

# BIOS 7345: Advanced Regression for Independent Data

**Andrew J. Spieker, Ph.D.**

Associate Professor of Biostatistics  
Vanderbilt University

Set 12: Sandwich and bootstrap methods for GLMs

Version: 06/11/2025

# TABLE OF CONTENTS

- 1 Variance based on likelihood theory
- 2 Variance based on theory of estimating equations
- 3 The nonparametric bootstrap

## Recall:

- A GLM involves specifying a parametric form for  $Y$  (given  $\mathbf{x}$ ) that can be factored into a nice exponential form.
- The score equations take the form:

$$\mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu})/\phi = \mathbf{0}.$$

- From the prior set of notes, we saw that for a GLM,

$$\mathcal{I}(\boldsymbol{\beta}, \phi) = \begin{bmatrix} \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \phi & \mathbf{0}^T \\ \mathbf{0} & \dots \end{bmatrix}$$

- The block-diagonal structure of the information tells us we need not propagate uncertainty in estimation of  $\phi$  in estimating variance of  $\hat{\boldsymbol{\beta}}$ .

## Asymptotic distribution:

- Likelihood theory (suitable regularity conditions):

$$\widehat{\boldsymbol{\beta}} \underset{\sim}{\sim} \mathcal{N}(\boldsymbol{\beta}, \phi(\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D})^{-1}).$$

- To estimate  $\text{Cov}[\widehat{\boldsymbol{\beta}}]$ , we could be in one of two cases:
  - 1  $\phi = 1$ , in which case  $\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = (\mathbb{A}_N(\widehat{\boldsymbol{\beta}}))^{-1}$ .
  - 2 Otherwise, we require a consistent estimate of  $\phi$ . This one will do:

$$\widehat{\phi} = \frac{1}{N - K} \sum_{i=1}^N \frac{(y_i - \widehat{\mu}_i)^2}{V(\widehat{\mu}_i)},$$

and we have that  $\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = \widehat{\phi}(\mathbb{A}_N(\widehat{\boldsymbol{\beta}}))^{-1}$

# TABLE OF CONTENTS

- 1 Variance based on likelihood theory
- 2 Variance based on theory of estimating equations
- 3 The nonparametric bootstrap

## Ideas:

- We don't need to solve for  $\phi$  to solve the score equations for  $\beta$ . We simply solve the estimating equations  $\mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}$ .
- These are referred to as unbiased estimating equations.
  - ▶ When mean model is correct,  $E[\mathbb{G}_N(\beta; \mathbf{X}, \mathbf{y})] = \mathbf{0}$ , where expectation can be either over  $\mathbf{y}|\mathbf{X}$  or  $(\mathbf{X}, \mathbf{y})$ —recall “fixed” vs. “random.”
- What if the likelihood is not correctly specified?
  - ▶ For instance, what if the mean model is not correct?
  - ▶ For instance, what if the mean-variance relationship is not correct?
  - ▶ For instance, what if the third (or higher) moment is not correct?
- Can we derive an expression for  $\text{Cov}[\hat{\beta}]$  based on the theory of estimating equations rather than likelihood theory?
  - ▶ Must assume  $\mathbf{X}$  is random (sampling over  $(\mathbf{X}, \mathbf{y})$ ).

## Notation:

- $\hat{\boldsymbol{\beta}}_N$ : solution to estimating equations.
- $\boldsymbol{\beta}_0$ : the true, unknown value to be estimated.
  - ▶ If the mean model is not correctly specified, then  $\boldsymbol{\beta}_0$  can be understood as “the quantity for which  $\hat{\boldsymbol{\beta}}_N$  is consistent.”
  - ▶ Of course, consistency of  $\hat{\boldsymbol{\beta}}_N$  for *any* quantity at all is not a given—I am assuming all the regularity conditions surrounding  $\mathbf{X}$  and  $\mathbf{y}$  that are necessary for  $\hat{\boldsymbol{\beta}}_N$  to converge in probability hold.
- $\mathbb{G}_N(\boldsymbol{\beta}; \mathbf{X}, \mathbf{y}) = \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i)$ .
- $\mathbf{A}(\boldsymbol{\beta}) = E\left[-\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{G}(\boldsymbol{\theta}; \mathbf{x}, Y) \Big|_{\boldsymbol{\theta}=\boldsymbol{\beta}}\right]$
- $\mathbf{B}(\boldsymbol{\beta}) = E[\mathbf{G}(\boldsymbol{\beta}; \mathbf{x}, Y)\mathbf{G}(\boldsymbol{\beta}; \mathbf{x}, Y)^\top]$

## Taylor expansion:

- Because  $\hat{\boldsymbol{\beta}}_N$  solves the estimating equations, it follows that:

$$\mathbf{0} = \frac{1}{N} \mathbb{G}_N(\hat{\boldsymbol{\beta}}_N; \mathbf{X}, \mathbf{y})$$

- If  $\mathbf{G}$  is analytic (has a Taylor series), we can expand about  $\boldsymbol{\beta}_0$ :

$$\begin{aligned} \mathbf{0} &\approx \frac{1}{N} \mathbb{G}_N(\boldsymbol{\beta}_0; \mathbf{X}, \mathbf{y}) + \left. \frac{\partial}{\partial \boldsymbol{\beta}} \left[ \frac{1}{N} \mathbb{G}_N(\boldsymbol{\beta}; \mathbf{X}, \mathbf{y}) \right] \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} (\hat{\boldsymbol{\beta}}_N - \boldsymbol{\beta}_0) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}_i, Y_i) + \left[ \frac{1}{N} \sum_{i=1}^N \left. \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i) \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right] (\hat{\boldsymbol{\beta}}_N - \boldsymbol{\beta}_0) \end{aligned}$$

## Rearrangement:

- Assume that  $\sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i) \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0}$  is invertible.
- Rearranging the equation on the prior slide (and leaving the details surrounding the regularity conditions on the remainder term of the Taylor expansion to a more theory-oriented course):

$$(\hat{\boldsymbol{\beta}}_N - \boldsymbol{\beta}_0) \approx \left[ -\frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i) \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right]^{-1} \left[ \frac{1}{N} \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}_i, Y_i) \right]$$

## Invoking asymptotics:

- Multiplying both sides by  $\sqrt{N}$ , we then have:

$$\sqrt{N}(\hat{\boldsymbol{\beta}}_N - \boldsymbol{\beta}_0) \approx \left[ -\frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i) \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right]^{-1} \left[ \frac{\sqrt{N}}{N} \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}_i, Y_i) \right]$$

- By the weak law of large numbers,

$$\left[ -\frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i) \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right] \xrightarrow{P} \mathbf{A}(\boldsymbol{\beta}_0)$$

- By the central limit theorem\*,

$$\frac{\sqrt{N}}{N} \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}_i, Y_i) = \left[ \sqrt{N} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}_i, Y_i) - \mathbf{0} \right) \right] \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{B}(\boldsymbol{\beta}_0)).$$

## Point of nuance:

- There is a nuanced point here that's easy to miss.
- Previous slide appears to assume  $E[\mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}, \mathbf{y})] = \mathbf{0}$  in this step:

$$\left[ \sqrt{N} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}_i, Y_i) - \mathbf{0} \right) \right] \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{B}(\boldsymbol{\beta}_0)).$$

- We're in the middle of deriving the asymptotic distribution for  $\hat{\boldsymbol{\beta}}_N$  in the setting that the mean model may *not* be correct. Is this fair?
- When we refer to  $\boldsymbol{\beta}_0$  as the “true value of the parameter,” it is more accurate to think of it as the value for which the implicit solution to the estimating equations is consistent.
- Since  $\mathbb{G}_N(\hat{\boldsymbol{\beta}}_N) = \mathbf{0}$ , we have  $E[\mathbb{G}_N(\boldsymbol{\beta}_0; \mathbf{X}, \mathbf{y})] = E[\mathbf{G}(\boldsymbol{\beta}_0; \mathbf{x}, Y)] = \mathbf{0}$  even if the mean model is not correct (assume  $\hat{\boldsymbol{\beta}}_N \xrightarrow{P} \boldsymbol{\beta}$ ).

## More asymptotics:

- Returning to the derivation, it follows from Slutsky's theorem that

$$\hat{\boldsymbol{\beta}}_N \sim \mathcal{N}\left(\boldsymbol{\beta}_0, \frac{1}{N}[\mathbf{A}(\boldsymbol{\beta}_0)]^{-1}\mathbf{B}(\boldsymbol{\beta}_0)[\mathbf{A}(\boldsymbol{\beta}_0)]^{-1}\right)$$

- To estimate  $\text{Cov}[\hat{\boldsymbol{\beta}}]$ , we can plug in estimators of  $\mathbf{A}(\boldsymbol{\beta}_0)$  and  $\mathbf{B}(\boldsymbol{\beta}_0)$  (such estimators are known as *sandwich* estimators).

## Plug-in estimators:

- How do we estimate  $\mathbf{A}(\boldsymbol{\beta}_0)$ ? Recall that:

$$\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta}) = -\frac{\partial}{\partial \boldsymbol{\beta}} \mathbb{G}_N(\boldsymbol{\beta}; \mathbf{X}, \mathbf{y}).$$

- By the weak law of large numbers and continuous mapping theorem,

$$\frac{1}{N} \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) = -\frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{G}(\boldsymbol{\beta}, \mathbf{x}_i, Y_i) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \xrightarrow{P} \mathbf{A}(\boldsymbol{\beta}_0).$$

- This statement is valid even if the mean model is not correct.

## Plug-in estimators:

- How do we estimate  $\mathbf{B}(\boldsymbol{\beta}_0)$ ? Let

$$\begin{aligned}\mathbb{B}_N^{\text{obs}}(\boldsymbol{\beta}) &= \sum_{i=1}^N \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i) \mathbf{G}(\boldsymbol{\beta}; \mathbf{x}_i, Y_i)^\top \\ &= \mathbf{D}^\top(\boldsymbol{\beta}) \mathbf{V}^{-1}(\boldsymbol{\beta}) \text{diag}(Y_i - \mu_i(\boldsymbol{\beta}))^2 \mathbf{V}^{-1}(\boldsymbol{\beta}) \mathbf{D}(\boldsymbol{\beta}).\end{aligned}$$

- By the weak law of large numbers and continuous mapping theorem,

$$\frac{1}{N} \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) = \frac{1}{N} \sum_{i=1}^N \mathbf{G}(\hat{\boldsymbol{\beta}}; \mathbf{x}_i, Y_i) \mathbf{G}(\hat{\boldsymbol{\beta}}; \mathbf{x}_i, Y_i)^\top \xrightarrow{p} \mathbf{B}(\boldsymbol{\beta}_0).$$

- This statement is valid even if the mean-variance relationship is not correctly specified.

## Plug-in estimators:

- We can therefore estimate  $\text{Cov}[\hat{\boldsymbol{\beta}}]$  as follows:

$$\begin{aligned}\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] &= \frac{1}{N} \left( \frac{1}{N} \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right)^{-1} \left( \frac{1}{N} \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right) \left( \frac{1}{N} \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right)^{-1} \\ &= \left( \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right)^{-1} \left( \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right) \left( \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right)^{-1},\end{aligned}$$

where  $\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$  and  $\mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$  are presented on the prior slides.

- This is *one* version of the sandwich estimator, and it is asymptotically valid even if neither aspect of the GLM (meaning the mean model and the mean-variance relationship) is correctly specified.

## Correct mean model:

- If the mean model is correct, it is straightforward to show:

$$\mathbf{A}(\boldsymbol{\beta}_0) = E_x[\mathbf{x}w(\boldsymbol{\beta}_0)\mathbf{x}^T].$$

In fact, this was determined in our presentation of the Gauss-Newton algorithm.

- If we believe the mean model to hold, it therefore seems sensible to estimate  $\mathbf{A}(\boldsymbol{\beta}_0)$  based on  $\mathbb{A}_N(\hat{\boldsymbol{\beta}})$  rather than  $\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$ . By the weak law of large numbers (if the mean model is correct):

$$\frac{1}{N}\mathbb{A}_N(\hat{\boldsymbol{\beta}}) = \frac{1}{N}\mathbf{X}^T\mathbf{W}(\hat{\boldsymbol{\beta}})\mathbf{X} \xrightarrow{P} \mathbf{A}(\boldsymbol{\beta}_0).$$

- Keep in mind: sometimes  $\mathbb{A}_N(\boldsymbol{\beta})$  and  $\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta})$  are the same.
  - ▶ When, in particular?

**Plug-in estimators:** When the mean model is correct

- We can estimate  $\text{Cov}[\hat{\boldsymbol{\beta}}]$  as follows:

$$\begin{aligned}\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] &= \frac{1}{N} \left( \frac{1}{N} \mathbb{A}_N(\hat{\boldsymbol{\beta}}) \right)^{-1} \left( \frac{1}{N} \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right) \left( \frac{1}{N} \mathbb{A}_N(\hat{\boldsymbol{\beta}}) \right)^{-1} \\ &= \left( \mathbb{A}_N(\hat{\boldsymbol{\beta}}) \right)^{-1} \left( \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) \right) \left( \mathbb{A}_N(\hat{\boldsymbol{\beta}}) \right)^{-1}\end{aligned}$$

- This is another version of the sandwich estimator; it is asymptotically valid if the mean-variance relationship of the GLM is misspecified.
- Under the canonical link, this will be equivalent to the previous sandwich and validity will not depend upon mean model being correct.
- Under a non-canonical link, validity depends upon the mean model being correctly specified.

## Correct mean model and mean-variance relationship:

- If both the mean model and mean-variance relationship are correct, it is straightforward to show:

$$\mathbf{B}(\boldsymbol{\beta}_0) = \phi \mathbf{A}(\boldsymbol{\beta}_0)$$

- If we believe the mean model to be correct and the mean-variance relationship to hold, it therefore seems sensible to estimate  $\mathbf{B}(\boldsymbol{\beta}_0)$  based on  $\mathbb{B}_N(\hat{\boldsymbol{\beta}})$  rather than  $\mathbb{B}_N^{\text{obs}}(\boldsymbol{\beta})$ :

$$\frac{1}{N} \mathbb{B}_N(\hat{\boldsymbol{\beta}}) := \frac{1}{N} \hat{\phi} \mathbf{A}_n(\hat{\boldsymbol{\beta}}) \xrightarrow{p} \mathbf{B}(\boldsymbol{\beta}_0)$$

## Model-based estimation:

- Take note that if we believe the mean model and the mean-variance relationship are correct, there is cancellation:  $\mathbb{B}_N(\hat{\boldsymbol{\beta}}) = \hat{\phi} \mathbb{A}_N(\hat{\boldsymbol{\beta}})$ .
- The variance estimator based on the assumption of both a correct mean model and mean-variance relationship collapses as follows:

$$\begin{aligned} \widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] &= (\mathbb{A}_N(\hat{\boldsymbol{\beta}}))^{-1} \mathbb{B}_N(\hat{\boldsymbol{\beta}}) (\mathbb{A}_N(\hat{\boldsymbol{\beta}}))^{-1} \\ &= \hat{\phi} (\mathbb{A}_N(\hat{\boldsymbol{\beta}}))^{-1}. \end{aligned}$$

which is exactly the formula based on the Fisher information!

## Model-based estimation:

- We previously saw that the solution to the GLM is determined by a user-specified mean model and mean-variance relationship. We have just argued the same for an asymptotically valid variance estimator.
- We are working within the theory of estimating equations—*not* likelihood theory—so what we have effectively just argued is that the likelihood-based approach will be valid so long as the mean model and mean-variance relationship are correctly specified.
- Though we often use likelihood language to describe a GLM, we don't rely on the third (and higher) moments implied by that likelihood.
- This is not your first exposure to the concept described on the previous slide; we already have seen this in OLS (which, as we also know, is a specific example of a GLM).

## Example 12.1: Derivation of variance for OLS

- Consider OLS linear regression, which may be thought of:
  - ▶ Parametrically: A normal GLM with the canonical link.
  - ▶ Semi-parametrically: A GLM based on the mean model  $E[\mathbf{y}|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$  and a working mean-variance relationship  $\mathbf{V} = \mathbf{I}$ .
- Let's use our findings to derive the sandwich variance estimator:
  - ▶  $\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta}) = \mathbb{A}_N(\boldsymbol{\beta}) = \mathbf{X}^T\mathbf{X}$ .
  - ▶  $\mathbb{B}_N^{\text{obs}}(\boldsymbol{\beta}) = \mathbf{X}^T \text{diag}(y_i - \mu_i(\boldsymbol{\beta}))^2 \mathbf{X}$ .
  - ▶ Sandwich:  $\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \text{diag}(y_i - \mathbf{x}_i^T\widehat{\boldsymbol{\beta}})^2 \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}$
- If we believe the mean model and the mean-variance relationship, we would replace the meat (or cheese, or veggies) of the sandwich with  $\widehat{\phi}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X}) = \widehat{\sigma}^2(\mathbf{X}^T\mathbf{I}\mathbf{X})$ :

$$\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = \widehat{\sigma}^2(\mathbf{X}^T\mathbf{X})^{-1}.$$

- We recognize this formula!

## Example 12.2: Derivation of variance for WLS

- Consider WLS linear regression based on  $\text{Var}[Y|\mathbf{X}] \propto \mathbf{V}(\mathbf{X})$ , which may be thought of:
  - ▶ Parametrically: A normal GLM with the canonical link, and a nuisance parameter  $\phi_i$  that depends upon  $\mathbf{X}_i$ .
  - ▶ Semi-parametrically: A GLM based on the mean model  $E[y|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$  and a working mean-variance relationship,  $\mathbf{V} = \mathbf{V}(\mathbf{X})$ .
- I leave it to you to verify the following sandwich variance estimator:

$$\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \text{diag}(y_i - \mathbf{x}_i^T \widehat{\boldsymbol{\beta}})^2 \mathbf{W} \mathbf{X} (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1},$$

where  $\mathbf{W} = \mathbf{V}^{-1}(\mathbf{X})$ .

- If we believe the mean model and the mean-variance relationship, we replace the meat/cheese/veggies with  $\widehat{\phi}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}) = \widehat{\phi}(\mathbf{X}^T \mathbf{W} \mathbf{X})$ :

$$\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = \left( \frac{1}{N - K} \sum_{i=1}^N w_i (y_i - \mathbf{x}_i^T \widehat{\boldsymbol{\beta}})^2 \right) (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}.$$

## Variance formulas so far: Canonical link

- $\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = \widehat{\phi}(\mathbb{A}_N(\widehat{\boldsymbol{\beta}}))^{-1}$ .
  - ▶ Relies on correct mean model.
  - ▶ Relies on correct mean-variance relationship.
  - ▶ Estimation of  $\widehat{\phi}$  not necessary if there is no nuisance parameter.
- $\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = (\mathbb{A}_N(\widehat{\boldsymbol{\beta}}))^{-1} \mathbb{B}_N^{\text{obs}}(\widehat{\boldsymbol{\beta}}) (\mathbb{A}_N(\widehat{\boldsymbol{\beta}}))^{-1}$ .
  - ▶  $\mathbb{A}_N(\widehat{\boldsymbol{\beta}}) = \mathbb{A}_N^{\text{obs}}(\widehat{\boldsymbol{\beta}})$ .
  - ▶ Does not rely on correct mean model.
  - ▶ Does not rely on correct mean-variance relationship.

## Variance formulas so far: Non-canonical link

- $\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = \hat{\phi}(\mathbb{A}_N(\hat{\boldsymbol{\beta}}))^{-1}$ .
  - ▶ Relies on correct mean model.
  - ▶ Relies on correct mean-variance relationship.
  - ▶ Estimation of  $\hat{\phi}$  not necessary if there is no nuisance parameter.
- $\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = (\mathbb{A}_N(\hat{\boldsymbol{\beta}}))^{-1} \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) (\mathbb{A}_N(\hat{\boldsymbol{\beta}}))^{-1}$ .
  - ▶  $\mathbb{A}_N(\hat{\boldsymbol{\beta}}) \neq \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$ .
  - ▶ Relies on correct mean model.
  - ▶ Does not rely on correct mean-variance relationship.
- $\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = (\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}))^{-1} \mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) (\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}))^{-1}$ .
  - ▶  $\mathbb{A}_N(\hat{\boldsymbol{\beta}}) \neq \mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$ .
  - ▶ Does not rely on correct mean model.
  - ▶ Does not rely on correct mean-variance relationship.

## Other considerations:

- There are other versions of the sandwich variance; most are simply modifications to the ones we've already discussed.
  - ▶ In the language of the `sandwich()` package in R, the ones we've discussed fall under the category of `HCO`.
- Some re-scale by a factor of  $N/(N - K)$  to add a correction for degrees of freedom.
- Some studentize the residuals of  $\mathbb{B}_N^{\text{obs}}(\hat{\beta})$ .
- In large samples, discrepancies across these versions are comparatively minor.

## Example 12.3: Normal distribution (log link)

- Revisiting a prior example, suppose our GLM is based on:
  - $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ .
  - $g(\cdot)$  given by the log link (i.e.,  $g(\mu) = \log(\mu)$ ).
- We derived the following estimating equations for  $\boldsymbol{\beta}$ :

$$\mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta})) (\mathbf{y} - \text{vec}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}.$$

- We determined the following:

$$\mathbb{A}_N(\hat{\boldsymbol{\beta}}) = \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))^2 \mathbf{X}.$$

$$\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}}) = \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))^2 \mathbf{X} - \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}})) (y_i - \exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}})) \mathbf{X}.$$

- I leave it for you to show that:

$$\mathbb{B}_N^{\text{obs}}(\boldsymbol{\beta}) = \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}})) (y_i - \exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))^2 \mathbf{X}.$$

## Example 12.3: Normal distribution (log link)

- We continue from our previous fit of the (simulated) data from this GLM (i.e., in Set 13).
- Consider the following variance estimators:
  - ① A “model-based” estimator.
    - ★ Assumes mean model and mean-variance relationship correct.
  - ② An estimator that allows a misspecified mean-variance relationship.
    - ★ But assumes the mean model is correct!
  - ③ An estimator that allows a misspecified mean model and mean-variance relationship.

## Example 12.3: Normal distribution (log link)

```
1 ## Store final iteration
2 betahat <- betaj
3
4 ## Linear predictor
5 etahat <- c(X %*% betahat)
6
7 ## Estimating function
8 Gn <- t(X * c(exp(etahat))) %*% (y - exp(etahat))
9
10 ## An
11 An <- t(X * c(exp(etahat)^2)) %*% X
12
13 ## W
14 W <- c(exp(etahat))
15
16 ## AnObs
17 AnObs <- An - t(X * W * c(y - exp(etahat))) %*% X
18
19 ## Bn
20 Bn <- t(X * W^2 * c(y - exp(etahat))^2) %*% X
```

## Example 12.3: Normal distribution (log link)

```
1 V1 <- phi * solve (An)
2 V2 <- solve (An) %*% Bn %*% solve (An)
3 V3 <- solve (AnObs) %*% Bn %*% solve (AnObs)
4 V2star <- V2 * (n) / (n - 2)
5 V3star <- V3 * (n) / (n - 2)
6
7 > sqrt (diag (V1))
8 [1] 0.1804808 0.1373220
9
10 > sqrt (diag (V2))
11 [1] 0.1649462 0.1077723
12
13 > sqrt (diag (V3))
14 [1] 0.1650517 0.1079106
15
16 > sqrt (diag (V2star))
17 [1] 0.1683475 0.1099946
18
19 > sqrt (diag (V3star))
20 [1] 0.1684552 0.1101358
```

## Example 12.3: Normal distribution (log link)

```
1 zz <- glm(y ~ X[,2], start = c(log(mean(y)), 0), family =  
  gaussian(link = "log"))  
2 V4 <- sandwich(zz)  
3  
4 > sqrt(diag(V4))  
5 (Intercept)      X[, 2]  
6 0.1649462      0.1077723
```

- And just like that, we've taken some of the magic away from the `sandwich()` function in R! What have we learned about which version of the sandwich is being used?

## Design matrix: Fixed vs. random

- The sandwich variance estimator(s) were derived based on the theory of estimating equations under the premise that sampling is from the joint distribution  $(\mathbf{X}, \mathbf{y})$ .
- If the mean model is correct, the validity of the sandwich still holds even if  $\mathbf{X}$  is fixed.
  - ▶ The argument for this lies in showing that when the mean model is correct,  $\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$ ,  $\mathbb{A}_N(\hat{\boldsymbol{\beta}})$ , and  $\mathbb{B}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$  are all consistent for the same (respective) quantities for which they are consistent when  $\mathbf{X}$  is random.
  - ▶ This will *not* be the case when the mean model is misspecified.
- We didn't have to care so much about this when we were thinking of GLMs through the likelihood framework, in which we were always assuming the model to be correctly specified.

# VARIANCE BASED ON THEORY OF ESTIMATING EQUATIONS

Summarizing what we know so far:

<b>X</b>	Link	MM	<b>V</b>	$\hat{\phi} \mathbb{A}_N^{-1}$	$(\mathbb{A}_N^{-1}) \mathbb{B}_N^{\text{obs}} (\mathbb{A}_N)^{-1}$	$(\mathbb{A}_N^{\text{obs}})^{-1} \mathbb{B}_N^{\text{obs}} (\mathbb{A}_N^{\text{obs}})^{-1}$
Fixed	C	✓	✓	✓	✓	✓
Fixed	C	✓	✗	✗	✓	✓
Fixed	C	✗	✓	✗	✗	✗
Fixed	C	✗	✗	✗	✗	✗
Fixed	NC	✓	✓	✓	✓	✓
Fixed	NC	✓	✗	✗	✓	✓
Fixed	NC	✗	✓	✗	✗	✗
Fixed	NC	✗	✗	✗	✗	✗
Random	C	✓	✓	✓	✓	✓
Random	C	✓	✗	✗	✓	✓
Random	C	✗	✓	✗	✓	✓
Random	C	✗	✗	✗	✓	✓
Random	NC	✓	✓	✓	✓	✓
Random	NC	✓	✗	✗	✓	✓
Random	NC	✗	✓	✗	✗	✓
Random	NC	✗	✗	✗	✗	✓

Link function,  $g(\cdot)$ : C=Canonical; NC=Non-canonical

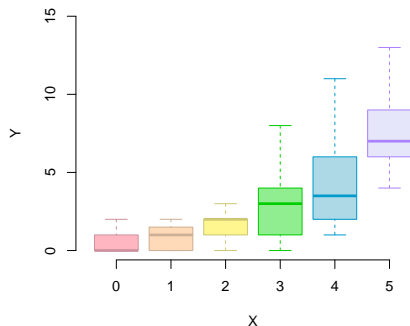
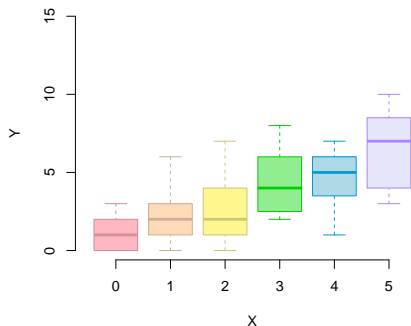
MM: Mean model correctly specified? ✓=Yes; ✗=No

**V**: Mean-variance relationship correctly specified? ✓=Yes; ✗=No

## Mean models study design: A simulation study

- Since seeing is believing, let's further conduct a simulation study to illustrate the interaction between mean model misspecification and the fixed/random nature of the design matrix,  $\mathbf{X}$ .
- Suppose that  $X$  takes on discrete values 0 through 5, and that  $Y \sim \text{Poisson}(\lambda = \exp(\beta_0 + \beta_1 1(x=1) + \dots + \beta_5 1(x=5)))$
- Study design scenarios:
  - ▶ Random: Equal probability  $p = 1/6$  of allocation to each category.
  - ▶ Fixed: Even allocation (1/6 observations in each category).
- Mean model scenarios:
  - ▶ Linear model correct:  $\boldsymbol{\beta} = (\log(1), \log(2), \log(3), \log(4), \log(5), \log(6))$ .
  - ▶ Linear model incorrect:  $\boldsymbol{\beta} = (-1/2, 1/2, 2/2, 3/2, 4/2, 5/2)$ .
- Suppose we fit OLS (Gaussian GLM with identity link) treating  $X$  continuously/linearly in each of the four scenarios.

## Mean models study design: A simulation study



- Regardless of whether  $X$  is fixed or random,  $E[Y|X = x] = \beta_0 + \beta_1 x$  is correct on the left, but incorrect on the right. The mean-variance relationship is not correctly captured by OLS in any case.

## Mean models study design: Random $X$ and correct mean model

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Sample size
5 n <- 120
6
7 ## Set true value of beta (mean model correct)
8 beta <- c(log(1), log(2), log(3), log(4), log(5), log(6))
9
10 ## Number of simulation replicates
11 nsim <- 10000
12
13 ## Store coefficients (column 1) and sandwich SEs (column 2)
14 res <- matrix(0, nrow = nsim, ncol = 2)
```

## Mean models study design: Random $X$ and correct mean model

```

1 for (j in 1:nsim) {
2   ## Generate X (random)
3   x <- factor(sample(c(0:5), size = n, replace = TRUE))
4   X <- matrix(model.matrix(~x), ncol = 6)
5
6   ## Generate Y
7   Y <- rpois(n, lambda = exp(X %*% beta))
8
9   ## Treat X continuously
10  x.c <- as.numeric(x)
11  X.c <- matrix(cbind(1, x.c), ncol = 2)
12
13  ## OLS
14  zz <- lm(Y ~ x.c)
15
16  ## Sandwich
17  An <- t(X.c) %*% X.c
18  Bn <- t(X.c * c(Y - c(X.c %*% coef(zz)))^2) %*% X.c
19  vhat <- solve(An) %*% Bn %*% solve(An)
20
21  ## Extract results
22  res[j,] <- c(coef(zz)[2], sqrt(diag(vhat))[2])
23 }
24
25 ## No problem!
26 > apply(res, 2, mean)
27 [1] 0.9991057 0.0990592
28
29 > apply(res, 2, sd)
30 [1] 0.1013802 0.0114567

```

## Mean models study design: Random $X$ and incorrect mean model

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Sample size
5 n <- 120
6
7 ## Set true value of beta (mean model incorrect)
8 beta <- c(-1/2, 1/2, 2/2, 3/2, 4/2, 5/2)
9
10 ## Number of simulation replicates
11 nsim <- 10000
12
13 ## Store coefficients (column 1) and sandwich SEs (column 2)
14 res <- matrix(0, nrow = nsim, ncol = 2)
```

## Mean models study design: Random $X$ and incorrect mean model

```

1 for (j in 1:nsim) {
2   ## Generate X (random)
3   x <- factor(sample(c(0:5), size = n, replace = TRUE))
4   X <- matrix(model.matrix(~x), ncol = 6)
5
6   ## Generate Y
7   Y <- rpois(n, lambda = exp(X %*% beta))
8
9   ## Treat X continuously
10  x.c <- as.numeric(x)
11  X.c <- matrix(cbind(1, x.c), ncol = 2)
12
13  ## OLS
14  zz <- lm(Y ~ x.c)
15
16  ## Sandwich
17  An <- t(X.c) %*% X.c
18  Bn <- t(X.c * c(Y - c(X.c %*% coef(zz)))^2) %*% X.c
19  vhat <- solve(An) %*% Bn %*% solve(An)
20
21  ## Extract results
22  res[j,] <- c(coef(zz)[2], sqrt(diag(vhat))[2])
23 }
24
25 ## No problem!
26 > apply(res, 2, mean)
27 [1] 1.295252 0.111242
28
29 > apply(res, 2, sd)
30 [1] 0.1134845 0.0141525

```

## Mean models study design: Fixed $\mathbf{X}$ and correct mean model

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Sample size
5 n <- 120
6
7 ## Fixed design matrix (even allocation)
8 x <- factor(c(rep(0, n/6), rep(1, n/6), rep(2, n/6),
9             rep(3, n/6), rep(4, n/6), rep(5, n/6)))
10 X <- matrix(model.matrix(~x), ncol = 6)
11
12 ## Treat continuously X
13 x.c <- as.numeric(x)
14 X.c <- matrix(cbind(1, x.c), ncol = 2)
15
16 ## Set true value of beta (mean model correct)
17 beta <- c(log(1), log(2), log(3), log(4), log(5), log(6))
18
19 ## Number of simulation replicates
20 nsim <- 10000
21
22 ## Store coefficients (column 1) and sandwich SEs (column 2)
23 res <- matrix(0, nrow = nsim, ncol = 2)
```

## Mean models study design: Fixed $X$ and correct mean model

```

1 for (j in 1:nsim) {
2   ## Generate Y
3   Y <- rpois(n, lambda = exp(X %*% beta))
4
5   ## OLS
6   zz <- lm(Y ~ x.c)
7
8   ## Sandwich
9   An <- t(X.c) %*% X.c
10  Bn <- t(X.c * c(Y - c(X.c %*% coef(zz)))^2) %*% X.c
11  vhat <- solve(An) %*% Bn %*% solve(An)
12
13  ## Extract results
14  res[j,] <- c(coef(zz)[2], sqrt(diag(vhat))[2])
15 }
16
17 ## No problem!
18 > apply(res, 2, mean)
19 [1] 1.0001101 0.0982242
20
21 > apply(res, 2, sd)
22 [1] 0.0991858 0.0105809

```

## Mean models study design: Fixed $\mathbf{X}$ and incorrect mean model

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Sample size
5 n <- 120
6
7 ## Fixed design matrix (even allocation)
8 x <- factor(c(rep(0, n/6), rep(1, n/6), rep(2, n/6),
9             rep(3, n/6), rep(4, n/6), rep(5, n/6)))
10 X <- matrix(model.matrix(~x), ncol = 6)
11
12 ## Treat continuously X
13 x.c <- as.numeric(x)
14 X.c <- matrix(cbind(1, x.c), ncol = 2)
15
16 ## Set true value of beta (mean model incorrect)
17 beta <- c(-1/2, 1/2, 2/2, 3/2, 4/2, 5/2)
18
19 ## Number of simulation replicates
20 nsim <- 10000
21
22 ## Store coefficients (column 1) and sandwich SEs (column 2)
23 res <- matrix(0, nrow = nsim, ncol = 2)
```

## Mean models study design: Fixed $X$ and incorrect mean model

```

1 for (j in 1:nsim) {
2   ## Generate Y
3   Y <- rpois(n, lambda = exp(X %*% beta))
4
5   ## OLS
6   zz <- lm(Y ~ x.c)
7
8   ## Sandwich
9   An <- t(X.c) %*% X.c
10  Bn <- t(X.c * c(Y - c(X.c %*% coef(zz)))^2) %*% X.c
11  vhat <- solve(An) %*% Bn %*% solve(An)
12
13  ## Extract results
14  res[j,] <- c(coef(zz)[2], sqrt(diag(vhat))[2])
15 }
16
17 ## Problem!!!
18 > apply(res, 2, mean)
19 [1] 1.298859 0.110562
20
21 > apply(res, 2, sd)
22 [1] 0.1011173 0.0130098
23
24 ## This difference may not seem that big, but recall that the violation to linearity
25 ## reflected by this example was pretty modest...

```

# TABLE OF CONTENTS

- 1 Variance based on likelihood theory
- 2 Variance based on theory of estimating equations
- 3 The nonparametric bootstrap

# THE NONPARAMETRIC BOOTSTRAP

## Main ideas:

- Let  $F$  denote CDF for  $(\mathbf{X}, Y)$  or  $(Y|\mathbf{X})$ , depending on context; let  $\mathbb{F}_N$  denote empirical CDF based on  $N$  observations.
  - ▶  $\boldsymbol{\beta} = T(F)$ , and hence  $\hat{\boldsymbol{\beta}} = T(\mathbb{F}_N)$ .
  - ▶ Absent parametric form,  $\mathbb{F}_N$  is our best estimate of  $F$ .
- Repeat-sampling from  $\mathbb{F}_N$  with replacement gives information on distribution of  $\hat{\boldsymbol{\beta}}^* = T(\mathbb{F}_N^*)$ ; asterisk denotes fixed  $\mathbb{F}_N$ .
- Let  $\{\hat{\boldsymbol{\beta}}_b^*\}_{b=1}^B$  denote the (bootstrap) samples.
- Note two layers of variation:
  - ▶ How well  $\mathbb{F}_N$  approximates  $F$ .
    - ★ Glivenko-Cantelli:  $\sup_t |F(t) - \mathbb{F}_N(t)| \xrightarrow{\text{a.s.}} 0$  as  $N \nearrow \infty$ .
  - ▶ How well  $\{\hat{\boldsymbol{\beta}}_b^*\}_{b=1}^B$  approximates  $T(\mathbb{F}_N^*)$ .
    - ★ Better as  $B \nearrow \infty$ .
- Which source of variation can we control once given the data?

## Estimator-attributed bias:

- Let  $\hat{\beta}_b^* = T(\mathbb{F}_{N:b}^*)$  denote estimate based on  $b^{\text{th}}$  bootstrap sample. We may estimate bias as follows:

$$\begin{aligned}\widehat{\text{Bias}} &= \frac{1}{B} \sum_{b=1}^B (T(\mathbb{F}_{N:b}^*) - T(\mathbb{F}_N)) \\ &= \frac{1}{B} \sum_{b=1}^B \hat{\beta}_b^* - \hat{\beta} = \hat{\beta}^* - \hat{\beta} \approx \hat{\beta} - \beta,\end{aligned}$$

where  $\hat{\beta}^* = \frac{1}{B} \sum_{b=1}^B \hat{\beta}_b^*$ .

- Correction won't catch external sources of bias; be warned.

## Covariance:

- We may estimate the covariance as well:

$$\widehat{\text{Cov}}[\widehat{\boldsymbol{\beta}}] = \frac{1}{B-1} \sum_{b=1}^B (\widehat{\boldsymbol{\beta}}_b^* - \widehat{\boldsymbol{\beta}}^*)(\widehat{\boldsymbol{\beta}}_b^* - \widehat{\boldsymbol{\beta}}^*)^\top$$

- For the  $k^{\text{th}}$  coefficient, we have:

$$\widehat{v}_k = \widehat{\text{Var}}[\widehat{\beta}_k] = \frac{1}{B-1} \sum_{b=1}^B ([\widehat{\boldsymbol{\beta}}_b^*]_k - \widehat{\beta}_k^*)^2$$

## Confidence intervals: Normal approximation (bias-correction)

- Symmetric  $(1 - \alpha)$  CI:

$$(\widehat{\beta}_k - \widehat{\text{Bias}}_k) \pm \sqrt{\widehat{v}_k} z_{1-\alpha/2}.$$

- Assumptions:

- ▶  $\widehat{\beta}_k - \beta_k \sim \mathcal{N}(\widehat{\text{Bias}}_k, \widehat{v}_k)$ , which is symmetric and pivotal.
- ▶  $\widehat{\text{Bias}}_k$  and  $\widehat{v}_k$  are good estimates of  $\text{Bias}_k$  and  $\sigma^2$ .

- Good for cases where  $N$  is large enough that normal approximation holds, but no known theoretical formula for asymptotic variance.
- Can use QQ-plots to evaluate departures from normality.

## Confidence intervals: Pivot based

- Let  $\hat{\beta}_{k(p)}^*$  denote  $p^{\text{th}}$  quantile of  $k^{\text{th}}$  coefficient of  $\{\hat{\beta}_b^*\}_{b=1}^B$ .
- Behavior of  $\beta_k - \hat{\beta}_k$  approximately that of  $\hat{\beta}_k - \hat{\beta}_k^*$ :

$$\begin{aligned}
 0.95 &\approx P\left(\hat{\beta}_{k(\alpha/2)}^* \leq \hat{\beta}_k^* \leq \hat{\beta}_{k(1-\alpha/2)}^*\right) \\
 &= P\left(\hat{\beta}_k - \hat{\beta}_{k(1-\alpha/2)}^* \leq \hat{\beta}_k - \hat{\beta}_k^* \leq \hat{\beta}_k - \hat{\beta}_{k(\alpha/2)}^*\right) \\
 &\approx P\left(\hat{\beta}_k - \hat{\beta}_{k(1-\alpha/2)}^* \leq \beta_k - \hat{\beta}_k \leq \hat{\beta}_k - \hat{\beta}_{k(\alpha/2)}^*\right) \\
 &= P\left(2\hat{\beta}_k - \hat{\beta}_{k(1-\alpha/2)}^* \leq \beta_k \leq 2\hat{\beta}_k - \hat{\beta}_{k(\alpha/2)}^*\right)
 \end{aligned}$$

- Assumptions:
  - ▶  $\hat{\beta}_k - \beta_k$  asymptotically pivotal (not necessarily symmetric).

## Confidence intervals: Quantile-based

- Let  $\hat{\beta}_{k(p)}^*$  denote  $p^{\text{th}}$  quantile of  $k^{\text{th}}$  coefficient of  $\{\hat{\beta}_b^*\}_{b=1}^B$ .
- One can form a  $100(1 - \alpha)\%$  CI as:

$$[\hat{\beta}_{k(\alpha/2)}^*, \hat{\beta}_{k(1-\alpha/2)}^*].$$

- Assumptions:
  - ▶ There is a monotone  $h(\cdot)$  for which the distribution of  $h(\hat{\beta}_k^*)$  is symmetric, and that  $h(\hat{\beta}_k^*)$  is pivotal.
  - ▶  $h(\hat{\beta}_k)$  is unbiased.

## **Linear regression:** Bootstrap procedures

- The following three slides outline reasonable bootstrap procedures for linear regression; all but one will generalize to GLMs.

## **Linear regression:** Random design

- Re-sample pairs  $(\mathbf{x}_i^*, y_i^*)$  from existing observations  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  with replacement.
- Estimate  $\hat{\boldsymbol{\beta}}_b^*$  for  $b = 1, \dots, B$ ; form estimates/confidence intervals of your choosing from prior methods.
- Design changes with each sample.
- Consistent with an observational study with random sampling irrespective of exposure/outcome.
- Consistent with fully/purely randomized experiment (like a coin toss).

## Linear regression: Fixed design

- Fit model  $E[\mathbf{y}|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$ .
- Stratify unconditional bootstrap procedure by discrete subgroups defined by  $\mathbf{X}$ .
- Estimate  $\hat{\boldsymbol{\beta}}_b^*$  for  $b = 1, \dots, B$ ; form estimates/confidence intervals of your choosing from prior methods.
- Allows heteroscedasticity; allows mean model misspecification.
- Example: designed experiment with a small number of large groups.

**Linear regression:** Fixed design (correct mean model, homoscedasticity)

- Fit model  $E[\mathbf{y}|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$  and extract residuals  $\{\widehat{\epsilon}_i\}_{i=1}^N$ .
- Re-sample  $N$  residuals  $\widehat{\epsilon}_i^*$  with replacement.
- Keep  $\mathbf{x}_i$  intact; form new outcomes  $y_i^* = \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}} + \widehat{\epsilon}_i^*$  for  $i = 1, \dots, N$ .
- Estimate  $\widehat{\boldsymbol{\beta}}_b^*$  for  $b = 1, \dots, B$ ; form estimates/confidence intervals of your choosing from prior methods.
- Assumptions:
  - ▶ Homoscedasticity of errors.
  - ▶ Correct mean-model.
- Example: designed experiment with many discrete categories of  $\mathbf{X}$  that each have relatively small samples.
- Does not generalize to GLMs (can't necessarily form  $y_i^*$  based on  $\widehat{\epsilon}_i^*$ ).

## **More generally:** Fixed vs. random design

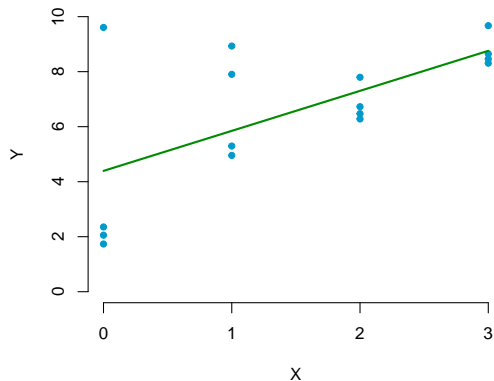
- If the mean model is correct, either version of the bootstrap should perform well regardless of whether  $\mathbf{X}$  is fixed or random.
- If  $\mathbf{X}$  is fixed, mean-model misspecification will tend to result in an overstated variance if you treat  $\mathbf{X}$  as random.
- If  $\mathbf{X}$  is random, mean-model misspecification will tend to result in an understated variance if you treat  $\mathbf{X}$  as fixed.

## Example 12.4: Bootstrap under high error skewness

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Set sample size
5 n <- 16
6
7 ## Generate predictor
8 x <- c(rep(0,n/4), rep(1,n/4), rep(2,n/4), rep(3,n/4))
9
10 ## Generate outcome (linearity correct)
11 y <- 1 + 3*(x == 1) + 5*(x == 2) + 7*(x == 3) + rexp(n, 1/2)
12
13 ## Create data frame
14 dat <- data.frame(cbind(x, y))
15
16 ## Analysis on original data
17 zz <- lm(y ~ x, data = dat)
18
19 ## Estimated slope
20 bhat <- as.numeric(zz$coef[2])
21
22 ## Set bootstrap replicates
23 B <- 5000
```

# THE NONPARAMETRIC BOOTSTRAP

**Example 12.4:** Bootstrap under high error skewness



# THE NONPARAMETRIC BOOTSTRAP

## Example 12.4: Bootstrap under high error skewness

```
1 ## Model-based and sandwich-based standard errors
2 se.model <- as.numeric(sqrt(diag(vcov(zz)))[2])
3 se.sandwich <- as.numeric(sqrt(diag(sandwich(zz)))[2])
4
5 ## Model-based CI
6 > c(EST = bhat,
7 +   CILO = bhat - qnorm(0.975)*se.model,
8 +   CIHI = bhat + qnorm(0.975)*se.model)
9     EST      CILO      CIHI
10 1.4550308 0.5314966 2.3785650
11
12 ## Sandwich-based CI
13 > c(EST = bhat,
14 +   CILO = bhat - qnorm(0.975)*se.sandwich,
15 +   CIHI = bhat + qnorm(0.975)*se.sandwich)
16     EST      CILO      CIHI
17 1.4550308 0.4480281 2.4620335
```

# THE NONPARAMETRIC BOOTSTRAP

## Example: Bootstrap under high error skewness

```
1 ## BOOTSTRAP METHOD 1: FULL-RESAMPLING
2
3 ## Create a place to store results
4 b.results <- matrix(0, nrow = B, ncol = 1)
5
6 ## Conduct bootstrap samples
7 for (j in 1:B)
8 {
9   ## Random sample with replacement with original sample size in mind
10  samp <- sample(1:n, replace = TRUE)
11  bdat <- dat[samp,]
12
13  ## Run model on bootstrap sample
14  bzz <- lm(y ~ x, data = bdat)
15
16  ## Extract results
17  b.results[j,1] <- coef(bzz)[2]
18 }
```

# THE NONPARAMETRIC BOOTSTRAP

## Example 12.4: Bootstrap under high error skewness

```
1 ## Bootstrap standard error
2 > sd(b.results)
3 [1] 0.5314059
4
5 ## Symmetric large-sample-justified CI
6 > c(CILO = mean(b.results) - qnorm(0.975)*sd(b.results),
7 +   CIHI = mean(b.results) + qnorm(0.975)*sd(b.results))
8     CILO      CIHI
9 0.4049048 2.4879778
10
11 ## Asymmetric pivot-based CI
12 qlo <- quantile(b.results, 0.025)
13 qhi <- quantile(b.results, 0.975)
14 > c(CILOW = as.numeric(2*mean(b.results) - qhi),
15 +   CIHI = as.numeric(2*mean(b.results) - qlo))
16     CILOW      CIHI
17 0.6392062 2.6397876
18
19 ## Quantile-based CI
20 > c(CILOW = as.numeric(quantile(b.results, c(0.025))),
21 +   CIHI = as.numeric(quantile(b.results, c(0.975))))
22     CILOW      CIHI
23 0.253095 2.253676
```

# THE NONPARAMETRIC BOOTSTRAP

## Example 12.4: Bootstrap under high error skewness

```
1 ## BOOTSTRAP METHOD 2: CONDITIONAL (FIXED X)
2
3 ## Set bootstrap replicates
4 B <- 5000
5
6 ## Create a place to store results
7 b.results <- matrix(0, nrow = B, ncol = 1)
8
9 ## Extract residuals from fitted model
10 rsdls <- zz$residuals
11
12 ## Keep a "fixed" version of the exposure
13 x.fixed <- dat$x
14
15 ## Extract estimate of beta
16 bhat <- as.numeric(zz$coef)
```

# THE NONPARAMETRIC BOOTSTRAP

## Example 12.4: Bootstrap under high error skewness

```
1 for (j in 1:B)
2 {
3   ## Random sample of residuals with replacement
4   samp <- sample(1:n, replace = TRUE)
5   brsdls <- rsdls[samp]
6
7   ## Append residuals to create a bootstrap FEV
8   by <- bhat[1] + bhat[2]*x.fixed + brsdls
9
10  ## Create bootstrap data set
11  bdat <- data.frame(cbind(x.fixed, by))
12
13  ## Run model on bootstrap sample
14  bzz <- lm(by ~ x.fixed, data = bdat)
15
16  ## Extract results
17  b.results[j,1] <- coef(bzz)[2]
18 }
```

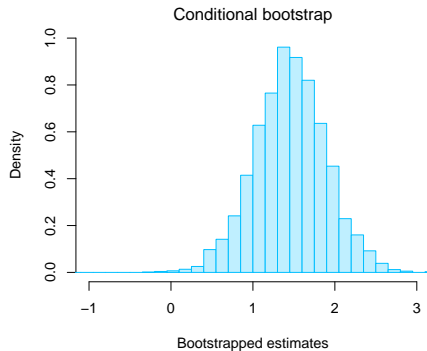
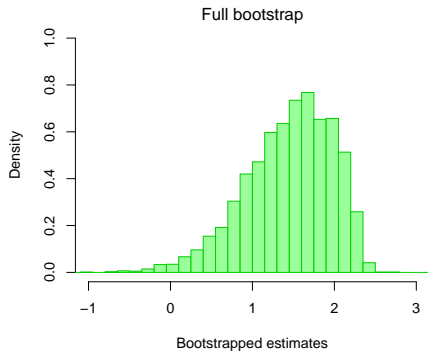
# THE NONPARAMETRIC BOOTSTRAP

## Example 12.4: Bootstrap under high error skewness

```
1 ## Bootstrap standard error
2 > as.numeric(sd(b.results))
3 [1] 0.4352239
4
5 ## Symmetric large-sample-justified CI
6 > c(CILO = mean(b.results) - qnorm(0.975)*sd(b.results),
7 +   CIHI = mean(b.results) + qnorm(0.975)*sd(b.results))
8     CILO      CIHI
9 0.603114 2.309160
10
11 ## Asymmetric pivot-based CI
12 qlo <- quantile(b.results, 0.025)
13 qhi <- quantile(b.results, 0.975)
14 > c(CILOW = as.numeric(2*mean(b.results) - qhi),
15 +   CIHI = as.numeric(2*mean(b.results) - qlo))
16     CILOW      CIHI
17 0.5893383 2.3304830
18
19 ## Quantile-based CI
20 > c(CILOW = as.numeric(quantile(b.results, c(0.025))),
21 +   CIHI = as.numeric(quantile(b.results, c(0.975))))
22     CILOW      CIHI
23 0.5817912 2.3229358
```

# THE NONPARAMETRIC BOOTSTRAP

## Example 12.4: Bootstrap under high error skewness



## Example 12.4: Bootstrap under high error skewness

- Why does the distribution of the bootstrapped estimates look so different between the two approaches?
- One point is clearly highly influential!
- Probability of inclusion in a single unconditional bootstrap replicate:

$$P(\geq 1 \text{ Inclusion}) = 1 - (15/16)^{16} \approx 0.64.$$

$$P(\geq 2 \text{ Inclusions}) = 1 - (15/16)^{16} - 15(15/16)^{15}(1/16) \approx 0.26.$$

$$P(\geq 3 \text{ Inclusions}) \approx 0.074.$$

$$P(\geq 4 \text{ Inclusions}) \approx 0.015.$$

$$P(\geq 5 \text{ Inclusions}) \approx 0.0023.$$

- Probability of inclusion in a single conditional bootstrap replicate:

$$P(\geq 1 \text{ Inclusion}) = 1 - (3/4)^4 \approx 0.68.$$

$$P(\geq 2 \text{ Inclusions}) = 1 - (3/4)^4 - 4(3/4)^3(1/4) \approx 0.26.$$

$$P(\geq 3 \text{ Inclusions}) \approx 0.051.$$

$$P(4 \text{ Inclusions}) \approx 0.0039.$$

$$P(\geq 5 \text{ Inclusions}) = 0$$

## Example 12.5: Clever uses of the bootstrap

- We can use the bootstrap to answer questions that would otherwise be difficult or impossible to analytically answer.
- As an example, consider the following two models (unadjusted and adjusted) using OLS linear regression:

$$\begin{aligned}E[Y|X = x] &= \alpha_0 + \alpha_1 x \\E[Y|X = x, Z = z] &= \beta_0 + \beta_1 x + \beta_2 z.\end{aligned}$$

- What is  $\text{Cov}[\hat{\alpha}_1, \hat{\beta}_1]$ ?

## Example 12.5: Set up simulation

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Set number of simulations
5 nsim <- 500
6
7 ## Set sample size
8 n <- 500
9
10 ## Set number of bootstrap replicates
11 B <- 100
12
13 ## Store results
14 res <- matrix(0, nrow = nsim, ncol = 3)
```

## Example 12.5: Simulation (part 1 - original data)

```
1 ## Conduct simulation
2 for (j in 1:nsim)
3 {
4   ## Generate predictors and outcomes
5   X <- runif(n, 0, 5)
6   Z <- runif(n, 0, 5)
7   Y <- 1 + X + Z + rnorm(n, 0, 5)
8
9   ## Fit adjusted and unadjusted models
10  XU <- cbind(1, X)
11  XA <- cbind(1, X, Z)
12  res[j,1] <- (solve(t(XU) %*% XU) %*% (t(XU) %*% Y) [2])
13  res[j,2] <- (solve(t(XA) %*% XA) %*% (t(XA) %*% Y) [2])
```

# THE NONPARAMETRIC BOOTSTRAP

## Example 12.5: Simulation (part 2 - bootstrap)

```
1  ## Store bootstrapped results
2  bres <- matrix(0, nrow = B, ncol = 2)
3
4  for (b in 1:B)
5  {
6    ## Full-size sample with replacement
7    samp <- sample(1:n, size = n, replace = TRUE)
8    bX <- X[samp]
9    bZ <- Z[samp]
10   bY <- Y[samp]
11
12   ## Fit adjusted and unadjusted models on bootstrapped data
13   bXU <- cbind(1, bX)
14   bXA <- cbind(1, bX, bZ)
15   bzz1 <- (solve(t(bXU) %*% bXU) %*% t(bXU) %*% bY) [2]
16   bzz2 <- (solve(t(bXA) %*% bXA) %*% t(bXA) %*% bY) [2]
17
18   ## Extracted data
19   bres[b,1] <- bzz1
20   bres[b,2] <- bzz2
21 }
```

## Example 12.5: Report results

```
1  ## Store estimated covariance
2  res[j,3] <- cov(bres[,1],bres[,2])
3
4  ## Track progress
5  if (round(j/50) == (j/50)) {print(paste(j, "sims complete!"))}
6 }
7
8 ## Average estimated covariance by bootstrap
9 > colMeans(res)[3]
10 [1] 0.02374372
11
12 ## Actual covariance by simulation
13 > cov(res[,1],res[,2])
14 [1] 0.0229229
```

## **This unit:**

- Theory of estimating equations.
- Sandwich variance methods.
- Bootstrap methods.

# SUMMARY: SO FAR

- Random vectors and matrices; multivariate normal theory.
- Ordinary least squares.
- Hypothesis testing and ANOVA.
- Weighted least squares.
- Misspecification.
- Confidence regions and prediction.
- Diagnostics.
- Regularization.
- Bayesian regression.
- Exponential families.
- Generalized linear models.
- Sandwich and bootstrap.

# SUMMARY: COMING UP

- Quasi-likelihood.
- Hypothesis testing for GLMs.
- Diagnostics for GLMs.
- Further considerations for binary outcomes.
- Nonlinear least squares.