

Linear Regression Models

P8111

Lecture 19

Jeff Goldsmith
April 5, 2016



THE DEPARTMENT OF
BIostatISTICS



Columbia University
**MAILMAN SCHOOL
OF PUBLIC HEALTH**

Today's Lecture

$$\boxed{\text{Var}(\epsilon) = \sigma^2 \mathbf{I}}$$

OLS

PemLS

- Weighted least squares
- Generalized least squares

Multiple regression model

We typically pose a model of the form

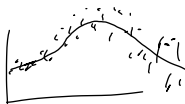
$$\underline{y_i} | \underline{x_i} = \underline{x_i} \beta + \epsilon_i$$

and assume $\underline{Var(\epsilon_i)} = \sigma^2$

- Today we're concerned with $Var(\epsilon_i) = \frac{\sigma^2}{w_i}$ ~~$\propto w_i \propto \frac{1}{Var(\epsilon_i)}$~~
- More generally, we'll look at $Var(\epsilon) = \sigma^2 \mathbf{W}$ or $Var(\epsilon) = \Sigma$
- Contexts include non-constant variance, sampling data (survey weights), proportional data (sample size in groups), meta-analysis (variance of effects in each study)

Handwritten diagram showing a vector of error terms ϵ_i with a σ^2 label and a vector of weights w_i .

Weighted least squares



- One way to handle non-constant variance is a variance stabilizing transformation, which works well if the variance depends on the mean
- Weighted least squares builds the weighting terms directly into the criterion to be minimized

- Let \mathbf{W} be the matrix with $(i, i)^{th}$ entry $\frac{1}{w_i}$ and 0 elsewhere
- Then $Var(\epsilon) = \sigma^2 \mathbf{W}$

Weighted least squares

$$\frac{1}{w_i}$$

- For weighted least squares, we minimize the RSS with terms weighted according to their variance

$$\begin{aligned}\underline{RSS_W(\beta)} &= \sum w_i (y_i - x_i^T \beta)^2 \quad \text{wt'd mean!} \quad \checkmark \\ &= (y - X\beta)^T \underline{W}^{-1} (y - X\beta)\end{aligned}$$

- We weight more heavily terms with low variance (small $\frac{\sigma^2}{w_i}$) and less heavily terms with high variance (big $\frac{\sigma^2}{w_i}$)
- Basic plan – differentiate $\underline{RSS_W(\beta)}$ wrt β and find the minimum

Weighted least squares estimator

$$RSS_W(\beta) = (y - X\beta)^T W^{-1} (y - X\beta)$$

$$= (y - X\beta)^T (W^{-1} y - W^{-1} X\beta)$$

$$y^T W^{-1} y - \beta^T X^T W^{-1} y - y^T W^{-1} X \beta + \beta^T X^T W^{-1} X \beta$$

$$= -2 \beta^T X^T W^{-1} y + \beta^T X^T W^{-1} X \beta$$

$$\frac{\partial}{\partial \beta} = -2 X^T W^{-1} y + 2 X^T W^{-1} X \beta = 0$$

$$X^T W^{-1} X \beta = X^T W^{-1} y \Rightarrow \beta_{WLS} = (X^T W^{-1} X)^{-1} X^T W^{-1} y$$

A note about MLE

We have the model

$$\underline{y = X\beta + \epsilon}$$

where $E(\epsilon) = 0$ and $Var(\epsilon) = \sigma^2 W$.

- Additionally, assume $\epsilon \sim N(0, \sigma^2 W)$
- Put differently, we're imposing the model

$$\underline{y \sim N(X\beta, \sigma^2 W)}$$

- y is multivariate Normal

Maximum likelihood estimation

Using matrix notation:

$$L(\beta; y) \propto \exp \left\{ \underbrace{-\frac{1}{2\sigma^2} (\underline{y} - \underline{X}\beta)^T \underline{W}^{-1} (\underline{y} - \underline{X}\beta)}_{RSS_w(\beta)} \right\}$$

Pre-whitening data

$$\begin{bmatrix} \frac{1}{w_1} & & \\ & \ddots & \\ & & \frac{1}{w_n} \end{bmatrix}$$

- Let $\underline{W}^{1/2}$ be the diagonal matrix with $(i, i)^{th} \frac{1}{\sqrt{w_i}}$ and 0 elsewhere
- So $\underline{W}^{-1/2} \stackrel{def}{=} (\underline{W}^{1/2})^{-1}$ is a diagonal matrix with $\sqrt{w_i}$ on the main diagonal and 0 elsewhere
- Note $\underline{W} = \underline{W}^{1/2}(\underline{W}^{1/2})^T$ and $\underline{W}^{1/2}\underline{W}^{-1/2} = I$
- So $Var(\underline{W}^{-1/2}\epsilon) =$

$$\begin{aligned} & (\underline{W}^{-1/2})^T Var(\epsilon) \underline{W}^{-1/2} \\ &= \sigma^2 \underbrace{(\underline{W}^{-1/2})^T (\underline{W}^{1/2} \underline{W}^{1/2})}_{I} \underline{W}^{-1/2} \\ &= \sigma^2 I \\ &= \sigma^2 I \end{aligned}$$

Pre-whitening data

- Let's pre-multiply everything by $W^{-1/2}$:

$$\left\{ \begin{array}{l} \triangleright z = W^{-1/2} y \\ \triangleright M = W^{-1/2} X \\ \triangleright \delta = W^{-1/2} \epsilon \end{array} \right.$$

$$y = (x\beta + \epsilon) \quad \checkmark$$

$\epsilon \sim (0, \sigma^2 \omega)$

- Our model is now

$$z = M\beta + \delta$$

- The OLS estimate of β is

$$(M^T M)^{-1} M^T z$$

$$\delta \sim (0, \sigma^2 I)$$

$$\begin{aligned} & (x^T \omega^{-1/2} \omega^{-1/2} x) x^T \omega^{-1/2} \omega^{-1/2} y \\ & (x^T \omega^{-1} x)^{-1} x^T \omega^{-1} y \end{aligned}$$

WLS example

- Data from a physics experiment, available as `physics` from the library `alr3`
- y : scattering cross-section, s : square of total energy,
 $x = \underline{s^{-1/2}}$
- Theoretical model:
 $E(y|s) = \beta_0 + \beta_1 s^{-1/2} + \text{relatively small terms}$
- Regression model: $\underline{y = \beta_0 + \beta_1 x + \epsilon}$
- $\underline{SD = \sqrt{Var(y|x)}}$ are known from the experiment

WLS example

↓ ↓
> library(alr3)
> data(physics)
> physics

↓
x y SD
1 0.345 367 17
2 0.287 311 9
3 0.251 295 9
4 0.225 268 7
5 0.207 253 7
6 0.186 239 6
7 0.161 220 6
8 0.132 213 6
9 0.084 193 5
10 0.060 192 5

↙

$$SD(\epsilon_i) =$$

$$\text{Var}(\epsilon_i) = \frac{\sigma^2}{w_i}$$

$$\Rightarrow w_i \propto \frac{1}{SD(\epsilon_i)^2}$$

WLS example

```
> lm.physics.wls <- lm(y~x, weights=1/SD^2, data=physics)
> summary(lm.physics.wls)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	148.473	8.079	18.38	7.91e-08 ***
x	<u>230.835</u>	<u>47.550</u>	11.16	3.71e-06 ***

Residual standard error: 1.657 on 8 degrees of freedom
Multiple R-squared: 0.9397, Adjusted R-squared: 0.9321
F-statistic: 124.6 on 1 and 8 DF, p-value: 3.710e-06

$$530 \pm 90$$

$$(440, 620)$$

WLS example

↓
> lm.physics.ols <- lm(y~x, data=physics)
> summary(lm.physics.ols)

Coefficients:

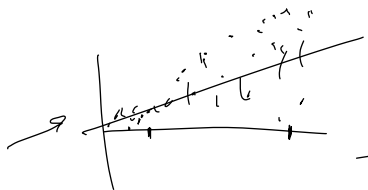
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	135.00	10.08	13.4	9.21e-07 ***
<u>x</u>	619.71	47.68	13.0	1.16e-06 ***

Residual standard error: 12.69 on 8 degrees of freedom

Multiple R-squared: 0.9548, Adjusted R-squared: 0.9491

F-statistic: 168.9 on 1 and 8 DF, p-value: 1.165e-06

WLS in practice



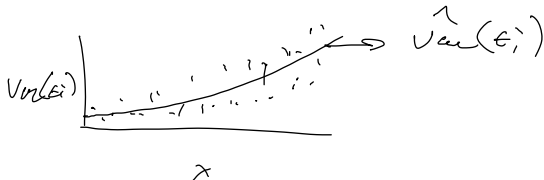
- Real life is rarely nice enough to give you the right weight
- Try to obtain an estimate of $var(\epsilon_i)$, plug that into W ...

① OLS

$$var(\epsilon_i) \approx \frac{\sum \hat{\epsilon}_i^2}{n}$$

② $\hat{\epsilon}_i$

$$\frac{\sum \hat{\epsilon}_i^2}{n}$$



Generalized least squares

$$\sigma^2 \mathbb{I} \Rightarrow \sigma^2 \underbrace{\mathbb{W}}_{\text{diag}} \Rightarrow \underline{\sigma^2 \Sigma}$$

- Weighted least squares can help a lot, but what if errors are correlated?
- That is, suppose our model is

$$\underline{y = X\beta + \epsilon}$$

where $\underline{E(\epsilon) = 0}$ and $\underline{Var(\epsilon) = \sigma^2 \Sigma}$

- (By analogy with WLS, suppose Σ is known but σ^2 is not; in general, one usually writes $\underline{Var(\epsilon) = \Sigma}$)
- Note, in terms of generality, GLS > WLS > OLS

Generalized least squares

- Writing out $RSS_G(\beta)$ as a sum is hard; possible using vector notation.
- Possibilities:
 - ▶ MLE (equivalent to minimizing RSS)
 - ▶ Pre-whiten

MLE

We have the model

$$\underline{y = X\beta + \epsilon}$$

where $\underline{E(\epsilon) = 0}$ and $\underline{Var(\epsilon) = \sigma^2 \Sigma}$.

- Additionally, assume $\underline{\epsilon \sim N(0, \underline{\sigma^2 \Sigma})}$
- Put differently, we're imposing the model

$$y \sim N(X\beta, \underline{\sigma^2 \Sigma})$$

- y is multivariate Normal

MLE

Using matrix notation:

$$L(\beta; y) \propto \exp \left\{ \underbrace{-\frac{1}{2\sigma^2} (y - X\beta)^T Z^{-1} (y - X\beta)} \right\}$$

$$\frac{\partial}{\partial \beta} = \dots = 0$$

$$\Rightarrow \hat{\beta}_{OLS} = (X^T Z^{-1} X)^{-1} X^T Z^{-1} y$$

Pre-whitening data

- Let $\Sigma = SS^T$ be the *Cholesky decomposition* of Σ
- Let's pre-multiply everything by S^{-1} :

$$\begin{aligned} \triangleright z &= \cancel{S^{-1}S} y & S^{-1} y \\ \triangleright M &= \cancel{S^{-1}S} X & S^{-1} X \\ \triangleright \delta &= \cancel{S^{-1}S} \epsilon & S^{-1} \epsilon \end{aligned}$$

- Our model is now

$$z = M\beta + \delta$$

- The OLS estimate of β is

$$\underline{(M^T M)^{-1} M^T z}$$

In practice

"Feasible" GLS

$$OLS \Rightarrow \hat{\Sigma}^{-1}$$

$$\Rightarrow GLS_{\hat{\Sigma}}$$

Some useful notes on GLS

Using $\hat{\beta}_{GLS} = (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{y}$, it turns out that

- $E(\hat{\beta}_{GLS}) = \beta$

$$E\left(\underbrace{(\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1}}_{\text{matrix}} \underbrace{\mathbf{X}^T \Sigma^{-1}}_{\text{matrix}} \mathbf{y}\right) =$$

$$\underbrace{(\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1}}_{\text{matrix}} \underbrace{\mathbf{X}^T \Sigma^{-1}}_{\text{matrix}} \mathbf{X} \beta = \beta$$

- $Var(\hat{\beta}_{GLS}) = \sigma^2 (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1}$

Diagram illustrating the relationship between GLS and OLS:

On the left, the GLS variance formula is shown with annotations:

- $Var(y)$ is indicated by a line pointing to Σ in $\sigma^2 \Sigma$.
- $\sigma^2 \Sigma$ is crossed out with a large 'X'.

On the right, the OLS formulas are shown:

- $E(\hat{\beta}_{OLS}) = \beta$
- $Var(\hat{\beta}_{OLS}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

A large 'X' is drawn over the OLS formulas, indicating they are not applicable in the GLS context.

Some less useful notes on GLS

- Typically we don't really know Σ and have to estimate it too

$$OLS \Rightarrow \hat{\beta} \Rightarrow GLS_{\hat{\Sigma}}$$

- A common approach is to parameterize Σ using a small number of parameters
- Comes up a lot for longitudinal and multilevel data

$$Var(\varepsilon) = \underline{\underline{\sigma^2 \Sigma}}$$

$$Var(\varepsilon) = \Sigma(\sigma^2, \rho, k)$$

Today's big ideas

- Weighted and generalized least squares
-

- Suggested reading: Ch. 5