

BIOS 7345: Advanced Regression for Independent Data

Andrew J. Spieker, Ph.D.

Associate Professor of Biostatistics
Vanderbilt University

Set 4: Weighted least squares

Version: 09/12/2025

TABLE OF CONTENTS

1 Weighted least squares

2 Iteratively re-weighted least squares

Recall: Ordinary least squares

- For this set of notes, assume \mathbf{X} is fixed and of full rank (though key results do not depend upon this).
- Model: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, with $\text{Cov}[\boldsymbol{\epsilon}] = \sigma^2\mathbf{I}$.
- OLS minimizes sum of squared errors (convex optimization problem):

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^K}{\text{minimize}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

- Leads to the normal equations:

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}.$$

- Closed-form expression for least squares estimator: $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$.
- Closed-form expression for variance: $\text{Cov}[\hat{\boldsymbol{\beta}}] = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$

WEIGHTED LEAST SQUARES

New setup: Diagonal variance

- Model: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, with $E[\boldsymbol{\epsilon}] = \mathbf{0}$ and $\text{Cov}[\boldsymbol{\epsilon}] = \sigma^2\mathbf{V}$.
- Specifically, $\text{Cov}[\boldsymbol{\epsilon}] = \sigma^2\mathbf{V}$ with $\sum_{i=1}^N \mathbf{V}_{ii} = 1$ (indeed, $\mathbf{V}_{ii} > 0$, and $\mathbf{V}_{ij} = 0$ for $i \neq j$, at least in BIOS 7345).
- The OLS estimate is unbiased, but OLS-based variance is incorrect.
- Solution: Let $\mathbf{K} = \mathbf{V}^{-1/2}$. Then:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \iff \mathbf{K}\mathbf{y} = \mathbf{K}\mathbf{X}\boldsymbol{\beta} + \mathbf{K}\boldsymbol{\epsilon}.$$

- Let $\mathbf{y}^* = \mathbf{K}\mathbf{y}$, $\mathbf{X}^* = \mathbf{K}\mathbf{X}$, and $\boldsymbol{\epsilon}^* = \mathbf{K}\boldsymbol{\epsilon}$.
- Then, $\mathbf{y}^* = \mathbf{X}^*\boldsymbol{\beta} + \boldsymbol{\epsilon}^*$, with:
 - ▶ $E[\boldsymbol{\epsilon}^*] = E[\mathbf{K}\boldsymbol{\epsilon}] = \mathbf{K}E[\boldsymbol{\epsilon}] = \mathbf{0}$.
 - ▶ $\text{Cov}[\boldsymbol{\epsilon}^*] = \text{Cov}[\mathbf{K}\boldsymbol{\epsilon}] = \mathbf{K}\text{Cov}[\boldsymbol{\epsilon}]\mathbf{K}^T = \mathbf{V}^{-1/2}\sigma^2\mathbf{V}\mathbf{V}^{-1/2} = \sigma^2\mathbf{I}$.
- Applying OLS on the transformed data gives:

$$\hat{\boldsymbol{\beta}}^* = ((\mathbf{X}^*)^T\mathbf{X}^*)^{-1}(\mathbf{X}^*)^T\mathbf{y}^* = (\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}.$$

Properties:

- Straightforward to show:
 - ▶ $\mathbf{P} = \mathbf{X}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}$ is an oblique projection matrix
 - ★ Idempotent but not necessarily symmetric.
 - ▶ $E[\hat{\boldsymbol{\beta}}^*] = \boldsymbol{\beta}$.
 - ▶ $\text{Cov}[\hat{\boldsymbol{\beta}}^*] = \sigma^2(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}$.
 - ▶ $E[\text{RSS}] = \sigma^2(N - K)$, where $\text{RSS} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^*)^T\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^*)$.
 - ★ This suggests estimating σ^2 as:

$$\hat{\sigma}^2 = \frac{1}{N - K} \sum_{i=1}^N \frac{(y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^*)^2}{V_i}$$

- If $\hat{\boldsymbol{\beta}}$ denotes the OLS estimate based on the original data, we also have $E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ and $\text{Cov}[\hat{\boldsymbol{\beta}}] = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}$.
- Since $\hat{\boldsymbol{\beta}}$ is unbiased, what is the motivation for $\hat{\boldsymbol{\beta}}^*$?

Theorem 4.1: The Gauss-Markov theorem (version 4)

Let \mathbf{X} be of full rank, and suppose $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, with $E[\boldsymbol{\epsilon}] = \mathbf{0}$ and $\text{Cov}[\boldsymbol{\epsilon}] = \sigma^2\mathbf{V}$. Then, for any constant vector \mathbf{a} , $\mathbf{a}^\top \hat{\boldsymbol{\beta}}^*$ is the unique BLUE of $\mathbf{a}^\top \boldsymbol{\beta}$ (where $\hat{\boldsymbol{\beta}}^*$ is the WLS estimator having weights given by $\mathbf{W} \propto \mathbf{V}^{-1}$).

- The proof of this version of the Gauss-Markov theorem relies on the version for OLS.

TABLE OF CONTENTS

- 1 Weighted least squares
- 2 Iteratively re-weighted least squares

Choosing weights:

- The premise of the previous results is, in some sense, a little bit cheap in the sense that we almost never know the form of \mathbf{V} exactly.
 - ▶ To be fair, OLS can be thought of as making the specific choice $\mathbf{V} = \mathbf{I}$, which we also almost never know to be true. Specifically having to correctly choose a \mathbf{V} that has a potentially unique entry along each entry of the main diagonal seems like even more of a stretch.
- The slides that follow suggest a procedure on how to deal with the case where you are willing to posit (and estimate) a form for the mean-variance relationship.

Unknown variance:

- Consider $\mathbf{V} = \mathbf{V}(\boldsymbol{\mu}(\boldsymbol{\beta}); \boldsymbol{\theta})$, so the mean is a function of $\boldsymbol{\mu} = E[\mathbf{y}]$ that depends upon $\boldsymbol{\beta}$ and possibly other unknown parameters $\boldsymbol{\theta}$.
 - ▶ We often focus on the case where $\mathbf{V} = \mathbf{V}(\boldsymbol{\mu})$.
- The structure of \mathbf{V} is assumed known, specified by the user.
- Beginning with some particular initializer allows us to *estimate* $\text{Cov}[\boldsymbol{\epsilon}]$, use that to inform our weights on the next round, and cycle through until apparent convergence. This procedure is known as iteratively re-weighted least squares (IRLS).
- If your model for $\text{Cov}[\boldsymbol{\epsilon}]$ is correct, resulting $\hat{\boldsymbol{\beta}}$ from IRLS is asymptotically BLUE for $\boldsymbol{\beta}$.

ITERATIVELY RE-WEIGHTED LEAST SQUARES

IRLS: Estimating $\boldsymbol{\beta}$ when $\text{Cov}[\boldsymbol{\epsilon}] = \mathbf{V}(\boldsymbol{\mu})$.

- 1 Set convergence threshold, $\delta_{\text{stop}} > 0$ (small).
- 2 If you initialize the value of $\boldsymbol{\beta}$ as $\boldsymbol{\beta}^{(0)}$, the next step in a first-order Taylor approximation is given by:

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

where $\mathbf{W}(\boldsymbol{\beta}^{(j)}) = \mathbf{V}^{-1}(\boldsymbol{\mu}(\boldsymbol{\beta}^{(j)}))$.

- 3 Iterate Step 2 until convergence: $\|\boldsymbol{\beta}^{(j+1)} - \boldsymbol{\beta}^{(j)}\|^2 < \delta_{\text{stop}}$
 - ▶ If $\boldsymbol{\beta}^{(0)}$ is initialized as a consistent estimator of $\boldsymbol{\beta}$ (e.g., OLS), you only need to take a single step in order to achieve an asymptotically efficient estimator. You'll see this again in BIOS 7346.

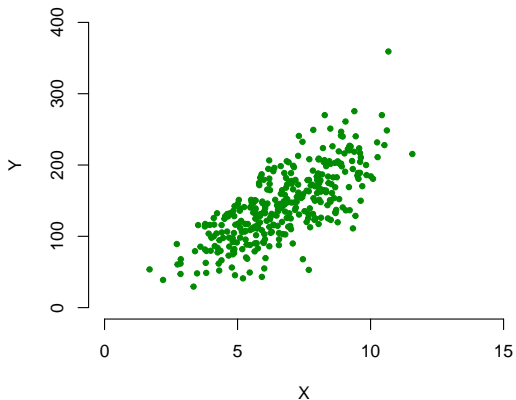
Example 4.1: IRLS

- Suppose $X \sim \mathcal{N}(\mu = 6.5, \sigma^2 = 4)$.
- Let $\boldsymbol{\mu} = 10 + 20\mathbf{x}$
- Suppose $\mathbf{y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$, where $\epsilon_i \sim \mathcal{N}(\mu = 0, \sigma_i^2 = 8\mu_i)$.
- Put another way, $\mathbf{V}(\boldsymbol{\mu}) = \text{diag}(8\boldsymbol{\mu})$.

Example 4.1: Coding setup of example in R

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3 n <- 350
4 X <- rnorm(n, 6.5, 2)
5 mu <- 10 + 20*X
6 Y <- mu + rnorm(n, 0, sqrt(8*mu))
```

Example 4.1: Scatter plot



ITERATIVELY RE-WEIGHTED LEAST SQUARES

Example 4.1: IRLS in R

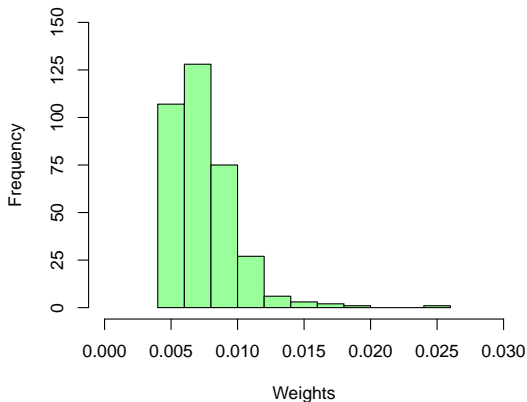
```
1 ## Set threshold and initialize
2 iter <- 0
3 tol <- 1e-16
4 delta <- 1
5
6 zz <- lm(Y ~ X)
7 betaj <- as.numeric(zz$coefficients)
8 zz0 <- zz
9
10 while (delta > tol)
11 {
12   ## Store previous iteration
13   prior.betaj <- betaj
14
15   ## Weights
16   W <- 1/predict(zz)
17
18   ## Run WLS
19   zz <- lm(Y ~ X, weights = W)
20
21   ## Extract coefficient
22   betaj <- as.numeric(zz$coefficients)
23
24   ## Check for convergence and track number of iterations
25   delta <- sum((betaj - prior.betaj)^2)
26   iter <- iter + 1
27 }
```

Example 4.1: Results

```
1 ## Report number of iterations
2 > print(iter)
3 [1] 6
4
5 ## Verify convergence
6 > print(delta)
7 [1] 3.849917e-18
8
9 > summary(zz0)
10
11 Coefficients:
12             Estimate Std. Error t value Pr(>|t|)
13 (Intercept)   4.6414     6.9195   0.671   0.503
14 X             20.7031     0.9971  20.764 <2e-16 ***
15
16 > summary(zz)
17
18             Estimate Std. Error t value Pr(>|t|)
19 (Intercept)   5.878      5.965   0.985   0.325
20 X             20.519     0.924  22.207 <2e-16 ***
```

ITERATIVELY RE-WEIGHTED LEAST SQUARES

Example 4.1: Check to make sure weights are sensible



IRLS: More sophisticated functions in R

- With the `glm()` function in R, you can specify more complex variance functions.
- Example: $\text{Var}[\epsilon_j] = (\theta_1 + |\mu_j|^{\theta_2})^2$, where θ_1 and θ_2 are also unknown parameters.
- θ_1 and θ_2 are also estimated iteratively (details omitted).

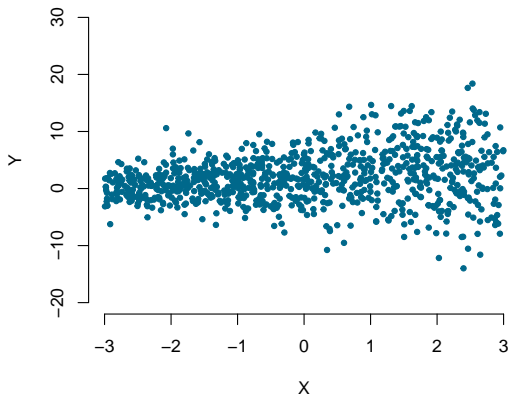
Example 4.2:

- Let $X \sim \text{Uniform}(-3, 3)$.
- Let $Y = 2 + 0.7X + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 2 + |\mu|^{1.1})$.
 - ▶ Mean model: $Y = \beta_0 + \beta_1 x + \epsilon$.
 - ▶ Variance model: $\text{Var}[\epsilon_i] = (\theta_1 + |\mu_i|_2^{\theta_2})^2$, where θ_1 and θ_2 are the unknown parameters of the variance model.
- One possibility: IRLS based on $W(\mathbf{x}) = (\theta_1 + |\hat{Y}(\mathbf{x})|_2^{\theta_2})^{-2}$, where θ_1 , θ_2 , and $\boldsymbol{\beta}$ are estimated iteratively (details omitted).
- Another possibility: Maximum likelihood (although variance estimate is biased).
- Most common: REML (restricted maximum likelihood).

Example 4.2: Coding setup of example in R

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3 n <- 900
4 X <- runif(n, -3, 3)
5 Y <- 2 + 0.7*X + rnorm(n, 0, 2 + abs(2 + 0.7*X)^(1.1))
```

Example 4.2: Scatter plot



Example 4.2: GLS (generalized least squares)

```
1 ## Required library for GLS
2 library(nlme)
3
4 zz1 <- gls(Y~X, weights = varConstPower())
5
6 ## Salient part of R output
7
8 Variance function:
9 Structure: Constant plus power of variance covariate
10 Formula: ~fitted(.)
11 Parameter estimates:
12   const   power
13 2.902585 1.458620
14
15 Coefficients:
16               Value Std.Error   t-value p-value
17 (Intercept)  1.9545581 0.14316666  13.652327    0
18 X            0.6223332 0.07673958   8.109676    0
```

Example 4.2: Compare to results from `lm()`

```
1 Coefficients:
2           Estimate Std. Error t value Pr(>|t|)
3 (Intercept)  1.93307    0.14587  13.252 < 2e-16 ***
4 X            0.57323    0.08525   6.724 3.14e-11 ***
```

Notes from examples:

- WLS (whether fit through IRLS through as in the first example or the `glm()` function as in the second example) produced lower variance, which is not surprising as the variance function was correctly specified.
- If you fit a model using the `glm()` function in R, you'll notice it's using something called *restricted maximum likelihood* (REML) to obtain estimates. You'll talk about this in BIOS 7346, so I'm electing to avoid an in-depth discussion on it in this course.
- Keep in mind that $\text{Cov}[\hat{\boldsymbol{\beta}}]$ can be estimated in a plug-in fashion:

$$\begin{aligned} \widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] &= \hat{\sigma}^2 (\mathbf{X}^T \mathbf{V}^{-1}(\hat{\boldsymbol{\beta}}) \mathbf{X})^{-1} \\ &= \left(\frac{1}{N-K} \sum_{i=1}^N \frac{(y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})^2}{V_i(\hat{\boldsymbol{\beta}})} \right) (\mathbf{X}^T \mathbf{V}^{-1}(\hat{\boldsymbol{\beta}}) \mathbf{X})^{-1} \end{aligned}$$

Things to come:

- We have not yet considered what happens when the weights are “misspecified” (meaning that they are *not* actually inversely proportional to the variance), or how to correct this issue. We will eventually do so!
- Further, these ideas will generalize next semester, in which \mathbf{V} will be block-diagonal (rather than diagonal) to account for clusters.

This unit:

- Weighted least squares.
- Gauss-Markov theorem for weighted least squares.
- Iteratively re-weighted least squares.

SUMMARY: SO FAR

- Random vectors and matrices; multivariate normal theory.
- Ordinary least squares.
- Hypothesis testing and ANOVA.
- Weighted least squares.

SUMMARY: COMING UP

- Misspecification.
- Confidence regions and prediction.
- Diagnostics.
- Regularization.
- Bayesian regression.
- Exponential families.
- Generalized linear models.
- Sandwich and bootstrap.
- Quasi-likelihood.
- Hypothesis testing for GLMs.
- Diagnostics for GLMs.
- Further considerations for binary outcomes.
- Nonlinear least squares.