

Advanced Regression Methods for Independent Data

STAT/BIOST 570, 2020

Regression Models with Weaker Assumptions

Mauricio Sadinle

Department of Biostatistics

W UNIVERSITY *of* WASHINGTON

Regression Models with Weaker Assumptions

We now study certain types of inference under weaker sets of assumptions:

- ▶ Least squares without assuming a fully parametric model for $Y \mid x$
- ▶ Assuming only a regression and variance function for $Y \mid x$
- ▶ Assuming only a regression function for $Y \mid x$
- ▶ What happens with fully parametric models when they are wrong?

We first study some key results from *estimating equations*¹

¹See Chapter 7 of Boos and Stefanski (2013) *Essential Statistical Inference*, Springer.

Estimating Equations

Estimating equations provide a very general framework for deriving estimators with desirable properties

- ▶ We assume that Y_i , $i = 1, \dots, n$, are i.i.d. from distribution F
- ▶ We seek to estimate

$$\theta_0 : \quad E_F[\mathbf{G}(\theta_0, Y)] = 0,$$

where

- ▶ $Y \sim F$: generic random variable (vector)
- ▶ θ : vector of p parameters
- ▶ \mathbf{G} : p -dimensional, continuously differentiable function of θ and Y

Estimating Equations

$$\mathbf{G}(\boldsymbol{\theta}, Y) = \begin{pmatrix} G_1(\boldsymbol{\theta}, Y) \\ \vdots \\ G_p(\boldsymbol{\theta}, Y) \end{pmatrix} : \Theta \times \mathcal{Y} \rightarrow \mathbb{R}^p$$

Estimating Equations

- ▶ Given the data $\mathbf{Y} = (Y_1, \dots, Y_n)$, an *estimating function* is defined as

$$\mathbf{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\boldsymbol{\theta}, Y_i)$$

- ▶ If we have that

$$\mathbb{E}[\mathbf{G}_n(\boldsymbol{\theta}_0)] = \mathbf{0},$$

we refer to $\mathbf{G}_n(\boldsymbol{\theta})$ as an *unbiased* estimating function, as it's unbiased for estimating $\mathbf{0}$

Estimating Equations

- ▶ The corresponding *estimating equation* is given by

$$\mathbf{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\boldsymbol{\theta}, Y_i) = \mathbf{0}$$

- ▶ The solution to this system of equations

$$\hat{\boldsymbol{\theta}}_n : \quad \mathbf{G}_n(\hat{\boldsymbol{\theta}}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\hat{\boldsymbol{\theta}}_n, Y_i) = \mathbf{0}$$

is often referred to as *M-estimator* or *Z-estimator*

M-Estimators

M-estimator is the term used when the above formulation is derived from a **M**aximization problem

- ▶ For known function h , say we seek to maximize

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n h(\theta, Y_i)$$

- ▶ Often we find the maximizer by taking derivatives and setting them equal to zero

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial h(\theta, Y_i)}{\partial \theta} = \mathbf{0}$$

in which case $\mathbf{G}(\theta, Y_i) = \partial h(\theta, Y_i) / \partial \theta$

- ▶ Examples: maximum likelihood estimation: $h(\theta, Y_i) = \log p(Y_i | \theta)$

Z-Estimators

- ▶ However, the theory of estimating equations does not care about where the estimating equation comes from
- ▶ *Z-estimator* is the term used more generally, as the estimator is derived as a **Z**ero or solution to an equation
- ▶ Example: method of moments:

$$\mathbf{G}(\boldsymbol{\theta}, Y_i) = \begin{pmatrix} Y_i - \mathbb{E}[Y_i | \boldsymbol{\theta}] \\ Y_i^2 - \mathbb{E}[Y_i^2 | \boldsymbol{\theta}] \\ \vdots \\ Y_i^p - \mathbb{E}[Y_i^p | \boldsymbol{\theta}] \end{pmatrix}$$

when $\boldsymbol{\theta}$ is p -dimensional

Comments on Parameters

- ▶ A parameter θ is a characteristic of a probability distribution F (a super-population)
- ▶ If I give you F , then you would know how to compute θ (in principle)
- ▶ A parameter θ can be seen as a functional that maps distributions to, say, the real line, and we often write $\theta := \theta(F)$
 - ▶ The median is defined as (similarly for other quantiles)

$$m : F\{(-\infty, m]\} = 1/2$$

- ▶ If F has a pdf $f(\cdot)$, then

$$\int h(y)f(y) dy$$

is a parameter for any integrable function h (in particular, moments)

Comments on Parameters

- ▶ In parametric models, a parameter θ appears explicitly in the functional form of the pdf $f(\cdot)$ of F
- ▶ With estimating equations, however, we can define parameters based on

$$\theta_0 : \quad E_F[\mathbf{G}(\theta_0, Y)] = 0,$$

where E_F means that the expected value is taken w.r.t. Y , where $Y \sim F$

- ▶ Ultimately, we can think that θ_0 is a function of F , i.e., $\theta_0 := \theta_0(F)$
- ▶ *The power of estimating functions*: we can estimate θ_0 without having to fully specify a model for F

Consistency of Z-estimators

Result (Consistency of Z-estimators): Suppose that $\hat{\theta}_n$ is a solution to the estimating equation $\mathbf{G}_n(\theta) = \mathbf{0}$, i.e. $\mathbf{G}_n(\hat{\theta}_n) = \mathbf{0}$. Then $\hat{\theta}_n \rightarrow_p \theta_0$.

See Theorem 7.1 of Boos and Stefanski (2013), similar to the proof of consistency of MLEs.

Asymptotic Normality of Z-Estimators

Result (Asymptotic Normality of Z-estimators): Suppose that $\hat{\theta}_n$ is a solution to the estimating equation $\mathbf{G}_n(\theta) = \mathbf{0}$, i.e. $\mathbf{G}_n(\hat{\theta}_n) = \mathbf{0}$. Then

$$\sqrt{n} (\hat{\theta}_n - \theta_0) \rightarrow_d N_p[\mathbf{0}, \mathbf{A}^{-1} \mathbf{B} (\mathbf{A}^{-1})^\top]$$

where $\mathbf{A} = \mathbf{A}(\theta_0)$ with

$$\mathbf{A}(\theta) = E[\mathbf{G}'(\theta, Y)] = E\left\{\left[\frac{\partial}{\partial \theta_1} \mathbf{G}(\theta, Y), \dots, \frac{\partial}{\partial \theta_p} \mathbf{G}(\theta, Y)\right]_{p \times p}\right\}$$

and

$$\mathbf{B} = \mathbf{B}(\theta_0) = E[\mathbf{G}(\theta_0, Y) \mathbf{G}(\theta_0, Y)^\top] = \text{cov}\{\mathbf{G}(\theta_0, Y)\}.$$

Proof outline: derive asymptotic properties of the estimating function, and then transfer these to the estimator.

Remark: Boos and Stefanski (2013) define $\mathbf{A}(\theta) = E[-\mathbf{G}'(\theta, Y)]$, but the minus cancels in $\mathbf{A}^{-1} \mathbf{B} (\mathbf{A}^{-1})^\top$

Asymptotic Normality of Z-Estimators

Proof outline: From a Taylor series expansion around θ_0 we get:

$$\mathbf{0} = \mathbf{G}_n(\hat{\theta}_n) = \mathbf{G}_n(\theta_0) + \mathbf{G}'_n(\theta_0)(\hat{\theta}_n - \theta_0) + \mathbf{R}_n$$

which implies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = [-\mathbf{G}'_n(\theta_0)]^{-1} \sqrt{n} \mathbf{G}_n(\theta_0) + \sqrt{n} \mathbf{R}_n^*$$

As $n \rightarrow \infty$, by the WLLN,

$$\mathbf{G}'_n(\theta_0) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}'(\theta_0, Y_i) \xrightarrow{p} \mathbb{E} [\mathbf{G}'(\theta_0, Y)] = \mathbf{A}(\theta_0)$$

and by the CLT,

$$\sqrt{n} \mathbf{G}_n(\theta_0) \xrightarrow{d} N_p[\mathbf{0}, \mathbf{B}(\theta_0)]$$

where $\mathbf{B}(\theta_0) = \text{cov}\{\mathbf{G}(\theta_0, Y)\} = \mathbb{E}[\mathbf{G}(\theta_0, Y) \mathbf{G}(\theta_0, Y)^T]$.

Showing conditions under which $\sqrt{n} \mathbf{R}_n^* = o_p(1)$ is 580's material.

Paste these results with the help of Slutsky and you get the result.

Comments on Estimating Functions

This result is *very* important!

- ▶ Using only the definition of θ_0 and $\hat{\theta}_n$, it tells us that as $n \rightarrow \infty$, $\hat{\theta}_n$ is approximately normally distributed, around the 'right' mean θ_0
- ▶ The covariance $\text{cov}(\hat{\theta}_n) \approx \mathbf{A}^{-1} \mathbf{B} (\mathbf{A}^{-1})^\top / n$ is known as the *sandwich formula*, and it goes to $\mathbf{0}_{p \times p}$ as $n \rightarrow \infty$

Comments on Estimating Functions

This result is *very* important!

- ▶ It holds under very mild conditions; in particular, there is no assumed parametric model, and therefore no likelihood
- ▶ Whether asymptotic normality is accurate depends on the accuracy of the Central Limit Theorem: in many cases OK for n in the hundreds
- ▶ The $\mathbf{A}(\boldsymbol{\theta})$ and $\mathbf{B}(\boldsymbol{\theta})$ matrices are expectations over unknown distribution F . Different ways of evaluating these expectations lead to different estimating approaches!
- ▶ Note that $\boldsymbol{\theta}$ is unknown but we know that $\hat{\boldsymbol{\theta}}_n$ is consistent for $\boldsymbol{\theta}_0$ (more on this later)

Example: Mean

- ▶ Taking $\mathbf{G}(\theta, Y) = Y - \theta$, we obtain

$$\mathbb{E}_F[\mathbf{G}(\theta_0, Y)] = 0 \implies \mathbb{E}_F(Y) = \theta_0$$

so, for any distribution F , θ_0 will represent the mean

- ▶ The Z-estimator based on $\{Y_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} F$ is obtained from

$$\mathbf{G}_n(\theta, Y) = \frac{1}{n} \sum_{i=1}^n (Y_i - \theta) = 0 \implies \hat{\theta} = \bar{Y}$$

- ▶ On the other hand,

$$\mathbf{A}(\theta) = \mathbb{E}_F [\mathbf{G}'(\theta, Y)] = -1$$

and

$$\mathbf{B}(\theta) = \mathbb{E}_F[\mathbf{G}(\theta, Y)\mathbf{G}(\theta, Y)^\top] = \text{var}_F(Y)$$

- ▶ We finally have that

$$\bar{Y} \approx N[\mathbb{E}_F(Y), \text{var}_F(Y)/n]$$

Example: Mean and Variance

- ▶ Taking

$$E_F[\mathbf{G}(\boldsymbol{\theta}, Y)] = E_F \begin{pmatrix} Y - \theta_1 \\ (Y - \theta_1)^2 - \theta_2 \end{pmatrix} = \mathbf{0} \implies \begin{aligned} \theta_1 &= E_F(Y) \\ \theta_2 &= \text{var}_F(Y) \end{aligned}$$

- ▶ Based on this definition of $\mathbf{G}(\boldsymbol{\theta}, Y)$ we obtain

$$\mathbf{A} = E_F[\mathbf{G}'(\boldsymbol{\theta}, Y)] = E_F \begin{pmatrix} -1 & 0 \\ -2(Y - \theta_1) & -1 \end{pmatrix} = -\mathbf{I}_2$$

and

$$\mathbf{B} = E_F[\mathbf{G}(\boldsymbol{\theta}, Y)\mathbf{G}(\boldsymbol{\theta}, Y)^\top] = \begin{pmatrix} \theta_2 & E_F[(Y - \theta_1)^3] \\ E_F[(Y - \theta_1)^3] & E_F[(Y - \theta_1)^4] - \theta_2^2 \end{pmatrix}$$

- ▶ The Z -estimators of (θ_1, θ_2) are the sample mean and variance (\bar{Y}, s_n^2) , and

$$\begin{pmatrix} \bar{Y} \\ s_n^2 \end{pmatrix} \approx N \left(\begin{pmatrix} E_F(Y) \\ \text{var}_F(Y) \end{pmatrix}, \mathbf{B}/n \right)$$

- ▶ See Sec 7.2.2 of Boos and Stefanski (2013) for more details

Estimating Equations and Likelihood-Based Inference

The asymptotic distribution of MLEs can be derived as a particular case of the asymptotic distribution of Z -estimators

- ▶ Take the score as the basis for the estimating function

$$\mathbf{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\boldsymbol{\theta}, Y_i) = \frac{1}{n} \mathbf{S}(\boldsymbol{\theta}),$$

with $\mathbf{G}(\boldsymbol{\theta}, Y_i) = \partial \log p(Y_i \mid \boldsymbol{\theta}) / \partial \boldsymbol{\theta}$

- ▶ We know that, *under the model that leads to the likelihood*,

$$\mathbb{E}[\mathbf{G}_n(\boldsymbol{\theta})] = \frac{1}{n} \mathbb{E}[\mathbf{S}(\boldsymbol{\theta})] = \mathbf{0}$$

- ▶ And the MLE satisfies $\mathbf{G}_n(\hat{\boldsymbol{\theta}}_n) = \mathbf{0}$

Estimating Equations and Likelihood-Based Inference

Reminder. The proof of $E[\mathbf{S}(\boldsymbol{\theta})] = \mathbf{0}$:

$$\begin{aligned} E[\mathbf{S}(\boldsymbol{\theta})] &= E\left[\frac{\partial l}{\partial \boldsymbol{\theta}}\right] = nE\left[\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta})\right] \\ &= n \int \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log p(y | \boldsymbol{\theta})\right] p(y | \boldsymbol{\theta}) dy \quad (\text{if } Y \sim p(\cdot | \boldsymbol{\theta})) \\ &= n \int \left[\frac{\partial}{\partial \boldsymbol{\theta}} p(y | \boldsymbol{\theta})\right] \frac{p(y | \boldsymbol{\theta})}{p(y | \boldsymbol{\theta})} dy \\ &= n \frac{\partial}{\partial \boldsymbol{\theta}} \int p(y | \boldsymbol{\theta}) dy \quad (\text{if OK to exchange } \int \text{ and } \partial/\partial \boldsymbol{\theta}) \\ &= \mathbf{0}. \end{aligned}$$

So this result, in general, relies on your model being correctly specified,
i.e. $Y \sim p(\cdot | \boldsymbol{\theta})$

Estimating Equations and Likelihood-Based Inference

Taking $\mathbf{G}_n(\boldsymbol{\theta}) = \mathbf{S}(\boldsymbol{\theta})/n$, the \mathbf{A} and \mathbf{B} in the asymptotic covariance of Z -estimators take the following form:

$$\mathbf{A}(\boldsymbol{\theta}) = \text{E} [\mathbf{G}'(\boldsymbol{\theta}, Y)] = \text{E} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^\top \partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right],$$

and

$$\begin{aligned} \mathbf{B}(\boldsymbol{\theta}) &= \text{E} [\mathbf{G}(\boldsymbol{\theta}, Y) \mathbf{G}(\boldsymbol{\theta}, Y)^\top] \\ &= \text{E} \left[\left(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right) \left(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right)^\top \right], \end{aligned}$$

- ▶ From the definition of Fisher information $\mathcal{I}_1(\boldsymbol{\theta}) = \mathbf{B}(\boldsymbol{\theta})$
- ▶ Under the model that leads to the likelihood we obtain $\mathcal{I}_1(\boldsymbol{\theta}) = -\mathbf{A}(\boldsymbol{\theta})$

Estimating Equations and Likelihood-Based Inference

- Therefore, *if the model is correctly specified*, the asymptotic variance of the MLE obtained as a Z -estimator is

$$\mathbf{A}^{-1} \mathbf{B}(\mathbf{A}^{-1})^\top / n = \mathcal{I}_1(\boldsymbol{\theta})^{-1} / n$$

- Hence, *under the model*

$$\sqrt{n} (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \rightarrow_d N_p[\mathbf{0}, \mathcal{I}_1(\boldsymbol{\theta})^{-1}]$$

- This result relies on $\mathcal{I}_1(\boldsymbol{\theta}) = \mathbf{B}(\boldsymbol{\theta}) = -\mathbf{A}(\boldsymbol{\theta})$, which *holds under the parametric model leading to the likelihood*

Estimating Equations and Likelihood-Based Inference

Reminder. The proof² that

$$-E \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^\top \partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right] = E \left[\left(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right) \left(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right)^\top \right],$$

relies on being able to write down

$$E \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log p(y | \boldsymbol{\theta}) \right] = \int \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log p(y | \boldsymbol{\theta}) \right] p(y | \boldsymbol{\theta}) dy,$$

which is valid if $Y \sim p(\cdot | \boldsymbol{\theta})$

- ▶ So this result, in general, relies on your model being correctly specified, i.e. $Y \sim p(\cdot | \boldsymbol{\theta})$
- ▶ Nevertheless, the theory of estimating equations will allow us to study the asymptotic behavior of MLEs under model misspecification, among many other things!

²See, e.g., Theorems 6.6 and 6.7 of Boos and Stefanski (2013), or p. 38 of Wakefield(2013)

Looking Forward

The theory of estimating equations will allow us to study:

- ▶ What happens with fully parametric models when they are wrong?
- ▶ Assuming only a regression and variance function for $Y \mid x$
- ▶ Least squares without assuming a fully parametric model for $Y \mid x$
- ▶ Assuming only a regression function for $Y \mid x$

Behavior of MLEs Under Misspecification

In practice we can *compute* MLEs *regardless* of whether our model is correct – what are we estimating if the model is wrong?

- ▶ Let F be the true distribution, with density $f(\cdot)$
- ▶ Let $p(y | \theta)$ denote the density function of the *assumed* model
- ▶ In practice, we wouldn't know that the data come from F , so we still maximize the log-likelihood (dividing by n doesn't change the maximizer)

$$\frac{1}{n} \sum_{i=1}^n \log p(Y_i | \theta)$$

but now $\{Y_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} F$

Behavior of MLEs Under Misspecification

- In the limit we have that

$$\frac{1}{n} \sum_{i=1}^n \log p(Y_i | \boldsymbol{\theta}) \rightarrow_{a.s.} E_F[\log p(Y | \boldsymbol{\theta})] = \int \log p(y | \boldsymbol{\theta}) f(y) dy,$$

- We also obtain that $\hat{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}_T$, where

$$\boldsymbol{\theta}_T = \arg \max_{\boldsymbol{\theta}} E_F [\log p(Y | \boldsymbol{\theta})]$$

- We shall refer to $\boldsymbol{\theta}_T$ as the *pseudo-true* parameter, as this is the parameter that we will recover asymptotically

Behavior of MLEs Under Misspecification

- Now, note that we can write

$$\begin{aligned} E_F[\log p(Y | \theta)] &= E_F[\log f(Y) - \log f(Y) + \log p(Y | \theta)] \\ &= E_F[\log f(Y)] - \text{KL}[f(\cdot), p(\cdot | \theta)] \end{aligned}$$

where KL is the *Kullback-Leibler divergence*

$$\text{KL}[f(\cdot), p(\cdot | \theta)] = \int f(y) \log \frac{f(y)}{p(y | \theta)} dy$$

- Therefore,

$$\theta_T = \arg \min_{\theta} \text{KL}[f(\cdot), p(\cdot | \theta)],$$

that is, the MLE asymptotically minimizes KL as a function of θ

- The MLE is that value of θ that makes the *assumed model* diverge the least from the true distribution F

Behavior of MLEs Under Misspecification

- ▶ Taking derivatives of $E_F [\log p(Y | \theta)]$ and interchanging with the integral, we find that θ_T solves the system of equations at the population level:

$$\theta_T : E_F \left[\frac{\partial}{\partial \theta} \log p(Y | \theta) \right] \Big|_{\theta=\theta_T} = \mathbf{0}$$

- ▶ We can use this as the basis for treating $\hat{\theta}_n$ as a Z -estimator for θ_T

Behavior of MLEs Under Misspecification

Result: Suppose that $\hat{\boldsymbol{\theta}}_n$ is a solution to the estimating equation

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y_i | \boldsymbol{\theta}) = \mathbf{0}.$$

Then again

$$\sqrt{n} (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_T) \rightarrow_d N_p[\mathbf{0}, \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{\top -1}]$$

where

$$\mathbf{A} = \mathbf{A}(\boldsymbol{\theta}_T) = E_T \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^{\top} \partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right] \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_T}$$

and

$$\mathbf{B} = \mathbf{B}(\boldsymbol{\theta}_T) = E_T \left[\left(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right) \left(\frac{\partial}{\partial \boldsymbol{\theta}} \log p(Y | \boldsymbol{\theta}) \right)^{\top} \right] \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_T}$$

Behavior of MLEs Under Misspecification

- ▶ The above result characterizes the behavior of MLEs under model misspecification
- ▶ It can be used for inferences on the pseudo-true parameters θ_T
- ▶ However, θ_T is *in general* difficult to interpret
 - ▶ Huber (1967) presented these results from a purely mathematical point of view: what is the MLE's asymptotic distribution if the model is misspecified?
 - ▶ David A. Freedman (2006)³ criticizes using these results in practice: if your model is wrong, how do you know you care about θ_T ?

³On The So-Called “Huber Sandwich Estimator” and “Robust Standard Errors”
The American Statistician

Behavior of MLEs Under Misspecification

- ▶ With parametric models, we assume that the true distribution F has a density that can be written as $p(y \mid \theta)$ for some $\theta \in \Theta$
- ▶ Ideally, we choose this model based on knowledge of the data, and therefore the interpretation of θ is given when you specify the model
- ▶ On the other hand, reality is complicated, and so F may not have a density that can be written as $p(y \mid \theta)$ for some $\theta \in \Theta$
- ▶ We saw that the MLE estimates the value of θ_T that makes the assumed model “least divergent” from the true distribution
- ▶ The interpretation of θ_T can also be obtained from reorganizing:

$$\theta_T : \quad E_F \left[\frac{\partial}{\partial \theta} \log p(Y \mid \theta) \right] \Big|_{\theta=\theta_T} = \mathbf{0}$$

- ▶ This may or may not lead to an interpretation of θ that matches what you intended when you specified your model

Example: Misspecified Exponential Dispersion Family

- ▶ The generic term of the log-likelihood *obtained* from the exponential dispersion family is given by

$$\log p(Y | \theta, \alpha) = \frac{Y\theta - b(\theta)}{\alpha} + c(Y, \alpha)$$

- ▶ From this,

$$\frac{\partial}{\partial \theta} \log p(Y | \theta, \alpha) = \frac{Y - b'(\theta)}{\alpha}$$

- ▶ Using this in an estimating equation, with $Y \sim F$, leads to

$$E_F \left[\frac{\partial}{\partial \theta} \log p(Y | \theta, \alpha) \right] \Big|_{\theta=\theta_T} = 0 \implies E_F(Y) = b'(\theta_T)$$

- ▶ Therefore, using the exponential dispersion family to form the likelihood based on data $\{Y_i\}_{i=1}^n \overset{i.i.d.}{\sim} F$ leads to $b'(\hat{\theta})$ as a consistent estimator for $E_F(Y)$, regardless of whether the model is misspecified!

Example: Misspecified Exponential Dispersion Family

- ▶ We can also find that

$$\frac{\partial}{\partial \alpha} \log p(Y | \theta, \alpha) = -\frac{Y\theta - b(\theta)}{\alpha^2} + \frac{\partial}{\partial \alpha} c(Y, \alpha)$$

- ▶ Adding this to the previous estimating equation, with $Y \sim F$, leads to α_T being the value that solves

$$\theta_T E_F(Y) - b(\theta_T) = \alpha^2 E_F \left[\frac{\partial}{\partial \alpha} c(Y, \alpha) \right]$$

- ▶ Therefore, in general it is not clear whether α_T would be of interest when the model is misspecified

Example: Misspecified Weibull Model

- ▶ Jon A. Wellner in his Breslow Lecture (November 12, 2020)⁴ presented the example of a misspecified Weibull(α, β) model, and he showed that the MLE of α can be interpretable even if the model is misspecified, whereas the MLE of β is more difficult to interpret.
- ▶ *Problem:* Explain where the MLE of α and β converge to if the Weibull model is misspecified.

⁴[https:](https://sites.stat.washington.edu/jaw/RESEARCH/TALKS/Breslow-Lecture.pdf)

[//sites.stat.washington.edu/jaw/RESEARCH/TALKS/Breslow-Lecture.pdf](https://sites.stat.washington.edu/jaw/RESEARCH/TALKS/Breslow-Lecture.pdf)

Sandwich Estimation

- ▶ The Z-estimator $\hat{\theta}_n$ that solves

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

based on independent and identically distributed observations, has asymptotic variance $\mathbf{A}^{-1} \mathbf{B} (\mathbf{A}^{-1})^\top / n$

- ▶ We use this to define the *sandwich estimator* of $\text{var}(\hat{\theta}_n)$ as

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

where

$$\hat{\mathbf{A}} := \hat{\mathbf{A}}(\hat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}'(\hat{\theta}_n, Y_i),$$

and

$$\hat{\mathbf{B}} := \hat{\mathbf{B}}(\hat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\hat{\theta}_n, Y_i) \mathbf{G}(\hat{\theta}_n, Y_i)^\top.$$

Sandwich Estimation

- ▶ Note that by the weak law of large numbers $\hat{\mathbf{A}}(\boldsymbol{\theta}) \rightarrow_p \mathbf{A}(\boldsymbol{\theta})$ and $\hat{\mathbf{B}}(\boldsymbol{\theta}) \rightarrow_p \mathbf{B}(\boldsymbol{\theta})$, for any fixed $\boldsymbol{\theta}$
- ▶ However, our estimators are $\hat{\mathbf{A}} := \hat{\mathbf{A}}(\hat{\boldsymbol{\theta}}_n)$ and $\hat{\mathbf{B}} := \hat{\mathbf{B}}(\hat{\boldsymbol{\theta}}_n)$, which depend on the Z- estimator $\hat{\boldsymbol{\theta}}_n$
- ▶ Nevertheless, since $\hat{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}_0$, Boos and Stefanski (2013) present Theorems 7.3 and 7.4, which guarantee

$$\hat{\mathbf{A}}(\hat{\boldsymbol{\theta}}_n) \rightarrow_p \mathbf{A}(\boldsymbol{\theta}_0) := \mathbf{A}, \quad \hat{\mathbf{B}}(\hat{\boldsymbol{\theta}}_n) \rightarrow_p \mathbf{B}(\boldsymbol{\theta}_0) := \mathbf{B},$$

and therefore

$$\hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} (\hat{\mathbf{A}}^\top)^{-1} \rightarrow_p \mathbf{A}^{-1} \mathbf{B} (\mathbf{A}^\top)^{-1}$$

Sandwich Estimation and Robust Standard Errors

- ▶ The sandwich estimator provides a consistent estimator of the variance-covariance of Z -estimators in very broad situations
- ▶ For small sample sizes, the sandwich estimator may be unreliable, as it builds on an asymptotic argument: model-based estimators may be preferable for small to medium sized n , but you need to trust your model!
- ▶ The sandwich estimator is often called *robust* (people talk about *robust standard errors*), meaning *robustness to model departures*

Sandwich Estimation

How about when $\hat{\theta}_n$ is the MLE of a misspecified model?

- ▶ The estimating function arises from the score from a likelihood function:

$$\mathbf{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} l_i(\boldsymbol{\theta}),$$

with $l_i(\boldsymbol{\theta}) = \log p(Y_i | \boldsymbol{\theta})$

- ▶ Plugging into the formulae above, the sandwich estimator is based on

$$\hat{\mathbf{A}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta}^\top \partial \boldsymbol{\theta}} l_i(\boldsymbol{\theta}) \Big|_{\hat{\boldsymbol{\theta}}_n}$$

and

$$\hat{\mathbf{B}} = \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial}{\partial \boldsymbol{\theta}} l_i(\boldsymbol{\theta}) \right) \left(\frac{\partial}{\partial \boldsymbol{\theta}} l_i(\boldsymbol{\theta}) \right)^\top \Big|_{\hat{\boldsymbol{\theta}}_n}$$

- ▶ Inferences on the pseudo-true parameter $\boldsymbol{\theta}_T$ can be based on the sandwich variance estimator $\widehat{\text{var}}(\hat{\boldsymbol{\theta}}_n) = \hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} (\hat{\mathbf{A}}^{-1})^\top / n$

Regression and Estimating Equations

We now get back to the topic of this class: regression

- ▶ For Z -estimation and estimating equations, so far we relied on i.i.d. data
- ▶ In a regression context we focus on modeling $E(Y \mid x)$, and typically focus on the randomness of Y conditional on covariates x
- ▶ Nevertheless, if we can assume that our data are, say,

$$(Y_1, x_1), \dots, (Y_n, x_n) \stackrel{i.i.d.}{\sim} F,$$

then all the results on estimating equations apply to regression estimators

- ▶ If we want to treat the data as independent pairs $\{(Y_i, x_i)\}_{i=1}^n$ where the covariates x_i are fixed, then more careful treatment is required

Ordinary Least Squares, Best Fitting Plane

Let's start with the simplest regression approach

- ▶ Consider the least squares estimator

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \mathbf{x}_i \beta)^2 = \underset{\beta}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta),$$

where

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix}$$

- ▶ Taking derivatives with respect to β and setting them to zero, leads to

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

which we had also obtained as the MLE in the *normal linear model*

Ordinary Least Squares, Best Fitting Plane

- ▶ The system of equations that we solve to find the OLS estimator is

$$\sum_{i=1}^n (Y_i - \mathbf{x}_i \boldsymbol{\beta}) \mathbf{x}_i^\top = \mathbf{0}$$

which fits into the framework of estimating equations!

- ▶ We can take

$$\mathbf{G}(\boldsymbol{\beta}, Y, \mathbf{x}) = (Y - \mathbf{x}\boldsymbol{\beta})\mathbf{x}^\top$$

as the basis for a Z -estimator for $\boldsymbol{\beta}$ based on $\{(Y_i, \mathbf{x}_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} F$

- ▶ If we do not impose any assumptions on the distribution F , this leads to the target of inference being

$$\boldsymbol{\beta}_0 : E_F[(Y - \mathbf{x}\boldsymbol{\beta}_0)\mathbf{x}^\top] = \mathbf{0},$$

from which we find

$$\boldsymbol{\beta}_0 = E_F(\mathbf{x}^\top \mathbf{x})^{-1} E_F(Y \mathbf{x}^\top),$$

with $(Y, \mathbf{x}) \sim F$

Ordinary Least Squares, Best Fitting Plane

- ▶ It can be seen that β_0 is the solution to the population-level least squares problem

$$\beta_0 = \underset{\beta}{\operatorname{argmin}} \quad E_F[(Y - \mathbf{x}\beta_0)^2]$$

- ▶ *Important:* β_0 characterizes the best fitting plane, in a least squares sense
 - ▶ This is well defined across distributions for which $E_F(\mathbf{x}^\top \mathbf{x})^{-1} E_F(Y \mathbf{x}^\top)$ is finite
 - ▶ It *does not* rely on the assumption $E(Y | \mathbf{x}) = \mathbf{x}\beta$
 - ▶ This estimation approach is *nonparametric*, in the sense that it does not impose stringent assumptions on the distribution F

Ordinary Least Squares, Best Fitting Plane

- ▶ In this case, we obtain the \mathbf{A} and \mathbf{B} matrices as

$$\mathbf{A} = \mathbb{E}[\mathbf{G}'(\beta_0, Y, \mathbf{x})] = -\mathbb{E}[\mathbf{x}^\top \mathbf{x}]$$

$$\mathbf{B} = \mathbb{E}[\mathbf{G}(\beta_0, Y, \mathbf{x})\mathbf{G}(\beta_0, Y, \mathbf{x})^\top] = \mathbb{E}[(Y - \mathbf{x}\beta_0)^2 \mathbf{x}^\top \mathbf{x}]$$

- ▶ To compute the sandwich estimator we obtain

$$\hat{\mathbf{A}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i = \mathbf{X}^\top \mathbf{X} / n, \quad \hat{\mathbf{B}} = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i \hat{\beta})^2 \mathbf{x}_i^\top \mathbf{x}_i$$

- ▶ The sandwich variance estimator is

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

which can be used for inference on the parameters β_0 of the best fitting plane

- ▶ However, entries of β_0 are hard to interpret individually

Ordinary Least Squares, Best Fitting Line

Example: consider least squares with a single covariate: $\mathbf{x}\boldsymbol{\beta} = \beta_0 + \beta_1 x$

Problem: Show that the OLS leads to

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\&= \sum_{i=1}^n \sum_{j=1}^n \left(\frac{(x_i - x_j)^2}{\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2} \right) \left(\frac{y_i - y_j}{x_i - x_j} \right) \\&= \sum_{i=1}^n \sum_{j=1}^n w_{ij} \text{slope}_{ij}\end{aligned}$$

- ▶ This expression shows that $\hat{\beta}_1$ is a weighted sum of the slopes of lines connecting all pairs of points
- ▶ Weights are proportional to squared x -distances between pairs of points
- ▶ What does $\hat{\beta}_1$ estimate?

Ordinary Least Squares, Best Fitting Line

- ▶ We can think of the *estimand* $\beta_{1,T}$ as a weighted sum in the (super)population:

$$\beta_{1,T} = \iint \frac{(x - x')^2}{\iint (x - x')^2 dF(x, y) dF(x', y')} \cdot \frac{(y - y')}{(x - x')} dF(x, y) dF(x', y'),$$

where you can think of

- ▶ For X and Y continuous: $\int h(x, y) dF(x, y) = \iint h(x, y) f(x, y) dx dy$
- ▶ For X and Y discrete: $\int h(x, y) dF(x, y) = \sum_{x, y} h(x, y) f(x, y)$

Ordinary Least Squares, Best Fitting Line

- ▶ $\beta_{1,T}$ could be a reasonable target of inference
 - ▶ Exact linearity of relationship between x and $E(Y | x)$ seldom holds
 - ▶ The slope of the best fitting line is also interpretable as a weighted average slope, where points farther apart receive more weight
 - ▶ Could provide a reasonable answer to questions about average trend in the $x - Y$ relationship
- ▶ Again, this is a *nonparametric* approach:
 - ▶ No distribution-family assumption
 - ▶ No mean-model assumption
 - ▶ Estimand does not depend on models
- ▶ Note, however, that the term *nonparametric regression* more commonly refers to flexible ways of modeling the regression function $E(Y | x)$

OLS for the Linear Model

- ▶ In *linear regression* we work under the assumption

$$E(Y_i | \mathbf{x}_i) = \mathbf{x}_i \boldsymbol{\beta},$$

or equivalently, $Y_i = \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i$, $E(\epsilon_i | \mathbf{x}_i) = 0$

- ▶ If, in addition, we assume homoscedasticity, that is,

$$\text{var}(Y_i | \mathbf{x}_i) = \text{var}(\epsilon_i | \mathbf{x}_i) = \sigma^2$$

then we obtain a simplification of the previous result:

$$\mathbf{B} = E[(Y - \mathbf{x}\boldsymbol{\beta}_0)^2 \mathbf{x}^\top \mathbf{x}] = E[\epsilon^2 \mathbf{x}^\top \mathbf{x}] = E[E(\epsilon^2 | \mathbf{x}) \mathbf{x}^\top \mathbf{x}] = \sigma^2 E(\mathbf{x}^\top \mathbf{x})$$

- ▶ This allows us to simplify

$$\mathbf{A}^{-1} \mathbf{B} (\mathbf{A}^{-1})^\top / n = \sigma^2 E(\mathbf{x}^\top \mathbf{x})^{-1} / n,$$

and so we obtain the estimator

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

where the estimator $\hat{\sigma}^2$ can be taken as

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

OLS for the Linear Model

- ▶ Note that we had obtained the OLS and the variance estimator:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}, \quad \widehat{\text{var}}(\hat{\beta}) = \hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

under the *normal linear model*, where we in addition assumed normality of the response (of the errors)

- ▶ These estimators now appear using an *asymptotic* justification under a *semiparametric regression model*, where we assume

- ▶ $\{(Y_i, \mathbf{x}_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} F$

- ▶ $E(Y_i \mid \mathbf{x}_i) = \mathbf{x}_i \beta$

- ▶ $\text{var}(Y_i \mid \mathbf{x}_i) = \sigma^2$

- ▶ However, these estimators do not depend on the assumption of i.i.d. data (not even needed to think of the covariates as random), and actually they enjoy certain optimality properties in finite samples

OLS for the Linear Model

Unbiasedness:

- ▶ If $E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}\beta$, then

$$\begin{aligned}E(\hat{\beta} \mid \mathbf{X}) &= E[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \mid \mathbf{X}] \\&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\mathbf{Y} \mid \mathbf{X}) \\&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta\end{aligned}$$

- ▶ Notice that this only relies on correct specification of the mean model as $E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}\beta$

OLS for the Linear Model

Variance:

- ▶ If $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \sigma^2 \mathbf{I}_n$, then

$$\begin{aligned}\text{var}(\hat{\beta} \mid \mathbf{X}) &= \text{var}\{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \mid \mathbf{X}\} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{var}(\mathbf{Y} \mid \mathbf{X}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

- ▶ Notice that this only relies on the homoskedasticity assumption $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \sigma^2 \mathbf{I}_n$

OLS for the Linear Model: Gauss-Markov Theorem

Theorem (Gauss-Markov): Among all *linear, unbiased* estimators of β in the model that assumes

$$E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}\beta, \quad \text{var}(\mathbf{Y} \mid \mathbf{X}) = \sigma^2 \mathbf{I}_n, \quad \sigma^2 < \infty,$$

the OLS estimator $\hat{\beta}$ has minimum variance and is unique. Likewise, the estimator $\mathbf{z}^T \hat{\beta}$ of $\mathbf{z}^T \beta$ has minimum variance and is unique.

Note: no normality assumed, no i.i.d. data assumed, the result holds for finite sample sizes n , and the covariates \mathbf{X} are seen as fixed

OLS for the Linear Model: Gauss-Markov Theorem

- ▶ Note that to talk about “minimization” of $\text{var}(\hat{\beta})$, for $\dim(\beta) > 1$, we need a way of ordering positive semi-definite matrices
- ▶ We use the so-called *Loewner partial order*, which says that $\mathbf{A} \succeq \mathbf{B}$ if $\mathbf{A} - \mathbf{B}$ is positive semi-definite
- ▶ For covariance of estimators, $\text{var}(\tilde{\beta}) \succeq \text{var}(\hat{\beta})$ is equivalent to $\mathbf{z}^T \hat{\beta}$ being optimal for estimating univariate contrasts $\mathbf{z}^T \beta$:

$$\begin{aligned}\text{var}(\mathbf{z}^T \tilde{\beta}) &= \mathbf{z}^T \text{var}(\tilde{\beta}) \mathbf{z} \\ \text{var}(\mathbf{z}^T \hat{\beta}) &= \mathbf{z}^T \text{var}(\hat{\beta}) \mathbf{z},\end{aligned}$$

and so by definition of positive semi-definite matrices:

$$\text{var}(\tilde{\beta}) \succeq \text{var}(\hat{\beta})$$

if and only if, for all vectors \mathbf{z} ,

$$\text{var}(\mathbf{z}^T \hat{\beta}) - \text{var}(\mathbf{z}^T \tilde{\beta}) = \mathbf{z}^T [\text{var}(\tilde{\beta}) - \text{var}(\hat{\beta})] \mathbf{z} \geq 0$$

OLS for the Linear Model: Gauss-Markov Theorem

Proof of Gauss-Markov: Linearity and unbiasedness of an estimator $\tilde{\beta}$ means:

- ▶ $\tilde{\beta} = \mathbf{C}\mathbf{Y}$, where \mathbf{C} can depend on \mathbf{X} but not on \mathbf{Y}
- ▶ $E(\tilde{\beta} | \mathbf{X}) = \beta$

Let $\Delta = \mathbf{C} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. Then

$$\begin{aligned} E(\tilde{\beta} | \mathbf{X}) &= E(\mathbf{C}\mathbf{Y} | \mathbf{X}) = \mathbf{C}\mathbf{X}\beta = ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \Delta)\mathbf{X}\beta \\ &= (\mathbf{I}_{k+1} + \Delta\mathbf{X})\beta \end{aligned}$$

$\Rightarrow \Delta\mathbf{X} = \mathbf{0}$, since $E(\tilde{\beta} | \mathbf{X}) = \beta$.

OLS for the Linear Model: Gauss-Markov Theorem

Proof of Gauss-Markov, cont'd:

$$\begin{aligned}\text{var}(\tilde{\beta} | \mathbf{X}) &= \text{var}(\mathbf{C}\mathbf{Y} | \mathbf{X}) = \mathbf{C}\text{var}(\mathbf{Y} | \mathbf{X})\mathbf{C}^T = \mathbf{C}\mathbf{I}_n\mathbf{C}^T\sigma^2 = \mathbf{C}\mathbf{C}^T\sigma^2 \\&= ((\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T + \Delta)((\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T + \Delta)^T\sigma^2 \\&= ((\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\Delta^T \\&\quad + \Delta\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} + \Delta\Delta^T)\sigma^2 \\&= (\mathbf{X}^T\mathbf{X})^{-1}\sigma^2 + \Delta\Delta^T\sigma^2 \\&= \text{var}(\hat{\beta} | \mathbf{X}) + \Delta\Delta^T\sigma^2\end{aligned}$$

Since $\Delta\Delta^T$ is positive semidefinite, this is minimized when $\Delta = \mathbf{0}$ because

$$\text{var}(\mathbf{z}^T\tilde{\beta}) = \mathbf{z}^T[(\mathbf{X}^T\mathbf{X})^{-1} + \Delta\Delta^T]\sigma^2\mathbf{z}.$$

is minimized for $\Delta = \mathbf{0}$.

This also shows that if $\text{var}(\tilde{\beta} | \mathbf{X}) = \text{var}(\hat{\beta} | \mathbf{X})$ then

$$\Delta = \mathbf{C} - (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T = \mathbf{0} \implies \tilde{\beta} = \mathbf{C}\mathbf{Y} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \hat{\beta}.$$

OLS for the Linear Model: Gauss-Markov Theorem

- ▶ Thus, the OLS estimator $\hat{\beta}$ is the *Best* (i.e. lowest variance) *Linear Unbiased Estimator* (BLUE) of the coefficients in the linear regression model
- ▶ Normality of Y (or of errors) not required
- ▶ Unbiasedness is important: there exist biased estimators with lower $\text{MSE} = \text{bias}^2 + \text{variance}$
 - ▶ ridge regression
 - ▶ lasso
 - ▶ Bayesian estimators, etc.

OLS for the Linear Model: Gauss-Markov Theorem

The Gauss-Markov theorem can be interpreted as demonstrating the properties of OLS estimators in a *semiparametric* model:

- ▶ Minimal distributional assumptions:
 - ▶ $E(\mathbf{Y} \mid \mathbf{X}) < \infty$;
 - ▶ $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \sigma^2 \mathbf{I}_n$, where $\sigma^2 < \infty$
- ▶ Mean model: $E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}\beta$
- ▶ $\hat{\beta}$ estimates β in this mean model
- ▶ The Gauss-Markov Theorem says that among unbiased estimators of β in this semiparametric model, the OLS estimator has minimum variance

OLS for the Linear Model: Asymptotic Distribution

What is the asymptotic distribution of the OLS estimator when we see the covariates as fixed or as random but not i.i.d.?

We now consider using OLS for inference in a model with the following characteristics:

- ▶ Covariates $\mathbf{x}_1, \dots, \mathbf{x}_n$ fixed or not i.i.d.
- ▶ Mean model: $E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}\beta$
- ▶ Possibly heteroskedastic, but uncorrelated responses (errors):
 $\text{var}(\mathbf{Y} \mid \mathbf{X}) = \text{diag}\{\sigma_i^2\}$

Equivalently:

$$Y_i = \mathbf{x}_i\beta + \epsilon_i, \quad E(\epsilon_i \mid \mathbf{x}_i) = 0, \quad \text{var}(\epsilon_i \mid \mathbf{x}_i) = \sigma_i^2, \quad \text{cov}(\epsilon_i, \epsilon_j) = 0$$

OLS for the Linear Model: Asymptotic Distribution

Under this model, different authors provide conditions that guarantee asymptotic normality of the OLS $\hat{\beta}_n$:

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

where

$$\mathbf{A}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\mathbf{x}_i^\top \mathbf{x}_i), \quad \mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\sigma_i^2 \mathbf{x}_i^\top \mathbf{x}_i]$$

where $\sigma_i^2 = \mathbb{E}(\epsilon_i^2 \mid \mathbf{x}_i) = \mathbb{E}[(Y_i - \mathbf{x}_i\beta)^2 \mid \mathbf{x}_i]$

If the covariates are taken as fixed, then

$$\mathbf{A}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i = \mathbf{X}^\top \mathbf{X} / n, \quad \mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i^\top \mathbf{x}_i = \mathbf{X}^\top \text{diag}\{\sigma_i^2\} \mathbf{X} / n$$

OLS for the Linear Model: Asymptotic Distribution

This is used as the basis for estimating $\text{var}(\hat{\beta}_n)$ as:

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

which might be found under the names of *sandwich*, *robust*, *heteroscedasticity-consistent*, *Huber-Eicker-White*, or *HCO* estimator

- ▶ This estimator was developed separately by three authors: Friedhelm Eicker, Peter J. Huber, and Halbert White
- ▶ This sandwich estimator is consistent for $\text{var}(\hat{\beta})$ under the linear model with heteroscedastic errors
- ▶ Notice that this estimator has exactly the same form as the sandwich estimator under the “best fitting plane” approach: the two approaches are numerically the same!
 - ▶ *Question:* If both approaches are numerically the same, what do we gain by assuming $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i \beta$?

OLS for the Linear Model: Asymptotic Distribution

- ▶ The sandwich estimator $\widehat{\text{var}}(\hat{\beta}_n)$ is more variable than the estimator under homoscedasticity, $\hat{\sigma}^2(\mathbf{X}^\top \mathbf{X})^{-1}$, in small samples
- ▶ Some small-sample fixes have been proposed:
 - ▶ HC1: replace $e_i^2 := (Y_i - \mathbf{x}_i \hat{\beta})^2$ with $\frac{n}{n-k-1} e_i^2$
 - ▶ HC2: replace e_i^2 with $\frac{e_i^2}{1-h_i}$, where h_i is the i th diagonal element of the “hat” matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$
 - ▶ More options in the `sandwich` R package

Each leads to consistent estimators of $\text{var}(\hat{\beta})$, see MacKinnon and White (1985).

OLS for the Linear Model: Asymptotic Distribution

- ▶ Note that

$$\mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\sigma_i^2 \mathbf{x}_i^\top \mathbf{x}_i)$$

depends on n different variances $\sigma_i^2 = \mathbb{E}(\epsilon_i^2 \mid \mathbf{x}_i) = \mathbb{E}[(Y_i - \mathbf{x}_i \boldsymbol{\beta})^2 \mid \mathbf{x}_i]$

- ▶ Estimating $\hat{\mathbf{B}}_n = \mathbf{X}^\top \text{diag}\{(Y_i - \mathbf{x}_i \hat{\boldsymbol{\beta}})^2\} \mathbf{X} / n$ is as if we took $\hat{\sigma}_i^2 = (Y_i - \mathbf{x}_i \hat{\boldsymbol{\beta}})^2$
- ▶ However, estimating n different variances σ_i^2 with n data points is hopeless
- ▶ How do theoreticians handled this situation?
- ▶ *Idea:* focus on estimating \mathbf{B}_n , not each individual σ_i^2 , and set yourself up for success, i.e., put enough conditions under which you can do a good job at estimating \mathbf{B}_n (and \mathbf{A}_n)

OLS for the Linear Model: Asymptotic Distribution

Some conditions for success: control \mathbf{A}_n and \mathbf{B}_n as $n \rightarrow \infty$

- ▶ Boos and Stefanski (2013, p. 316) for non-random covariates, assume these limits exist:

$$\lim_{n \rightarrow \infty} \mathbf{B}_n = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i^\top \mathbf{x}_i = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(Y_i - \mathbf{x}_i \beta)^2] \mathbf{x}_i^\top \mathbf{x}_i$$

$$\lim_{n \rightarrow \infty} \mathbf{A}_n = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{X} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i$$

- ▶ White (1980) imposes conditions to control the behavior of

$$\mathbf{A}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\mathbf{x}_i^\top \mathbf{x}_i), \quad \mathbf{B}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\sigma_i^2 \mathbf{x}_i^\top \mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\epsilon_i^2 \mathbf{x}_i^\top \mathbf{x}_i)$$

(interestingly, White doesn't require that these converge to a limit)

OLS for the Linear Model: Asymptotic Distribution

White (1980) imposed the following assumptions:

- ▶ $E(\mathbf{x}_i \epsilon_i) = \mathbf{0}$ (particular cases: $E(\epsilon_i | \mathbf{x}_i) = 0$ and $\epsilon_i \perp\!\!\!\perp \mathbf{x}_i$)
- ▶ There exist constants $\delta_1, \delta_2 > 0$ such that for all i, j, k , $E(|\epsilon_i^2|^{1+\delta_1}) < \delta_2$, $E(|X_{ij}X_{ik}|^{1+\delta_1}) < \delta_2$ (\approx uniformly bounded error variances and covariance of covariates, stronger)
- ▶ \mathbf{A}_n is non-singular for all $n > n_0$, such that $\det \mathbf{A}_n > \delta_1 > 0$ (eventually the average covariance matrix of covariates is non-singular, and elements of its inverse are uniformly bounded)
- ▶ Similar conditions for $E(|\epsilon_i^2 X_{ij}X_{ik}|^{1+\delta_1}) < \delta_2$ and for \mathbf{B}_n (elements of \mathbf{B}_n and of its inverse are uniformly bounded)
- ▶ There exist $\delta_1, \delta_2 > 0$ such that for all i, j, k, l , $E(|X_{ij}^2 X_{ik}X_{il}|^{1+\delta_1}) < \delta_2$

OLS for the Linear Model: Asymptotic Distribution

- *Lemma 1 of White (1980, Econometrica Vol. 48 No. 4):*

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

- *Theorem 1 of White (1980, Econometrica Vol. 48 No. 4):*
Let $\hat{\mathbf{B}}_n = \mathbf{X}^\top \text{diag}\{(Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2\} \mathbf{X}/n$. Under the assumptions above

- i) $|\hat{\mathbf{B}}_n - \mathbf{B}_n| \xrightarrow{a.s.} 0$
- ii) $|(\mathbf{X}^\top \mathbf{X}/n)^{-1} \hat{\mathbf{B}}_n (\mathbf{X}^\top \mathbf{X}/n)^{-1} - \mathbf{A}_n^{-1} \mathbf{B}_n \mathbf{A}_n^{-1}| \xrightarrow{a.s.} 0$
- iii) Under $H_0 : \mathbf{C}\boldsymbol{\beta} = \mathbf{b}$, for \mathbf{C} of rank q

$$n(\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{b})^\top [\mathbf{C}(\mathbf{X}^\top \mathbf{X}/n)^{-1} \hat{\mathbf{B}}_n (\mathbf{X}^\top \mathbf{X}/n)^{-1} \mathbf{C}^\top]^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{b}) \xrightarrow{d} \chi_q^2$$

OLS for the Linear Model: Asymptotic Distribution

Key points in the proof of part i) $|\hat{\mathbf{B}}_n - \mathbf{B}_n| \xrightarrow{a.s.} 0$

- ▶ Lemma 2.3 of White (1980, Econometrica Vol. 48 No. 3): Let Z_1, \dots, Z_n be independent random variables. Under some conditions

$$\sup_{\theta} \left| \frac{1}{n} \sum_{i=1}^n q_i(Z_i, \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[q_i(Z_i, \theta)] \right| \xrightarrow{a.s.} 0$$

- ▶ Lemma 2.6 of the same paper: If in addition $\hat{\theta} \xrightarrow{a.s.} \theta$ then

$$\left| \frac{1}{n} \sum_{i=1}^n q_i(Z_i, \hat{\theta}) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[q_i(Z_i, \theta)] \right| \xrightarrow{a.s.} 0$$

- ▶ Since $\hat{\beta} \xrightarrow{a.s.} \beta$, these are combined to show

$$\left| \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i \hat{\beta})^2 \mathbf{x}_i^T \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\epsilon_i^2 \mathbf{x}_i^T \mathbf{x}_i] \right| \xrightarrow{a.s.} 0$$

where $\epsilon_i = Y_i - \mathbf{x}_i \beta$. That is, $|\hat{\mathbf{B}}_n - \mathbf{B}_n| \xrightarrow{a.s.} 0$.

OLS for the Linear Model: Asymptotic Distribution

Connection of part **iii)** of White's theorem with other tests:

► Under errors with equal variance σ^2

► $(\mathbf{X}^\top \mathbf{X})^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})/\sigma \xrightarrow{d} N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$

► Then under $H_0 : \mathbf{C}\boldsymbol{\beta} = \mathbf{b}$

$$(RSS_{H_0} - RSS)/\sigma^2 = (\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{b})^\top [\mathbf{C}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{C}^\top]^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{b})/\sigma^2 \xrightarrow{d} \chi_q^2$$

► Also, $\hat{\sigma}^2 = RSS/(n - k - 1) \xrightarrow{a.s.} \sigma^2$

► Then under H_0 we have $(RSS_{H_0} - RSS)/\hat{\sigma}^2 \xrightarrow{d} \chi_q^2$

► Under i.i.d. normal errors we had found $(RSS_{H_0} - RSS)/q\hat{\sigma}^2 \sim F_{q, n-k-1}$, so what's the connection?

► If $Z \sim F(m, n)$ then as $n \rightarrow \infty$, $mZ \xrightarrow{d} \chi_m^2$

► Under H_0 we have $(RSS_{H_0} - RSS)/\hat{\sigma}^2 \sim qF_{q, n-k-1} \xrightarrow{d} \chi_q^2$

OLS for the Linear Model: Asymptotic Distribution

Key important pieces of White's results for inference:

- ▶ Under $H_0 : \mathbf{C}\boldsymbol{\beta} = \mathbf{b}$, for \mathbf{C} of rank q

$$(\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{b})^T [\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \hat{\mathbf{B}}_n (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T]^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{b}) / n \xrightarrow{d} \chi_q^2$$

- ▶ Confidence regions for $\mathbf{C}\boldsymbol{\beta}$

$$\{\mathbf{C}\boldsymbol{\beta} : \mathbf{C}^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T [\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \hat{\mathbf{B}}_n (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T]^{-1} \mathbf{C} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) / n \leq \chi_q^2(1-\alpha)\}$$

OLS for the Linear Model: Asymptotic Distribution

Particular cases:

- Confidence interval for a single β_j , equivalent to

$$(\hat{\beta}_j - \text{se}(\hat{\beta}_j)z_{1-\alpha/2}, \hat{\beta}_j + \text{se}(\hat{\beta}_j)z_{1-\alpha/2})$$

with $\text{se}(\hat{\beta}_j)^2$ the j th element of the diagonal of $n(\mathbf{X}^\top \mathbf{X})^{-1} \hat{\mathbf{B}}_n(\mathbf{X}^\top \mathbf{X})^{-1}$

- Confidence region for $\beta_S = (\beta_j : j \in S, S \subseteq 0 : k)$ can be formed by

$$\{\beta_S : (\hat{\beta}_S - \beta_S)^\top \widehat{\text{var}}(\hat{\beta}_S)^{-1} (\hat{\beta}_S - \beta_S) \leq \chi_q^2(1 - \alpha)\}$$

where

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

with \mathbf{C} a matrix that indicates how to extract the entries S of a vector of length $k + 1$

Summary: Properties of OLS Under Different Assumptions

We covered properties of OLS for different purposes:

OLS for the best fitting plane as a summary of the relationship of response Y_i and covariates \mathbf{x}_i

- ▶ Estimand β characterizes best fitting plane, but individual elements in β hard to interpret
- ▶ OLS $\hat{\beta}_n$ satisfies

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

- ▶ Variance of $\hat{\beta}_n$ can be estimated as

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

Summary: Properties of OLS Under Different Assumptions

OLS for the parameters in the linear mean model (regression function), assumed to be $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i\beta$, allowing for heteroskedastic errors:

$$Y_i = \mathbf{x}_i\beta + \epsilon_i, \quad E(\epsilon_i | \mathbf{x}_i) = 0, \quad \text{var}(\epsilon_i | \mathbf{x}_i) = \sigma_i^2, \quad \text{cov}(\epsilon_i, \epsilon_j) = 0$$

- ▶ Estimand β characterizes the regression function, and individual elements in β can be interpretable as usual in linear regression
- ▶ OLS $\hat{\beta}_n$ satisfies

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

- ▶ Variance of $\hat{\beta}_n$ can be estimated as

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

Summary: Properties of OLS Under Different Assumptions

OLS for the parameters in the linear mean model (regression function), assumed to be $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i\beta$, with homoskedastic errors:

$$Y_i = \mathbf{x}_i\beta + \epsilon_i, \quad E(\epsilon_i | \mathbf{x}_i) = 0, \quad \text{var}(\epsilon_i | \mathbf{x}_i) = \sigma^2, \quad \text{cov}(\epsilon_i, \epsilon_j) = 0$$

- ▶ Individual elements in β can be interpretable as usual in linear regression since we assume $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i\beta$
- ▶ OLS $\hat{\beta}_n$ satisfies $(\mathbf{X}^\top \mathbf{X})^{1/2}(\hat{\beta}_n - \beta)/\sigma \xrightarrow{d} N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$
- ▶ Variance of $\hat{\beta}_n$ can be estimated as $\widehat{\text{var}}(\hat{\beta}) = \hat{\sigma}^2(\mathbf{X}^\top \mathbf{X})^{-1}$ where the estimator $\hat{\sigma}^2$ can be taken as

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

- ▶ OLS $\hat{\beta}_n$ is BLUE (Gauss-Markov)

Summary: Properties of OLS Under Different Assumptions

Back to the first part of the course, OLS for the parameters in the normal linear model:

$$Y_i = \mathbf{x}_i\boldsymbol{\beta} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

- ▶ Individual elements in $\boldsymbol{\beta}$ can be interpretable as usual in linear regression since we assume $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i\boldsymbol{\beta}$
- ▶ OLS $\hat{\boldsymbol{\beta}}_n$ derived as MLE
- ▶ OLS $\hat{\boldsymbol{\beta}}_n$ satisfies $(\mathbf{X}^\top \mathbf{X})^{1/2}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta})/\sigma \sim N_{k+1}(\mathbf{0}, \mathbf{I}_{k+1})$ exactly, in finite samples
- ▶ Variance of $\hat{\boldsymbol{\beta}}_n$ can be estimated as $\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2(\mathbf{X}^\top \mathbf{X})^{-1}$ where the estimator $\hat{\sigma}^2$ can be taken as

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

- ▶ Many other exact, finite sample distributional results useful for tests and confidence intervals