(\boldsymbol{\beta}) \mu_i'(\boldsymbol{\beta})^\top = -\mathbf{D}^\top \mathbf{D} / n$$

Least Squares for Nonlinear Models

- We obtain

$$\begin{aligned}\mathbf{A}_n &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{G}'(\boldsymbol{\beta}, Y_i, \mathbf{x}_i)] \\ &= \frac{1}{n} \sum_{i=1}^n \{ \mathbb{E}[Y_i - \mu_i(\boldsymbol{\beta})] \mu_i''(\boldsymbol{\beta}) - \mu_i'(\boldsymbol{\beta}) \mu_i'(\boldsymbol{\beta})^\top \}\end{aligned}$$

where

$$\mu_i''(\boldsymbol{\beta}) = \frac{\partial^2}{\partial \boldsymbol{\beta}^\top \partial \boldsymbol{\beta}} \mu_i(\boldsymbol{\beta})$$

- Under the model assumption $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = \mathbb{E}(Y_i \mid \mathbf{x}_i)$

$$\mathbf{A}_n = -\frac{1}{n} \sum_{i=1}^n \mu_i'(\boldsymbol{\beta}) \mu_i'(\boldsymbol{\beta})^\top = -\mathbf{D}^\top \mathbf{D} / n$$

Least Squares for Nonlinear Models

► Similarly,

$$\begin{aligned}\mathbf{B}_n &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{G}(\boldsymbol{\beta}, Y_i, \mathbf{x}_i) \mathbf{G}(\boldsymbol{\beta}, Y_i, \mathbf{x}_i)^\top] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(Y_i - \mu_i(\boldsymbol{\beta}))^2] \mu'_i(\boldsymbol{\beta}) \mu'_i(\boldsymbol{\beta})^\top \\ &= \mathbf{D}^\top \text{diag}\{\mathbb{E}[(Y_i - \mu_i(\boldsymbol{\beta}))^2]\} \mathbf{D} / n\end{aligned}$$

► If we assume homoskedasticity, that is, $Y_i = \mu(\mathbf{x}_i, \boldsymbol{\beta}) + \epsilon_i$ with $\text{var}(\epsilon_i | \mathbf{x}_i) = \sigma^2$ then

$$\mathbf{B}_n = \sigma^2 \frac{1}{n} \sum_{i=1}^n \mu'_i(\boldsymbol{\beta}) \mu'_i(\boldsymbol{\beta})^\top = \sigma^2 \mathbf{D}^\top \mathbf{D} / n$$

Least Squares for Nonlinear Models

- Similarly,

$$\begin{aligned}\mathbf{B}_n &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{G}(\boldsymbol{\beta}, Y_i, \mathbf{x}_i) \mathbf{G}(\boldsymbol{\beta}, Y_i, \mathbf{x}_i)^\top] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(Y_i - \mu_i(\boldsymbol{\beta}))^2] \mu'_i(\boldsymbol{\beta}) \mu'_i(\boldsymbol{\beta})^\top \\ &= \mathbf{D}^\top \text{diag}\{\mathbb{E}[(Y_i - \mu_i(\boldsymbol{\beta}))^2]\} \mathbf{D} / n\end{aligned}$$

- If we assume homoskedasticity, that is, $Y_i = \mu(\mathbf{x}_i, \boldsymbol{\beta}) + \epsilon_i$ with $\text{var}(\epsilon_i \mid \mathbf{x}_i) = \sigma^2$ then

$$\mathbf{B}_n = \sigma^2 \frac{1}{n} \sum_{i=1}^n \mu'_i(\boldsymbol{\beta}) \mu'_i(\boldsymbol{\beta})^\top = \sigma^2 \mathbf{D}^\top \mathbf{D} / n$$

Least Squares for Nonlinear Models

The sandwich estimator of the variance of $\hat{\beta}$ is then

$$\widehat{\text{var}}(\hat{\beta}) = \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{-1})^\top / n,$$

where for $\hat{\mathbf{A}}_n$ we have

- ▶ If we allow $\mu(\mathbf{x}_i, \beta) \neq \text{E}(Y_i \mid \mathbf{x}_i)$, then

$$\hat{\mathbf{A}}_n = \frac{1}{n} \sum_{i=1}^n \{ [Y_i - \mu_i(\hat{\beta})] \mu_i''(\hat{\beta}) - \mu_i'(\hat{\beta}) \mu_i'(\hat{\beta})^\top \},$$

for example, if we want to estimate $\beta_0 = \underset{\beta}{\operatorname{argmin}} \text{E}_F \{ [Y - \mu(\beta, \mathbf{x})]^2 \}$

- ▶ Under the model assumption $\mu(\mathbf{x}_i, \beta) = \text{E}(Y_i \mid \mathbf{x}_i)$

$$\hat{\mathbf{A}}_n = -\frac{1}{n} \sum_{i=1}^n \mu_i'(\hat{\beta}) \mu_i'(\hat{\beta})^\top = -\hat{\mathbf{D}}^\top \hat{\mathbf{D}} / n$$

Least Squares for Nonlinear Models

The sandwich estimator of the variance of $\hat{\beta}$ is then

$$\widehat{\text{var}}(\hat{\beta}) = \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{-1})^\top / n,$$

where for $\hat{\mathbf{A}}_n$ we have

- ▶ If we allow $\mu(\mathbf{x}_i, \beta) \neq \text{E}(Y_i \mid \mathbf{x}_i)$, then

$$\hat{\mathbf{A}}_n = \frac{1}{n} \sum_{i=1}^n \{ [Y_i - \mu_i(\hat{\beta})] \mu_i''(\hat{\beta}) - \mu_i'(\hat{\beta}) \mu_i'(\hat{\beta})^\top \},$$

for example, if we want to estimate $\beta_0 = \underset{\beta}{\text{argmin}} \text{E}_F \{ [Y - \mu(\beta, \mathbf{x})]^2 \}$

- ▶ Under the model assumption $\mu(\mathbf{x}_i, \beta) = \text{E}(Y_i \mid \mathbf{x}_i)$

$$\hat{\mathbf{A}}_n = -\frac{1}{n} \sum_{i=1}^n \mu_i'(\hat{\beta}) \mu_i'(\hat{\beta})^\top = -\hat{\mathbf{D}}^\top \hat{\mathbf{D}} / n$$

Least Squares for Nonlinear Models

$$\widehat{\text{var}}(\hat{\beta}) = \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{-1})^\top / n,$$

where for $\hat{\mathbf{B}}_n$ we have

► In general

$$\begin{aligned}\hat{\mathbf{B}}_n &= \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_i(\hat{\beta}))^2 \mu'_i(\hat{\beta}) \mu'_i(\hat{\beta})^\top \\ &= \hat{\mathbf{D}}^\top \text{diag}\{(Y_i - \mu_i(\hat{\beta}))^2\} \hat{\mathbf{D}} / n\end{aligned}$$

► If we assume homoskedasticity, that is, $Y_i = \mu(x_i, \beta) + \epsilon_i$ with $\text{var}(\epsilon_i \mid x_i) = \sigma^2$ then

$$\hat{\mathbf{B}}_n = \hat{\sigma}^2 \frac{1}{n} \sum_{i=1}^n \mu'_i(\hat{\beta}) \mu'_i(\hat{\beta})^\top = \hat{\sigma}^2 \hat{\mathbf{D}}^\top \hat{\mathbf{D}} / n$$

with

$$\hat{\sigma}^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (Y_i - \mu_i(\hat{\beta}))^2$$

Least Squares for Nonlinear Models

$$\widehat{\text{var}}(\hat{\beta}) = \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{-1})^\top / n,$$

where for $\hat{\mathbf{B}}_n$ we have

► In general

$$\begin{aligned}\hat{\mathbf{B}}_n &= \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_i(\hat{\beta}))^2 \mu'_i(\hat{\beta}) \mu'_i(\hat{\beta})^\top \\ &= \hat{\mathbf{D}}^\top \text{diag}\{(Y_i - \mu_i(\hat{\beta}))^2\} \hat{\mathbf{D}} / n\end{aligned}$$

► If we assume homoskedasticity, that is, $Y_i = \mu(\mathbf{x}_i, \beta) + \epsilon_i$ with $\text{var}(\epsilon_i \mid \mathbf{x}_i) = \sigma^2$ then

$$\hat{\mathbf{B}}_n = \hat{\sigma}^2 \frac{1}{n} \sum_{i=1}^n \mu'_i(\hat{\beta}) \mu'_i(\hat{\beta})^\top = \hat{\sigma}^2 \hat{\mathbf{D}}^\top \hat{\mathbf{D}} / n$$

with

$$\hat{\sigma}^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (Y_i - \mu_i(\hat{\beta}))^2$$

Least Squares for Nonlinear Models

- Under the assumption $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = E(Y_i | \mathbf{x}_i)$, we finally have

$$\begin{aligned}\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) &= \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{-1})^\top / n, \\ &= (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1} [\hat{\mathbf{D}}^\top \text{diag}\{(Y_i - \mu_i(\hat{\boldsymbol{\beta}}))^2\} \hat{\mathbf{D}}] (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1}\end{aligned}$$

- Under homoskedasticity

$$\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2 (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1}$$

with

$$\hat{\sigma}^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (Y_i - \mu_i(\hat{\boldsymbol{\beta}}))^2$$

Least Squares for Nonlinear Models

- Under the assumption $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = E(Y_i | \mathbf{x}_i)$, we finally have

$$\begin{aligned}\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) &= \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{-1})^\top / n, \\ &= (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1} [\hat{\mathbf{D}}^\top \text{diag}\{(Y_i - \mu_i(\hat{\boldsymbol{\beta}}))^2\} \hat{\mathbf{D}}] (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1}\end{aligned}$$

- Under homoskedasticity

$$\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2 (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1}$$

with

$$\hat{\sigma}^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (Y_i - \mu_i(\hat{\boldsymbol{\beta}}))^2$$

Estimating Equations for Nonlinear Models

Least squares is *not* the only option for obtaining estimating equations for nonlinear models

- Consider, for example,

$$\begin{aligned} \mathbf{G}_n(\boldsymbol{\beta}) &= \frac{1}{n} \mathbf{X}^\top (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top (Y_i - \mu_i(\boldsymbol{\beta})) \\ &= \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} (Y_i - \mu_i(\boldsymbol{\beta})) \\ x_{1i}(Y_i - \mu_i(\boldsymbol{\beta})) \\ \vdots \\ x_{ki}(Y_i - \mu_i(\boldsymbol{\beta})) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \end{aligned}$$

where $\mu_i(\boldsymbol{\beta})$ is a general nonlinear function in $\boldsymbol{\beta}$

- Note that these expressions correspond to the score equations of GLMs with canonical links and equal dispersion parameters $\alpha_i = \alpha$ (see slides4.pdf, p. 39), and also to the OLS equations with $\mu_i(\boldsymbol{\beta}) = \mathbf{x}_i \boldsymbol{\beta}$

Estimating Equations for Nonlinear Models

- ▶ Another estimating equation with this generic form we previously saw was

$$\mathbf{S}_n(\boldsymbol{\beta}) = \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) = \mathbf{0},$$

which we obtained from GLMs with general link functions, with specification of the mean and variance as

$$E(Y_i | \mathbf{x}_i) = \mu_i(\boldsymbol{\beta}), \quad \text{var}(Y_i | \mathbf{x}_i) = \alpha_i V(\mu_i),$$

where $\alpha_i = \alpha / \phi_i$ with ϕ_i known, or in matrix form

$$E(\mathbf{Y} | \mathbf{X}) = \boldsymbol{\mu}(\boldsymbol{\beta}), \quad \text{var}(\mathbf{Y} | \mathbf{X}) = \alpha \mathbf{V},$$

with $\mathbf{V} = \text{diag}\{V(\mu_i)/\phi_i\}$ (see slides4.pdf, p. 25)

Any of these approaches can be used to obtain a Z-estimator and its sandwich covariance matrix, analogously as before!

Estimating Equations for Nonlinear Models

- ▶ All the estimating equations seen above have the form:

$$\mathbf{G}_n(\boldsymbol{\beta}) = \frac{1}{n} \mathbf{Z}^T (\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})) = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i^T (Y_i - \mu_i(\boldsymbol{\beta})) = \mathbf{0},$$

where \mathbf{z}_i is a row vector that might depend on $\boldsymbol{\beta}$, and \mathbf{Z} is a matrix containing the \mathbf{z}_i 's as its rows

- ▶ How to choose \mathbf{Z} ?
 - ▶ We might hope to choose \mathbf{Z} in a way that assures good properties of our solutions/estimator $\hat{\boldsymbol{\beta}}$
 - ▶ One criterion is to select \mathbf{Z} to minimize the large sample approximate $\text{var}(\hat{\boldsymbol{\beta}})$

Optimal Estimating Equations

- ▶ *Question:* If we know both the form of $\mu(\beta)$ and also $\text{var}(\mathbf{Y} \mid \mathbf{X})$, can we derive estimating equations that produce consistent estimates of β with smaller variance?
- ▶ *Important Result:* Let $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$. Among estimating functions of the form $\mathbf{G}_n(\beta) = \mathbf{Z}^T (\mathbf{Y} - \mu(\beta))/n$, setting $\mathbf{Z}^T = \mathbf{D}^T \mathbf{V}^{-1}$ yields an estimator $\hat{\beta}$ with the smallest asymptotic variance.
- ▶ Therefore, if we knew $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$, to get more precise estimates of β , find $\hat{\beta}$ that solves

$$\mathbf{U}_n(\beta) = \frac{1}{n} \mathbf{D}^T(\beta) \mathbf{V}^{-1}(\beta) (\mathbf{Y} - \mu(\beta)) = \mathbf{0}$$

Optimal Estimating Equations

- ▶ *Question:* If we know both the form of $\mu(\beta)$ and also $\text{var}(\mathbf{Y} \mid \mathbf{X})$, can we derive estimating equations that produce consistent estimates of β with smaller variance?
- ▶ *Important Result:* Let $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$. Among estimating functions of the form $\mathbf{G}_n(\beta) = \mathbf{Z}^T (\mathbf{Y} - \mu(\beta))/n$, setting $\mathbf{Z}^T = \mathbf{D}^T \mathbf{V}^{-1}$ yields an estimator $\hat{\beta}$ with the smallest asymptotic variance.
- ▶ Therefore, if we knew $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$, to get more precise estimates of β , find $\hat{\beta}$ that solves

$$\mathbf{U}_n(\beta) = \frac{1}{n} \mathbf{D}^T(\beta) \mathbf{V}^{-1}(\beta) (\mathbf{Y} - \mu(\beta)) = \mathbf{0}$$

Optimal Estimating Equations

- ▶ *Question:* If we know both the form of $\mu(\beta)$ and also $\text{var}(\mathbf{Y} \mid \mathbf{X})$, can we derive estimating equations that produce consistent estimates of β with smaller variance?
- ▶ *Important Result:* Let $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$. Among estimating functions of the form $\mathbf{G}_n(\beta) = \mathbf{Z}^T (\mathbf{Y} - \mu(\beta))/n$, setting $\mathbf{Z}^T = \mathbf{D}^T \mathbf{V}^{-1}$ yields an estimator $\hat{\beta}$ with the smallest asymptotic variance.
- ▶ Therefore, if we knew $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$, to get more precise estimates of β , find $\hat{\beta}$ that solves

$$\mathbf{U}_n(\beta) = \frac{1}{n} \mathbf{D}^T(\beta) \mathbf{V}^{-1}(\beta) (\mathbf{Y} - \mu(\beta)) = \mathbf{0}$$

Optimal Estimating Equations

- ▶ This result is due to Godambe and Heyde (1987, *International Statistical Review*) who “state the Gauss-Markov theorem in the framework of estimating function theory.”
- ▶ See also Heyde (1997), *Quasi-Likelihood and its Application*, Springer.
- ▶ Additional assumptions by Godambe and Heyde (1987):
 - ▶ $E(Y | \mathbf{x}) < \infty$, $E(\mathbf{Y} | \mathbf{X}) = \boldsymbol{\mu}(\boldsymbol{\beta}, \mathbf{X})$, $\text{var}(\mathbf{Y} | \mathbf{X}) \propto \mathbf{V}$.
 - ▶ $E[\frac{\partial \mathbf{G}_n}{\partial \boldsymbol{\beta}}]$ is of full rank and $E[\mathbf{G}_n(\boldsymbol{\beta}) \mathbf{G}_n(\boldsymbol{\beta})^T]$ and $E[\mathbf{U}_n(\boldsymbol{\beta}) \mathbf{U}_n(\boldsymbol{\beta})^T]$ are positive definite.

Optimal Estimating Equations

Proof outline of the result of Godambe and Heyde:

- ▶ Remember, the asymptotic variance of $\hat{\beta}$, the solution of

$$\mathbf{G}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\beta, Y_i, \mathbf{x}_i) = \mathbf{0},$$

is given by $\mathbf{A}_n^{-1} \mathbf{B}_n (\mathbf{A}_n^{-1})^T / n$, where \mathbf{A}_n can be written as

$$\mathbf{A}_n := \mathbf{A}_G = \mathbb{E} \left[\frac{\partial \mathbf{G}_n}{\partial \beta} \right] = \left(\mathbb{E} \left[\frac{\partial \mathbf{G}_n}{\partial \beta_0} \right], \dots, \mathbb{E} \left[\frac{\partial \mathbf{G}_n}{\partial \beta_k} \right] \right),$$

with $\mathbb{E} \left[\frac{\partial \mathbf{G}_n}{\partial \beta_j} \right]$ a column of length $k+1$, and \mathbf{B}_n can be written as

$$\begin{aligned} \mathbf{B}_n := \mathbf{B}_G &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\mathbf{G}(\beta, Y_i, \mathbf{x}_i) \mathbf{G}(\beta, Y_i, \mathbf{x}_i)^T \right] \\ &= n \mathbb{E} \left[\mathbf{G}_n(\beta) \mathbf{G}_n(\beta)^T \right], \end{aligned}$$

since $\mathbb{E}[\mathbf{G}(\beta, Y_i, \mathbf{x}_i)] = \mathbf{0}$ by assumption, and

$$\mathbb{E}[\mathbf{G}(\beta, Y_i, \mathbf{x}_i) \mathbf{G}(\beta, Y_j, \mathbf{x}_j)^T] = \mathbf{0},$$

for all $i \neq j$ since the observations are independent.

Optimal Estimating Equations

Proof outline of the result of Godambe and Heyde, cont'd:

- ▶ The general linear unbiased estimating equation

$$\mathbf{G}_n(\boldsymbol{\beta}) = \mathbf{Z}^T(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))/n = \mathbf{0},$$

leads to

$$\mathbf{A}_G = \mathbf{Z}^T \mathbf{D} / n,$$

$$\mathbf{B}_G = \alpha \mathbf{Z}^T \mathbf{V} \mathbf{Z} / n,$$

and asymptotic variance of the solution of $\mathbf{G}_n(\boldsymbol{\beta}) = \mathbf{0}$, $\hat{\boldsymbol{\beta}}_G$, as

$$\text{var}(\hat{\boldsymbol{\beta}}_G) = \alpha (\mathbf{Z}^T \mathbf{D})^{-1} (\mathbf{Z}^T \mathbf{V} \mathbf{Z}) (\mathbf{Z}^T \mathbf{D})^{-1}$$

Optimal Estimating Equations

Proof outline of the result of Godambe and Heyde, cont'd:

- ▶ The optimal linear unbiased EE is actually defined up to an invertible constant matrix \mathbf{C}

$$\mathbf{U}_n(\beta) = \mathbf{C}\mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{Y} - \mu(\beta))/n = \mathbf{0},$$

which leads to

$$\mathbf{A}_U = \mathbf{C}\mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D}/n,$$

$$\mathbf{B}_U = \alpha \mathbf{C}\mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D}\mathbf{C}^\top/n,$$

and asymptotic variance of the solution of $\mathbf{U}_n(\beta) = \mathbf{0}$, $\hat{\beta}_U$, as

$$\text{var}(\hat{\beta}_U) = \alpha(\mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D})^{-1} = \alpha^2 \mathbf{C}^\top \mathbf{B}_U^{-1} \mathbf{C}/n$$

Optimal Estimating Equations

Proof outline of the result of Godambe and Heyde, cont'd:

- ▶ Also, we obtain

$$E \left[\mathbf{U}_n(\beta) \mathbf{G}_n(\beta)^T \right] = \alpha \mathbf{C} \mathbf{D}^T \mathbf{V}^{-1} \mathbf{V} \mathbf{Z} / n^2 = \alpha \mathbf{C} \mathbf{D}^T \mathbf{Z} / n^2 = \alpha \mathbf{C} \mathbf{A}_G^T / n$$

- ▶ Let

$$\mathbf{M} = \begin{pmatrix} E[\mathbf{U}_n(\beta) \mathbf{U}_n(\beta)^T] & E[\mathbf{U}_n(\beta) \mathbf{G}_n(\beta)^T] \\ E[\mathbf{U}_n(\beta) \mathbf{G}_n(\beta)^T]^T & E[\mathbf{G}_n(\beta) \mathbf{G}_n(\beta)^T] \end{pmatrix} = \begin{pmatrix} \mathbf{B}_U / n & \alpha \mathbf{C} \mathbf{A}_G^T / n \\ \alpha \mathbf{A}_G \mathbf{C}^T / n & \mathbf{B}_G / n \end{pmatrix}$$

be the variance-covariance matrix of $(\mathbf{U}_n(\beta), \mathbf{G}_n(\beta))^T$, since $E[\mathbf{U}_n(\beta)] = \mathbf{0}$, $E[\mathbf{G}_n(\beta)] = \mathbf{0}$.

- ▶ The Schur complement of \mathbf{B}_U / n in \mathbf{M} is

$$(\mathbf{B}_G - \alpha^2 \mathbf{A}_G \mathbf{C}^T \mathbf{B}_U^{-1} \mathbf{C} \mathbf{A}_G^T) / n,$$

which is non negative definite since \mathbf{M} is non negative definite and \mathbf{B}_U is assumed to be positive definite (property of the Schur complement)

Optimal Estimating Equations

Proof outline of the result of Godambe and Heyde, cont'd:

- Note that if

$$\mathbf{a}^T [\mathbf{B}_G - \alpha^2 \mathbf{A}_G \mathbf{C}^T \mathbf{B}_U^{-1} \mathbf{C} \mathbf{A}_G^T] \mathbf{a} \geq 0$$

for all \mathbf{a} , then

$$(\mathbf{A}_G^T \mathbf{a})^T [(\mathbf{A}_G)^{-1} \mathbf{B}_G (\mathbf{A}_G^T)^{-1} - \alpha^2 \mathbf{C}^T \mathbf{B}_U^{-1} \mathbf{C}] (\mathbf{A}_G^T \mathbf{a}) \geq 0,$$

for all $\mathbf{a}' = \mathbf{A}_G^T \mathbf{a}$. Since \mathbf{A}_G^T is invertible, \mathbf{a}' can be any vector in \mathbb{R}^{k+1} .

- Therefore,

$$(\mathbf{A}_G)^{-1} \mathbf{B}_G (\mathbf{A}_G^T)^{-1} / n - \alpha^2 \mathbf{C}^T \mathbf{B}_U^{-1} \mathbf{C} / n = \text{var}(\hat{\beta}_G) - \text{var}(\hat{\beta}_U)$$

is non negative definite.

This completes the proof of Godambe and Heyde's result.

Optimal Estimating Equations

These results lead to an interesting connection with GLMs:

- ▶ Regardless of where the Y_i 's live, we could think of a generative model

$$Y_i = \mu_i(\beta) + \epsilon_i,$$

which certainly holds for $\epsilon_i = Y_i - \mu_i(\beta)$, with $Y_i \mid x_i \sim H_i[\mu_i(\beta)]$

- ▶ If the Y_i 's are independent given covariates, so are the ϵ_i 's
- ▶ $\text{var}(\epsilon_i \mid x_i) = \text{var}(Y_i \mid x_i)$ and say we assume $\text{var}(Y_i \mid x_i) = \alpha V(\mu_i)$, where $\alpha V(\mu_i)$ coincides with one of the variances in a GLM
- ▶ So a GLM can be seen as a nonlinear model $Y_i = \mu_i(\beta) + \epsilon_i$ with a specific mean and variance structure
- ▶ The optimal estimating equations of Godambe and Heyde (1987)

$$U_n(\beta) = D^T V^{-1}(\mathbf{Y} - \boldsymbol{\mu})/n = \mathbf{0}$$

correspond to the score equations under maximum likelihood estimation

Optimal Estimating Equations

These results lead to an interesting connection with GLMs:

- ▶ Regardless of where the Y_i 's live, we could think of a generative model

$$Y_i = \mu_i(\beta) + \epsilon_i,$$

which certainly holds for $\epsilon_i = Y_i - \mu_i(\beta)$, with $Y_i \mid x_i \sim H_i[\mu_i(\beta)]$

- ▶ If the Y_i 's are independent given covariates, so are the ϵ_i 's
- ▶ $\text{var}(\epsilon_i \mid x_i) = \text{var}(Y_i \mid x_i)$ and say we assume $\text{var}(Y_i \mid x_i) = \alpha V(\mu_i)$, where $\alpha V(\mu_i)$ coincides with one of the variances in a GLM
- ▶ So a GLM can be seen as a nonlinear model $Y_i = \mu_i(\beta) + \epsilon_i$ with a specific mean and variance structure
- ▶ The optimal estimating equations of Godambe and Heyde (1987)

$$U_n(\beta) = D^T V^{-1}(\mathbf{Y} - \boldsymbol{\mu})/n = \mathbf{0}$$

correspond to the score equations under maximum likelihood estimation

Optimal Estimating Equations

These results lead to an interesting connection with GLMs:

- ▶ Regardless of where the Y_i 's live, we could think of a generative model

$$Y_i = \mu_i(\beta) + \epsilon_i,$$

which certainly holds for $\epsilon_i = Y_i - \mu_i(\beta)$, with $Y_i \mid \mathbf{x}_i \sim H_i[\mu_i(\beta)]$

- ▶ If the Y_i 's are independent given covariates, so are the ϵ_i 's
- ▶ $\text{var}(\epsilon_i \mid \mathbf{x}_i) = \text{var}(Y_i \mid \mathbf{x}_i)$ and say we assume $\text{var}(Y_i \mid \mathbf{x}_i) = \alpha V(\mu_i)$, where $\alpha V(\mu_i)$ coincides with one of the variances in a GLM
- ▶ So a GLM can be seen as a nonlinear model $Y_i = \mu_i(\beta) + \epsilon_i$ with a specific mean and variance structure
- ▶ The optimal estimating equations of Godambe and Heyde (1987)

$$U_n(\beta) = D^T V^{-1}(\mathbf{Y} - \boldsymbol{\mu})/n = \mathbf{0}$$

correspond to the score equations under maximum likelihood estimation

Optimal Estimating Equations

These results lead to an interesting connection with GLMs:

- ▶ Regardless of where the Y_i 's live, we could think of a generative model

$$Y_i = \mu_i(\beta) + \epsilon_i,$$

which certainly holds for $\epsilon_i = Y_i - \mu_i(\beta)$, with $Y_i \mid \mathbf{x}_i \sim H_i[\mu_i(\beta)]$

- ▶ If the Y_i 's are independent given covariates, so are the ϵ_i 's
- ▶ $\text{var}(\epsilon_i \mid \mathbf{x}_i) = \text{var}(Y_i \mid \mathbf{x}_i)$ and say we assume $\text{var}(Y_i \mid \mathbf{x}_i) = \alpha V(\mu_i)$, where $\alpha V(\mu_i)$ coincides with one of the variances in a GLM
- ▶ So a GLM can be seen as a nonlinear model $Y_i = \mu_i(\beta) + \epsilon_i$ with a specific mean and variance structure

- ▶ The optimal estimating equations of Godambe and Heyde (1987)

$$U_n(\beta) = D^\top V^{-1}(\mathbf{Y} - \boldsymbol{\mu})/n = \mathbf{0}$$

correspond to the score equations under maximum likelihood estimation

Optimal Estimating Equations

These results lead to an interesting connection with GLMs:

- ▶ Regardless of where the Y_i 's live, we could think of a generative model

$$Y_i = \mu_i(\beta) + \epsilon_i,$$

which certainly holds for $\epsilon_i = Y_i - \mu_i(\beta)$, with $Y_i \mid \mathbf{x}_i \sim H_i[\mu_i(\beta)]$

- ▶ If the Y_i 's are independent given covariates, so are the ϵ_i 's
- ▶ $\text{var}(\epsilon_i \mid \mathbf{x}_i) = \text{var}(Y_i \mid \mathbf{x}_i)$ and say we assume $\text{var}(Y_i \mid \mathbf{x}_i) = \alpha V(\mu_i)$, where $\alpha V(\mu_i)$ coincides with one of the variances in a GLM
- ▶ So a GLM can be seen as a nonlinear model $Y_i = \mu_i(\beta) + \epsilon_i$ with a specific mean and variance structure
- ▶ The optimal estimating equations of Godambe and Heyde (1987)

$$\mathbf{U}_n(\beta) = \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{Y} - \boldsymbol{\mu})/n = \mathbf{0}$$

correspond to the score equations under maximum likelihood estimation

Reminder: Key Parts of Likelihood Inference for GLMs

When it comes down to obtaining an MLE and its asymptotic distribution in GLMs, the following are the key parts:

- ▶ The MLE $\hat{\beta}_n$ is obtained as the solution to

$$\mathbf{S}(\beta) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha = \mathbf{0},$$

which is derived from the log-likelihood, where

- ▶ $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ (assumption)
 - ▶ $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$ with $\mathbf{V} = \text{diag}\{V(\mu_i)/\phi_i\}$ (assumption)
 - ▶ \mathbf{D} is $n \times (k+1)$ with (i, j) th entry $\partial \mu_i / \partial \beta_j$, $i = 1, \dots, n$, $j = 0, \dots, k$
- ▶ The proof that the MLE $\hat{\beta}_n$ has asymptotic distribution

$$\mathcal{I}_n(\beta)^{1/2}(\hat{\beta}_n - \beta) \rightarrow_d N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$$

relies on $E[\mathbf{S}(\beta)] = \mathbf{0}$ and $-E[\mathbf{S}'(\beta)] = E[\mathbf{S}(\beta)\mathbf{S}(\beta)^T] = \mathcal{I}_n(\beta)$, see, e.g., Theorem 6.6 in Boos and Stefanski (2013)

Reminder: Key Parts of Likelihood Inference for GLMs

When it comes down to obtaining an MLE and its asymptotic distribution in GLMs, the following are the key parts:

- ▶ The MLE $\hat{\beta}_n$ is obtained as the solution to

$$\mathbf{S}(\beta) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha = \mathbf{0},$$

which is derived from the log-likelihood, where

- ▶ $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ (assumption)
- ▶ $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$ with $\mathbf{V} = \text{diag}\{V(\mu_i)/\phi_i\}$ (assumption)
- ▶ \mathbf{D} is $n \times (k+1)$ with (i, j) th entry $\partial \mu_i / \partial \beta_j$, $i = 1, \dots, n$, $j = 0, \dots, k$
- ▶ The proof that the MLE $\hat{\beta}_n$ has asymptotic distribution

$$\mathcal{I}_n(\beta)^{1/2}(\hat{\beta}_n - \beta) \rightarrow_d \mathbf{N}_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$$

relies on $E[\mathbf{S}(\beta)] = \mathbf{0}$ and $-E[\mathbf{S}'(\beta)] = E[\mathbf{S}(\beta)\mathbf{S}(\beta)^T] = \mathcal{I}_n(\beta)$, see, e.g., Theorem 6.6 in Boos and Stefanski (2013)

Reminder: Key Parts of Likelihood Inference for GLMs

In the context of GLMs:

- To obtain

$$E[\mathbf{S}(\boldsymbol{\beta})] = E\{\mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] / \alpha\} = \mathbf{0},$$

we need $\boldsymbol{\mu}(\boldsymbol{\beta}) = E(\mathbf{Y} \mid \mathbf{X})$

- To obtain $-E[\mathbf{S}'(\boldsymbol{\beta})] = E[\mathbf{S}(\boldsymbol{\beta})\mathbf{S}(\boldsymbol{\beta})^T]$, note

$$\begin{aligned} E[\mathbf{S}(\boldsymbol{\beta})\mathbf{S}(\boldsymbol{\beta})^T] &= \mathbf{D}^T \mathbf{V}^{-1} E[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^T] \mathbf{V}^{-1} \mathbf{D} / \alpha^2 \\ &= \alpha \mathbf{D}^T \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{D} / \alpha^2 \\ &= \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha \end{aligned}$$

which holds if $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$

Reminder: Key Parts of Likelihood Inference for GLMs

In the context of GLMs:

- ▶ To obtain

$$E[\mathbf{S}(\boldsymbol{\beta})] = E\{\mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] / \alpha\} = \mathbf{0},$$

we need $\boldsymbol{\mu}(\boldsymbol{\beta}) = E(\mathbf{Y} \mid \mathbf{X})$

- ▶ To obtain $-E[\mathbf{S}'(\boldsymbol{\beta})] = E[\mathbf{S}(\boldsymbol{\beta})\mathbf{S}(\boldsymbol{\beta})^T]$, note

$$\begin{aligned} E[\mathbf{S}(\boldsymbol{\beta})\mathbf{S}(\boldsymbol{\beta})^T] &= \mathbf{D}^T \mathbf{V}^{-1} E[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^T] \mathbf{V}^{-1} \mathbf{D} / \alpha^2 \\ &= \alpha \mathbf{D}^T \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{D} / \alpha^2 \\ &= \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha \end{aligned}$$

which holds if $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$

Reminder: Key Parts of Likelihood Inference for GLMs

► Similarly,

$$\begin{aligned} -E[S'(\beta)] &= -E \left[\frac{\partial}{\partial \beta^\top} \{ \mathbf{D}^\top \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha \} \right] \\ &= -\mathbf{D}^\top \mathbf{V}^{-1} E \left[\frac{\partial}{\partial \beta^\top} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] \right] / \alpha \\ &\quad - \frac{\partial}{\partial \beta^\top} (\mathbf{D}^\top \mathbf{V}^{-1}) E[\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha \\ &= \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D} / \alpha \end{aligned}$$

which depends on $\boldsymbol{\mu}(\beta) = E(\mathbf{Y} \mid \mathbf{X})$

Reminder: Key Parts of Likelihood Inference for GLMs

Key observations:

- ▶ The theory that supports

$$(\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha)^{1/2} (\hat{\beta}_n - \beta) \rightarrow_d N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$$

for $\hat{\beta}_n$, the solution to

$$\mathbf{S}(\beta) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha = \mathbf{0},$$

relies mainly on the assumptions that $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ and $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$

- ▶ We don't really rely on a specific form for $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ and $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$
- ▶ When we use the asymptotic distribution of the MLEs in GLMs, we *don't really rely on the full model specification*: we only use the implied mean and variance functions!

Reminder: Key Parts of Likelihood Inference for GLMs

Key observations:

- ▶ The theory that supports

$$(\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha)^{1/2} (\hat{\beta}_n - \beta) \rightarrow_d N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$$

for $\hat{\beta}_n$, the solution to

$$\mathbf{S}(\beta) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha = \mathbf{0},$$

relies mainly on the assumptions that $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ and $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$

- ▶ We don't really rely on a specific form for $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ and $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$
- ▶ When we use the asymptotic distribution of the MLEs in GLMs, we *don't really rely on the full model specification*: we only use the implied mean and variance functions!

Reminder: Key Parts of Likelihood Inference for GLMs

Key observations:

- ▶ The theory that supports

$$(\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha)^{1/2} (\hat{\beta}_n - \beta) \rightarrow_d N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$$

for $\hat{\beta}_n$, the solution to

$$\mathbf{S}(\beta) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\beta)] / \alpha = \mathbf{0},$$

relies mainly on the assumptions that $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ and $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$

- ▶ We don't really rely on a specific form for $E(\mathbf{Y} \mid \mathbf{X}) = \boldsymbol{\mu}(\beta)$ and $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$
- ▶ When we use the asymptotic distribution of the MLEs in GLMs, we *don't really rely on the full model specification*: we only use the implied mean and variance functions!

Quasi-Likelihood

Given the previous point of view, there is no reason to restrict yourself to assuming $\text{var}(Y_i | \mathbf{x}_i) = \alpha V(\mu_i)$ with $\alpha V(\mu_i)$ coming from one of the variances in a GLM: *quasi-likelihood* builds on this idea!

- ▶ Proposed by R. W. M. Wedderburn (1974, *Biometrika*)
- ▶ An alternative to MLE, when we do not wish to commit to specifying the full distribution of the data

Quasi-Likelihood

- ▶ Let Y_i , $i = 1, \dots, n$, have expected values μ_i and variances $\alpha V(\mu_i)$, where $V(\cdot)$ is a known function and $\alpha > 0$ is a scalar
- ▶ We assume only a structure for the mean and the variance:

$$\begin{aligned} E(\mathbf{Y} \mid \mathbf{X}) &= \boldsymbol{\mu}(\boldsymbol{\beta}) \\ \text{var}(\mathbf{Y} \mid \mathbf{X}) &= \alpha \mathbf{V}[\boldsymbol{\mu}(\boldsymbol{\beta})] \end{aligned}$$

where

- ▶ $\boldsymbol{\mu}(\boldsymbol{\beta}) = [\mu_1(\boldsymbol{\beta}), \dots, \mu_n(\boldsymbol{\beta})]^\top$, with $\mu_i(\boldsymbol{\beta}) = \mu(\mathbf{x}_i, \boldsymbol{\beta})$ representing the regression function
- ▶ $\mathbf{V} := \mathbf{V}[\boldsymbol{\mu}(\boldsymbol{\beta})] = \text{diag}\{V[\mu_i(\boldsymbol{\beta})]\}$ so that the observations are uncorrelated, and

$$\text{var}(Y_i \mid \mathbf{x}_i) = \alpha V[\mu_i(\boldsymbol{\beta})]$$

Quasi-Likelihood

- ▶ We use the estimating function or *quasi-score*:

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \mathbf{D}^T \mathbf{V}^{-1}(\mathbf{Y} - \boldsymbol{\mu})/\alpha$$

- ▶ This quasi-score is such that:

- ▶ $E[\tilde{\mathbf{S}}(\boldsymbol{\beta})] = \mathbf{0}$

- ▶ $\text{var}[\tilde{\mathbf{S}}(\boldsymbol{\beta})] = \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D}/\alpha$

- ▶ $-E\left[\frac{\partial \tilde{\mathbf{S}}}{\partial \boldsymbol{\beta}}\right] = \text{var}[\tilde{\mathbf{S}}(\boldsymbol{\beta})] = \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D}/\alpha$

under the quasi-likelihood mean and variance assumptions

- ▶ These properties are analogous to those obtained under maximum likelihood, thereby the name *quasi*-likelihood.

Quasi-Likelihood

- ▶ The word “quasi” also refers to the fact that the score may or may not correspond to a likelihood derived from a probability distribution
- ▶ For example, the variance function $V(\mu) = \mu^2(1 - \mu)^2$ does not correspond to a probability distribution (see, Wakefield (2013, p. 52))
- ▶ Model misspecification can also happen with quasi-likelihood!

Quasi-Likelihood

- ▶ The quasi-likelihood estimator satisfies

$$(\mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D} / \alpha)^{1/2} (\hat{\beta}_n - \beta) \rightarrow_d N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1}),$$

which is identical to what we obtain for GLMs under maximum likelihood

- ▶ The asymptotic covariance matrix of $\hat{\beta}_n$ is

$$\alpha (\mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D})^{-1},$$

which we use to obtain the estimator

$$\widehat{\text{var}}(\hat{\beta}_n) = \hat{\alpha} (\hat{\mathbf{D}}^\top \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}})^{-1},$$

where $\hat{\mathbf{D}} = \mathbf{D}(\hat{\beta}_n)$, $\hat{\mathbf{V}} = \mathbf{V}(\hat{\mu})$, $\hat{\mu} = \mu(\hat{\beta}_n)$, where $\hat{\mu}_i = \hat{\mu}_i(\hat{\beta}_n)$, and $\hat{\alpha}$ is an estimator of α

Quasi-Likelihood

- ▶ Since¹

$$E[(\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{V}^{-1}(\boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})] = n\alpha,$$

an unbiased estimator of α would be (with $\boldsymbol{\mu}$ known)

$$\hat{\alpha} = (\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{V}^{-1}(\boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})/n$$

- ▶ For diagonal \mathbf{V} and replacing $\boldsymbol{\mu}$ by $\hat{\boldsymbol{\mu}}$:

$$\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)}$$

- ▶ As usual, it is common to use a degrees-of-freedom-corrected (but not, in general, unbiased) estimator

$$\hat{\alpha} = \frac{1}{n - k - 1} \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)}$$

¹Reminder: suppose \mathbf{Z} is an $n \times 1$ random variable with $E[\mathbf{Z}] = \boldsymbol{\mu}$, $\text{var}(\mathbf{Z}) = \Sigma$ and \mathbf{A} is a symmetric $n \times n$ matrix. Then

$$E[\mathbf{Z}^\top \mathbf{A} \mathbf{Z}] = \text{tr}(\mathbf{A} \Sigma) + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}$$

Quasi-Likelihood and Overdispersion

- ▶ Quasi-likelihood is often introduced in the context of *overdispersion*: greater variability than expected under a given statistical model
- ▶ GLMs with $\alpha = 1$ in $\text{var}(Y_i | x_i) = \alpha V(\mu_i)$ are particularly vulnerable to mean-variance model misspecification, since they do not have a way of absorbing extra variability of the response
- ▶ For instance, it is common to find overdispersion with respect to the Poisson ($\text{var}(Y_i | x_i) = \mu_i$) and binomial ($\text{var}(Y_i | x_i) = n_i \mu_i (1 - \mu_i)$) GLMs, since their variances are entirely determined by their means
- ▶ In such cases, a simple implementation of quasi-likelihood corresponds to taking
 - ▶ $\text{var}(Y_i | x_i) = \alpha \mu_i$ for overdispersed Poisson data
 - ▶ $\text{var}(Y_i | x_i) = \alpha n_i \mu_i (1 - \mu_i)$ for overdispersed binomial data

Quasi-Likelihood and Overdispersion

- ▶ Quasi-likelihood is often introduced in the context of *overdispersion*: greater variability than expected under a given statistical model
- ▶ GLMs with $\alpha = 1$ in $\text{var}(Y_i | \mathbf{x}_i) = \alpha V(\mu_i)$ are particularly vulnerable to mean-variance model misspecification, since they do not have a way of absorbing extra variability of the response
- ▶ For instance, it is common to find overdispersion with respect to the Poisson ($\text{var}(Y_i | \mathbf{x}_i) = \mu_i$) and binomial ($\text{var}(Y_i | \mathbf{x}_i) = n_i \mu_i (1 - \mu_i)$) GLMs, since their variances are entirely determined by their means
- ▶ In such cases, a simple implementation of quasi-likelihood corresponds to taking
 - ▶ $\text{var}(Y_i | \mathbf{x}_i) = \alpha \mu_i$ for overdispersed Poisson data
 - ▶ $\text{var}(Y_i | \mathbf{x}_i) = \alpha n_i \mu_i (1 - \mu_i)$ for overdispersed binomial data

Quasi-Likelihood and Overdispersion

- ▶ Quasi-likelihood is often introduced in the context of *overdispersion*: greater variability than expected under a given statistical model
- ▶ GLMs with $\alpha = 1$ in $\text{var}(Y_i | \mathbf{x}_i) = \alpha V(\mu_i)$ are particularly vulnerable to mean-variance model misspecification, since they do not have a way of absorbing extra variability of the response
- ▶ For instance, it is common to find overdispersion with respect to the Poisson ($\text{var}(Y_i | \mathbf{x}_i) = \mu_i$) and binomial ($\text{var}(Y_i | \mathbf{x}_i) = n_i \mu_i (1 - \mu_i)$) GLMs, since their variances are entirely determined by their means
- ▶ In such cases, a simple implementation of quasi-likelihood corresponds to taking
 - ▶ $\text{var}(Y_i | \mathbf{x}_i) = \alpha \mu_i$ for overdispersed Poisson data
 - ▶ $\text{var}(Y_i | \mathbf{x}_i) = \alpha n_i \mu_i (1 - \mu_i)$ for overdispersed binomial data

Quasi-Likelihood and Overdispersion

- ▶ Quasi-likelihood is often introduced in the context of *overdispersion*: greater variability than expected under a given statistical model
- ▶ GLMs with $\alpha = 1$ in $\text{var}(Y_i | \mathbf{x}_i) = \alpha V(\mu_i)$ are particularly vulnerable to mean-variance model misspecification, since they do not have a way of absorbing extra variability of the response
- ▶ For instance, it is common to find overdispersion with respect to the Poisson ($\text{var}(Y_i | \mathbf{x}_i) = \mu_i$) and binomial ($\text{var}(Y_i | \mathbf{x}_i) = n_i \mu_i (1 - \mu_i)$) GLMs, since their variances are entirely determined by their means
- ▶ In such cases, a simple implementation of quasi-likelihood corresponds to taking
 - ▶ $\text{var}(Y_i | \mathbf{x}_i) = \alpha \mu_i$ for overdispersed Poisson data
 - ▶ $\text{var}(Y_i | \mathbf{x}_i) = \alpha n_i \mu_i (1 - \mu_i)$ for overdispersed binomial data

Sandwich Estimation from Misspecified GLMs

- ▶ The asymptotic variance-covariance matrices for $\hat{\beta}_n$ under likelihood and quasi-likelihood are appropriate only if the first two moments are correctly specified
- ▶ Let us take the score or quasi-score to be

$$\tilde{\mathbf{S}}(\beta) = \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{Y} - \boldsymbol{\mu}) \propto \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{Y} - \boldsymbol{\mu})/\alpha$$

where we ignore α as it cancels in $\tilde{\mathbf{S}}(\beta) = \mathbf{0}$

- ▶ In general, the asymptotic variance-covariance matrix for $\hat{\beta}_n$ is

$$\mathbf{A}_n^{-1} \mathbf{B}_n (\mathbf{A}_n^\top)^{-1} / n$$

Sandwich Estimation from Misspecified GLMs

- In general,

$$\begin{aligned}n\mathbf{A}_n &= \mathbb{E}[\tilde{\mathbf{S}}'(\boldsymbol{\beta})] = \mathbb{E}\left[\frac{\partial}{\partial\boldsymbol{\beta}^\top}\{\mathbf{D}^\top\mathbf{V}^{-1}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})]\}\right] \\&= \mathbf{D}^\top\mathbf{V}^{-1}\mathbb{E}\left[\frac{\partial}{\partial\boldsymbol{\beta}^\top}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})]\right] \\&\quad + \frac{\partial}{\partial\boldsymbol{\beta}^\top}(\mathbf{D}^\top\mathbf{V}^{-1})\mathbb{E}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] \\&= -\mathbf{D}^\top\mathbf{V}^{-1}\mathbf{D} + \frac{\partial}{\partial\boldsymbol{\beta}^\top}(\mathbf{D}^\top\mathbf{V}^{-1})\mathbb{E}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})]\end{aligned}$$

- If we assume $\boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbb{E}(\mathbf{Y} \mid \mathbf{X})$, we can take \mathbf{A}_n as

$$\mathbf{A}_n = \mathbf{D}^\top\mathbf{V}^{-1}\mathbf{D}/n$$

since the -1 cancels in the sandwich formula

Sandwich Estimation from Misspecified GLMs

- In general,

$$\begin{aligned}n\mathbf{A}_n &= \mathbb{E}[\tilde{\mathbf{S}}'(\beta)] = \mathbb{E}\left[\frac{\partial}{\partial\beta^\top}\{\mathbf{D}^\top\mathbf{V}^{-1}[\mathbf{Y} - \mu(\beta)]\}\right] \\&= \mathbf{D}^\top\mathbf{V}^{-1}\mathbb{E}\left[\frac{\partial}{\partial\beta^\top}[\mathbf{Y} - \mu(\beta)]\right] \\&\quad + \frac{\partial}{\partial\beta^\top}(\mathbf{D}^\top\mathbf{V}^{-1})\mathbb{E}[\mathbf{Y} - \mu(\beta)] \\&= -\mathbf{D}^\top\mathbf{V}^{-1}\mathbf{D} + \frac{\partial}{\partial\beta^\top}(\mathbf{D}^\top\mathbf{V}^{-1})\mathbb{E}[\mathbf{Y} - \mu(\beta)]\end{aligned}$$

- If we assume $\mu(\beta) = \mathbb{E}(\mathbf{Y} \mid \mathbf{X})$, we can take \mathbf{A}_n as

$$\mathbf{A}_n = \mathbf{D}^\top\mathbf{V}^{-1}\mathbf{D}/n$$

since the -1 cancels in the sandwich formula

Sandwich Estimation from Misspecified GLMs

Note that in:

$$n\mathbf{A}_n = \mathbb{E}[\tilde{\mathbf{S}}'(\boldsymbol{\beta})] = \mathbb{E}\left[\frac{\partial}{\partial\boldsymbol{\beta}^T}\{\mathbf{D}^T\mathbf{V}^{-1}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})]\}\right] = -\mathbf{D}^T\mathbf{V}^{-1}\mathbf{D} + \frac{\partial}{\partial\boldsymbol{\beta}^T}(\mathbf{D}^T\mathbf{V}^{-1})\mathbb{E}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})],$$

we have that

$$\frac{\partial}{\partial\boldsymbol{\beta}^T}(\mathbf{D}^T\mathbf{V}^{-1}) := \mathbf{Q}$$

is a $(k+1) \times (k+1) \times n$ tensor, so that $\mathbf{Q}[:, , i]$ is a $(k+1) \times (k+1)$ matrix

$$\mathbf{Q}[:, , i] = \frac{\partial}{\partial\boldsymbol{\beta}^T} \left(\frac{1}{V[\mu_i(\boldsymbol{\beta})]} \frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right) = \left[\frac{\partial}{\partial\beta_0} \left(\frac{1}{V[\mu_i(\boldsymbol{\beta})]} \frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right), \dots, \frac{\partial}{\partial\beta_k} \left(\frac{1}{V[\mu_i(\boldsymbol{\beta})]} \frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right) \right].$$

To see this, we can write down

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \mathbf{D}^T\mathbf{V}^{-1}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] = \sum_{i=1}^n \frac{[Y_i - \mu_i(\boldsymbol{\beta})]}{V[\mu_i(\boldsymbol{\beta})]} \left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right)$$

so that

$$\begin{aligned} \mathbb{E}[\tilde{\mathbf{S}}'(\boldsymbol{\beta})] &= -\sum_{i=1}^n \frac{1}{V[\mu_i(\boldsymbol{\beta})]} \left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right) \left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right)^T + \sum_{i=1}^n \mathbb{E}[Y_i - \mu_i(\boldsymbol{\beta})] \frac{\partial}{\partial\boldsymbol{\beta}^T} \left(\frac{1}{V[\mu_i(\boldsymbol{\beta})]} \frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}} \right) \\ &= -\mathbf{D}^T\mathbf{V}^{-1}\mathbf{D} + \frac{\partial}{\partial\boldsymbol{\beta}^T}(\mathbf{D}^T\mathbf{V}^{-1})\mathbb{E}[\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})]. \end{aligned}$$

Sandwich Estimation from Misspecified GLMs

- ▶ We can write down

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{[Y_i - \mu_i(\boldsymbol{\beta})]}{V[\mu_i(\boldsymbol{\beta})]} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right) := \sum_{i=1}^n \tilde{S}_i(\boldsymbol{\beta})$$

- ▶ In general,

$$\mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\tilde{S}_i(\boldsymbol{\beta}) \tilde{S}_i(\boldsymbol{\beta})^\top]$$

- ▶ If we assume $\boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbb{E}(\mathbf{Y} \mid \mathbf{X})$, since we have uncorrelated data, we can write

$$\mathbf{B}_n = \frac{1}{n} \mathbb{E}[\tilde{\mathbf{S}}(\boldsymbol{\beta}) \tilde{\mathbf{S}}(\boldsymbol{\beta})^\top] = \frac{1}{n} \mathbf{D}^\top \mathbf{V}^{-1} \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top] \mathbf{V}^{-1} \mathbf{D}$$

and

$$\text{var}(\mathbf{Y} \mid \mathbf{X}) = \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top]$$

- ▶ If we assume $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$ then

$$\mathbf{B}_n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{D} / n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D} / n$$

Sandwich Estimation from Misspecified GLMs

- ▶ We can write down

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{[Y_i - \mu_i(\boldsymbol{\beta})]}{V[\mu_i(\boldsymbol{\beta})]} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right) := \sum_{i=1}^n \tilde{S}_i(\boldsymbol{\beta})$$

- ▶ In general,

$$\mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\tilde{S}_i(\boldsymbol{\beta}) \tilde{S}_i(\boldsymbol{\beta})^\top]$$

- ▶ If we assume $\boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbb{E}(\mathbf{Y} \mid \mathbf{X})$, since we have uncorrelated data, we can write

$$\mathbf{B}_n = \frac{1}{n} \mathbb{E}[\tilde{\mathbf{S}}(\boldsymbol{\beta}) \tilde{\mathbf{S}}(\boldsymbol{\beta})^\top] = \frac{1}{n} \mathbf{D}^\top \mathbf{V}^{-1} \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top] \mathbf{V}^{-1} \mathbf{D}$$

and

$$\text{var}(\mathbf{Y} \mid \mathbf{X}) = \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top]$$

- ▶ If we assume $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$ then

$$\mathbf{B}_n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{D} / n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D} / n$$

Sandwich Estimation from Misspecified GLMs

- ▶ We can write down

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{[Y_i - \mu_i(\boldsymbol{\beta})]}{V[\mu_i(\boldsymbol{\beta})]} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right) := \sum_{i=1}^n \tilde{S}_i(\boldsymbol{\beta})$$

- ▶ In general,

$$\mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\tilde{S}_i(\boldsymbol{\beta}) \tilde{S}_i(\boldsymbol{\beta})^\top]$$

- ▶ If we assume $\boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbb{E}(\mathbf{Y} \mid \mathbf{X})$, since we have uncorrelated data, we can write

$$\mathbf{B}_n = \frac{1}{n} \mathbb{E}[\tilde{\mathbf{S}}(\boldsymbol{\beta}) \tilde{\mathbf{S}}(\boldsymbol{\beta})^\top] = \frac{1}{n} \mathbf{D}^\top \mathbf{V}^{-1} \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top] \mathbf{V}^{-1} \mathbf{D}$$

and

$$\text{var}(\mathbf{Y} \mid \mathbf{X}) = \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top]$$

- ▶ If we assume $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$ then

$$\mathbf{B}_n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{D} / n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D} / n$$

Sandwich Estimation from Misspecified GLMs

- ▶ We can write down

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{[Y_i - \mu_i(\boldsymbol{\beta})]}{V[\mu_i(\boldsymbol{\beta})]} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right) := \sum_{i=1}^n \tilde{S}_i(\boldsymbol{\beta})$$

- ▶ In general,

$$\mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\tilde{S}_i(\boldsymbol{\beta}) \tilde{S}_i(\boldsymbol{\beta})^\top]$$

- ▶ If we assume $\boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbb{E}(\mathbf{Y} \mid \mathbf{X})$, since we have uncorrelated data, we can write

$$\mathbf{B}_n = \frac{1}{n} \mathbb{E}[\tilde{\mathbf{S}}(\boldsymbol{\beta}) \tilde{\mathbf{S}}(\boldsymbol{\beta})^\top] = \frac{1}{n} \mathbf{D}^\top \mathbf{V}^{-1} \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top] \mathbf{V}^{-1} \mathbf{D}$$

and

$$\text{var}(\mathbf{Y} \mid \mathbf{X}) = \mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))(\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta}))^\top]$$

- ▶ If we assume $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \alpha \mathbf{V}$ then

$$\mathbf{B}_n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{D} / n = \alpha \mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D} / n$$

Sandwich Estimation from Misspecified GLMs

- ▶ Defining $\hat{\mathbf{A}}_n$ and $\hat{\mathbf{B}}_n$ by plugging $\hat{\beta}_n$ into the expressions for \mathbf{A}_n and \mathbf{B}_n leads to the sandwich estimator of the variance-covariance matrix for $\hat{\beta}_n$ as

$$\widehat{\text{var}}(\hat{\beta}_n) = \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n (\hat{\mathbf{A}}_n^{\top})^{-1} / n$$

- ▶ This sandwich estimator $\widehat{\text{var}}(\hat{\beta}_n)$ provides a consistent estimator of the variance $\text{var}(\hat{\beta}_n)$ and therefore asymptotically correct confidence interval coverage (as long as independence of responses holds)
- ▶ For more details see Kauermann and Carroll (2001, *JASA*)

Sandwich Estimation from Misspecified GLMs

- For instance, the most general version of $\hat{\mathbf{B}}_n$ for uncorrelated data is

$$\begin{aligned}\hat{\mathbf{B}}_n &= \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial \mu_i}{\partial \boldsymbol{\beta}} \right) \bigg|_{\hat{\boldsymbol{\beta}}_n} \frac{[Y_i - \mu_i(\hat{\boldsymbol{\beta}}_n)]^2}{(V[\mu_i(\hat{\boldsymbol{\beta}}_n)])^2} \left(\frac{\partial \mu_i}{\partial \boldsymbol{\beta}} \right) \bigg|_{\hat{\boldsymbol{\beta}}_n}^T \\ &= \frac{1}{n} \hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \text{diag}\{[Y_i - \mu_i(\hat{\boldsymbol{\beta}}_n)]^2\} \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}}\end{aligned}$$

- And the most general version of $\hat{\mathbf{A}}_n$ is

$$\hat{\mathbf{A}}_n = -\frac{1}{n} \hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}} + \frac{1}{n} \left[\frac{\partial}{\partial \boldsymbol{\beta}^T} (\mathbf{D}^T \mathbf{V}^{-1}) \right] \bigg|_{\hat{\boldsymbol{\beta}}_n} [\mathbf{Y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}_n)]$$

which simplifies to $\hat{\mathbf{A}}_n = -\hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}}/n$ if the mean model is correctly specified, which we often assume for interpretability (again, we could ignore the -1 in $\hat{\mathbf{A}}_n$ since it cancels in the sandwich formula)

Sandwich Estimation from Misspecified GLMs

- ▶ For instance, the most general version of $\hat{\mathbf{B}}_n$ for uncorrelated data is

$$\begin{aligned}\hat{\mathbf{B}}_n &= \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial \mu_i}{\partial \boldsymbol{\beta}} \right) \bigg|_{\hat{\boldsymbol{\beta}}_n} \frac{[Y_i - \mu_i(\hat{\boldsymbol{\beta}}_n)]^2}{(V[\mu_i(\hat{\boldsymbol{\beta}}_n)])^2} \left(\frac{\partial \mu_i}{\partial \boldsymbol{\beta}} \right) \bigg|_{\hat{\boldsymbol{\beta}}_n}^T \\ &= \frac{1}{n} \hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \text{diag}\{[Y_i - \mu_i(\hat{\boldsymbol{\beta}}_n)]^2\} \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}}\end{aligned}$$

- ▶ And the most general version of $\hat{\mathbf{A}}_n$ is

$$\hat{\mathbf{A}}_n = -\frac{1}{n} \hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}} + \frac{1}{n} \left[\frac{\partial}{\partial \boldsymbol{\beta}^T} (\mathbf{D}^T \mathbf{V}^{-1}) \right] \bigg|_{\hat{\boldsymbol{\beta}}_n} [\mathbf{Y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}_n)]$$

which simplifies to $\hat{\mathbf{A}}_n = -\hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}}/n$ if the mean model is correctly specified, which we often assume for interpretability (again, we could ignore the -1 in $\hat{\mathbf{A}}_n$ since it cancels in the sandwich formula)

Final Comments

These results are very powerful!

- ▶ Say you want to fit a regression model where you assume $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = E(Y_i | \mathbf{x}_i)$
- ▶ The optimality theory of Godambe and Heyde (1987) tells you that, for the sake of efficiency, you should use as your estimating function the quasi-score:

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] = \mathbf{0},$$

where $E(\mathbf{Y} | \mathbf{X}) = \boldsymbol{\mu}(\boldsymbol{\beta})$ and $\text{var}(\mathbf{Y} | \mathbf{X}) = \alpha \mathbf{V}$

- ▶ However, the theory of estimating equations has your back in case you misspecify $\text{var}(\mathbf{Y} | \mathbf{X})$!
- ▶ So, you can choose a *working variance* $\alpha \mathbf{V}$ to specify your estimating function, but use sandwich estimators, just in case you are wrong!
- ▶ A more general version of this idea is heavily used in *generalized estimating equations* (STAT/BIOST 571)

Final Comments

These results are very powerful!

- ▶ Say you want to fit a regression model where you assume $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = E(Y_i | \mathbf{x}_i)$
- ▶ The optimality theory of Godambe and Heyde (1987) tells you that, for the sake of efficiency, you should use as your estimating function the quasi-score:

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] = \mathbf{0},$$

where $E(\mathbf{Y} | \mathbf{X}) = \boldsymbol{\mu}(\boldsymbol{\beta})$ and $\text{var}(\mathbf{Y} | \mathbf{X}) = \alpha \mathbf{V}$

- ▶ However, the theory of estimating equations has your back in case you misspecify $\text{var}(\mathbf{Y} | \mathbf{X})$!
- ▶ So, you can choose a *working variance* $\alpha \mathbf{V}$ to specify your estimating function, but use sandwich estimators, just in case you are wrong!
- ▶ A more general version of this idea is heavily used in *generalized estimating equations* (STAT/BIOST 571)

Final Comments

These results are very powerful!

- ▶ Say you want to fit a regression model where you assume $\mu(\mathbf{x}_i, \boldsymbol{\beta}) = E(Y_i | \mathbf{x}_i)$
- ▶ The optimality theory of Godambe and Heyde (1987) tells you that, for the sake of efficiency, you should use as your estimating function the quasi-score:

$$\tilde{\mathbf{S}}(\boldsymbol{\beta}) = \mathbf{D}^T \mathbf{V}^{-1} [\mathbf{Y} - \boldsymbol{\mu}(\boldsymbol{\beta})] = \mathbf{0},$$

where $E(\mathbf{Y} | \mathbf{X}) = \boldsymbol{\mu}(\boldsymbol{\beta})$ and $\text{var}(\mathbf{Y} | \mathbf{X}) = \alpha \mathbf{V}$

- ▶ However, the theory of estimating equations has your back in case you misspecify $\text{var}(\mathbf{Y} | \mathbf{X})$!
- ▶ So, you can choose a *working variance* $\alpha \mathbf{V}$ to specify your estimating function, but use sandwich estimators, just in case you are wrong!
- ▶ A more general version of this idea is heavily used in *generalized estimating equations* (STAT/BIOST 571)