

# Generalized Least Squares

FW8051 Statistics for Ecologists

Department of Fisheries, Wildlife and Conservation Biology



# Learning Objectives

Learn how to use generalized least squares (GLS) to model data where  $Y_i|X_i$  is normally distributed, but the variance of the residuals is not constant and may depend on one or more predictor variables.

# Generalized Least Squares

Can be used to model data where  $Y_i|X_i$  is normally distributed,  
but we have:

- Non-constant variance (Chapter 5)
- Data that are correlated
  - Multiple measurements on the same sample unit (Chapter 18)
  - Temporal dependence (Chapter 6 of Zuur et al)
  - Spatial dependence (Chapter 7 of Zuur et al)

# Generalized Least Squares

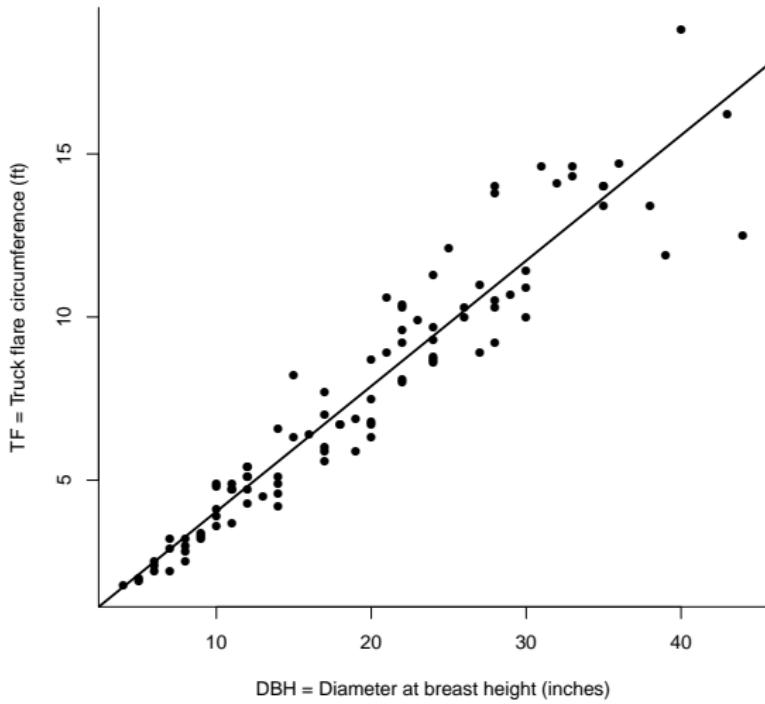
Can be used to model data where  $Y_i|X_i$  is normally distributed,  
but we have:

- Non-constant variance (Chapter 5)
- Data that are correlated
  - Multiple measurements on the same sample unit (Chapter 18)
  - Temporal dependence (Chapter 6 of Zuur et al)
  - Spatial dependence (Chapter 7 of Zuur et al)

For this class, we will focus on non-constant variance and  
multiple measurements on the same sample unit [later in the  
course]

# Trunk Flare Diameter





# Linear Model

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

Assume  $\epsilon_i$  are independent, normally distributed, with constant variance.

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

*iid* = independent and identically distributed

# Generalized Least Squares: Non-Constant Variance

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

$$\sigma_i^2 \sim f(X_i; \tau)$$

Model the mean and variance:

- $E[Y_i|X_i] = \beta_0 + X_i\beta_1$
- $Var[Y_i|X_i] = f(X_i; \tau)$ , where  $\tau$  are additional variance parameters.

# Generalized Least Squares: Non-Constant Variance

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

$$\sigma_i^2 \sim f(X_i; \tau)$$

Some options:

- $\sigma_i^2 = \sigma_g^2$  (different  $\sigma$  for each group,  $g$ , modeled using a set of multiplicative factors and a reference group)

# Generalized Least Squares: Non-Constant Variance

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

$$\sigma_i^2 \sim f(X_i; \tau)$$

Some options:

- $\sigma_i^2 = \sigma_g^2$  (different  $\sigma$  for each group,  $g$ , modeled using a set of multiplicative factors and a reference group)
- $\sigma_i^2 = \sigma^2 X_i$ , or  $= \sigma^2 |X_i|^{2\delta}$ , or  $= \sigma^2 e^{2\delta X_i}$  for continuous covariate,  $X$

# Generalized Least Squares: Non-Constant Variance

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

$$\sigma_i^2 \sim f(X_i; \tau)$$

Some options:

- $\sigma_i^2 = \sigma_g^2$  (different  $\sigma$  for each group,  $g$ , modeled using a set of multiplicative factors and a reference group)
- $\sigma_i^2 = \sigma^2 X_i$ , or  $= \sigma^2 |X_i|^{2\delta}$ , or  $= \sigma^2 e^{2\delta X_i}$  for continuous covariate,  $X$
- $\sigma_i^2 = \sigma^2 E[Y_i|X_i]^{2\theta} = \sigma^2 (\beta_0 + X_i\beta_1)^{2\theta}$ . This one is not in Zuur et al. (2009) and can be fit using:  
`varPower (form=~fitted(.))`

# Generalized Least Squares: Non-Constant Variance

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

$$\sigma_i^2 \sim f(X_i; \tau)$$

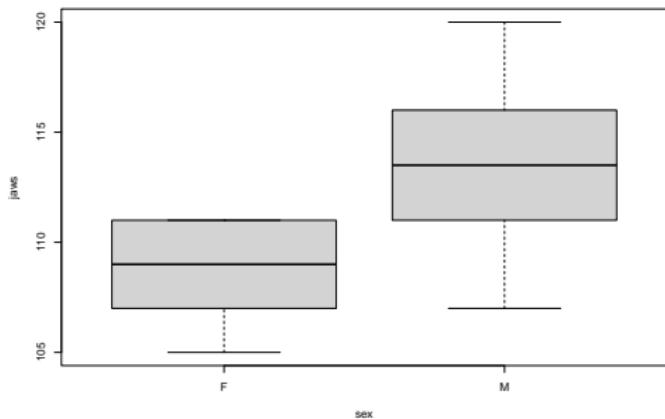
Some options:

- $\sigma_i^2 = \sigma_g^2$  (different  $\sigma$  for each group,  $g$ , modeled using a set of multiplicative factors and a reference group)
- $\sigma_i^2 = \sigma^2 X_i$ , or  $= \sigma^2 |X_i|^{2\delta}$ , or  $= \sigma^2 e^{2\delta X_i}$  for continuous covariate,  $X$
- $\sigma_i^2 = \sigma^2 E[Y_i|X_i]^{2\theta} = \sigma^2 (\beta_0 + X_i\beta_1)^{2\theta}$ . This one is not in Zuur et al. (2009) and can be fit using:  
`varPower (form=~fitted(.))`
- Some combination of the above + other options (see Ch. 4 of Zuur et al.)

# T-test with unequal variances: Jaw data

```
males<-c(120, 107, 110, 116, 114, 111, 113, 117, 114, 112)
females<-c(110, 111, 107, 108, 110, 105, 107, 106, 111, 111)
jawdat <- data.frame(jaws = c(males, females),
                      sex = c(rep("M", 10), rep("F", 10)))
```

```
boxplot(jaws~sex, data=jawdat)
```



## T-test with unequal variances: Jaw data

$Y_i$  = jaw length for jackal  $i$

$$Y_i \sim N(\mu_i, \sigma_i^2) \quad (1)$$

$$\mu_i = \beta_0 + \beta_1 I(\text{sex}=\text{male})_i \quad (2)$$

$$\log(\sigma_i) = \log(\sigma) + \log(\delta)I(\text{sex} = \text{female})_i \quad (3)$$

```
gls_ttest <- gls(jaws ~ sex,
                   weights = varIdent(form = ~ 1 | sex),
                   data = jawdat)
```

Estimates of regression parameters are obtained by minimizing:

$$\sum_{i=1}^n \frac{(Y - \mu_i)^2}{2\sigma_i^2}$$

## Summary output

Variance model:

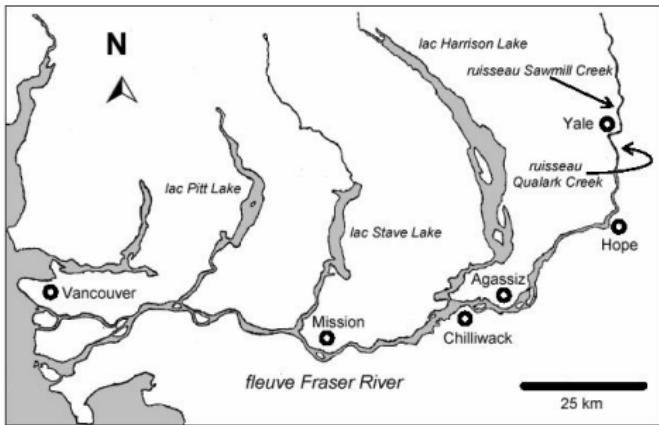
$$\log(\sigma_i) = \log(\sigma) + \log(\delta)I(\text{sex} = \text{female})_i$$

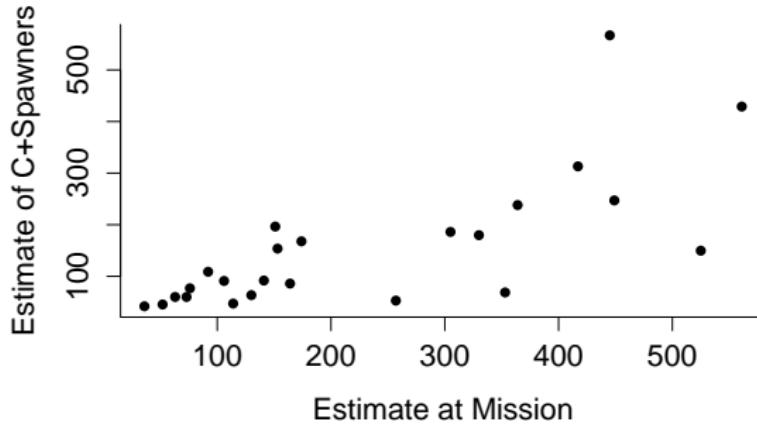
R returns:

- $\hat{\sigma}$
- $\hat{\delta}$  (for females) and  $\delta = 1$  for males

See textbook example via this link.

# Fraser River Sockeye





Use historical correlation between the *count at Mission* and  $S_t + C_t$  to manage the fishery

## Variance increasing with $X_i$ or $\mu_i$

1. Fixed variance model:  $\sigma_i^2 = \sigma^2 \text{MisEsc}_i$
2. Power variance model:  $\sigma_i^2 = \sigma^2 |\text{MisEsc}_i|^{2\delta}$
3. Exponential variance model:  $\sigma_i^2 = \sigma^2 e^{2\delta \text{MisEsc}_i}$
4. Constant + power variance model:  $\sigma_i^2 = \sigma^2 (\delta_1 + |\text{MisEsc}_i|^{\delta_2})^2$

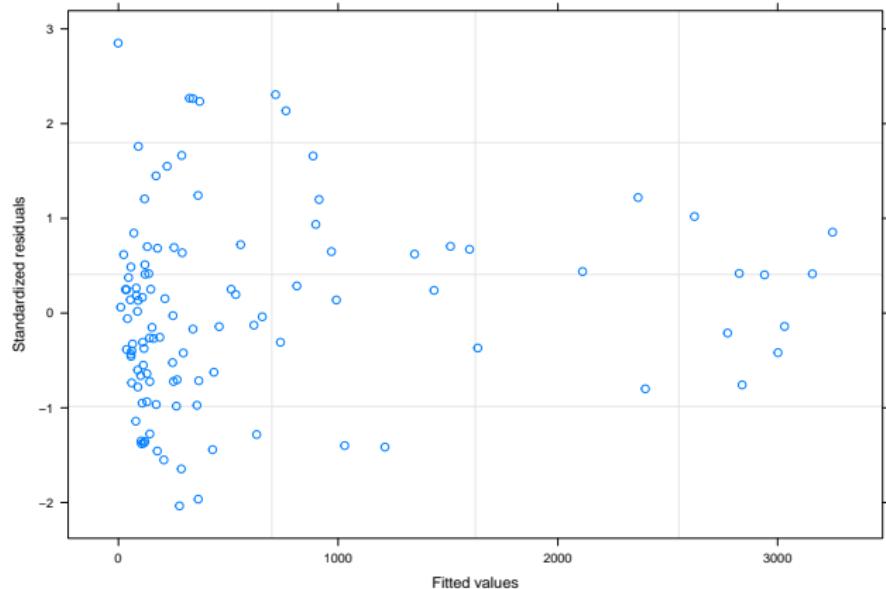
```
varconstp <- gls(SpnEsc ~ MisEsc,  
                    weights = varConstPower(form = ~ MisEsc),  
                    data = sockeye)
```

See textbook via this link.

# Standardized residuals

**Standardized residuals** =  $(Y_i - \hat{Y}_i)/\hat{\sigma}_i$ , should have approximately constant variance:

```
plot(varconstp)
```



## Variance depending on $\mu_i$

$$Y_i \sim N(\mu_i, \sigma_i^2)$$

$$\mu_i = \beta_0 + \beta_1 X_i$$

$$\sigma_i^2 = \mu_i^{2\delta}$$

Fit using `varPower(form = ~ fitted(.))`. See textbook via this link.

## Zuur et al.'s Strategy (Section 4.2.3)

- Start with a “full model” (containing as many predictors as possible) fit using `lm`

## Zuur et al.'s Strategy (Section 4.2.3)

- Start with a “full model” (containing as many predictors as possible) fit using `lm`
- Inspect residuals, consider alternative variance structures if necessary
  - Inspect normalized residuals =  $\frac{Y_i - \hat{Y}_i}{\sigma_i}$  (these should have constant variance if the variance model is adequate)
  - Compare models using AIC (after refitting the constant variance model using `gls`)
  - Settle on optimal variance model

## Zuur et al.'s Strategy (Section 4.2.3)

- Start with a “full model” (containing as many predictors as possible) fit using `lme`
- Inspect residuals, consider alternative variance structures if necessary
  - Inspect normalized residuals =  $\frac{Y_i - \hat{Y}_i}{\sigma_i}$  (these should have constant variance if the variance model is adequate)
  - Compare models using AIC (after refitting the constant variance model using `gls`)
  - Settle on optimal variance model
- Choose best set of predictor variables for the mean of  $Y_i$  (using the variance model, chosen above)

## Zuur et al.'s Strategy (Section 4.2.3)

- Start with a “full model” (containing as many predictors as possible) fit using `lm`
- Inspect residuals, consider alternative variance structures if necessary
  - Inspect normalized residuals =  $\frac{Y_i - \hat{Y}_i}{\sigma_i}$  (these should have constant variance if the variance model is adequate)
  - Compare models using AIC (after refitting the constant variance model using `gls`)
  - Settle on optimal variance model
- Choose best set of predictor variables for the mean of  $Y_i$  (using the variance model, chosen above)
- Check diagnostics again and pick best model

... Try to make sense of your results.

# Approximate Confidence and Prediction Intervals

For a confidence interval, we need to consider:

$$\text{var}(\widehat{E[Y|X]}) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1)$$

# Approximate Confidence and Prediction Intervals

For a confidence interval, we need to consider:

$$\text{var}(\widehat{E[Y|X]}) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1)$$

For a prediction interval, we need to consider:

$$\text{var}(\hat{Y}_i|X_i) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1 + \epsilon_i)$$

# Approximate Confidence and Prediction Intervals

For a confidence interval, we need to consider:

$$\text{var}(\widehat{E[Y|X]}) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1)$$

For a prediction interval, we need to consider:

$$\text{var}(\hat{Y}_i|X_i) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1 + \epsilon_i)$$

- Confidence interval = captures uncertainty regarding the average value of  $Y$  (given by the line)

# Approximate Confidence and Prediction Intervals

For a confidence interval, we need to consider:

$$\text{var}(\widehat{E[Y|X]}) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1)$$

For a prediction interval, we need to consider:

$$\text{var}(\hat{Y}_i|X_i) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1 + \epsilon_i)$$

- Confidence interval = captures uncertainty regarding the average value of  $Y$  (given by the line)
- Prediction interval = captures uncertainty regarding a *particular* value of  $Y$  (need to also consider spread about the line)

# Statistical Theory

- means, sums, regression coefficients will be normally distributed if our sample size,  $n$ , is “large” (thanks Central Limit Theorem!)

Consider  $(X, Y) \sim MVN(\mu, \Sigma)$ .

- $\mu = (E(X), E(Y))$
- $\Sigma = \begin{bmatrix} \sigma_x^2 & \sigma_{x,y} \\ \sigma_{x,y} & \sigma_y^2 \end{bmatrix}$

Then for constants  $a$  and  $b$ :

- $aX + bY$  will be  $MVN(\tilde{\mu}, \tilde{\sigma})$
- $\tilde{\mu} = E(aX + bY) = aE(X) + bE(Y)$
- $\tilde{\sigma} = var(aX + bY) = a^2 var(X) + b^2 var(Y) + 2abcov(X, Y)$

# Matrix Multiplication

Matrix multiplication offers a convenient way to calculate variances of linear combinations of random variables (i.e.,  $aX + bY$ )

Let  $\Sigma = \begin{bmatrix} \sigma_X^2 & \sigma_{X,Y}^2 \\ \sigma_{X,Y}^2 & \sigma_Y^2 \end{bmatrix}$  give the variance covariance matrix of  $(X, Y)$

- $\sigma_X^2, \sigma_Y^2$  = variance of  $X, Y$
- $\sigma_{X,Y}^2$  = covariance of  $X$  and  $Y$

# Matrix Multiplication

Matrix multiplication offers a convenient way to calculate variances of linear combinations of random variables (i.e.,  $aX + bY$ )

Let  $\Sigma = \begin{bmatrix} \sigma_X^2 & \sigma_{X,Y}^2 \\ \sigma_{X,Y}^2 & \sigma_Y^2 \end{bmatrix}$  give the variance covariance matrix of  $(X, Y)$

- $\sigma_X^2, \sigma_Y^2$  = variance of  $X, Y$
- $\sigma_{X,Y}^2$  = covariance of  $X$  and  $Y$

$$\begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} = aX + bY$$

# Matrix Multiplication

Matrix multiplication offers a convenient way to calculate variances of linear combinations of random variables (i.e.,  $aX + bY$ )

Let  $\Sigma = \begin{bmatrix} \sigma_X^2 & \sigma_{X,Y}^2 \\ \sigma_{X,Y}^2 & \sigma_Y^2 \end{bmatrix}$  give the variance covariance matrix of  $(X, Y)$

- $\sigma_X^2, \sigma_Y^2$  = variance of  $X, Y$
- $\sigma_{X,Y}^2$  = covariance of  $X$  and  $Y$

$$\begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} = aX + bY$$

$$\begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} \sigma_X^2 & \sigma_{X,Y}^2 \\ \sigma_{X,Y}^2 & \sigma_Y^2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\sigma_{X,Y}^2$$

# Confidence Intervals

$$var(\widehat{E[Y|X]}) = Var(\hat{\beta}_0 + X\hat{\beta}_1)$$

# Confidence Intervals

$$var(\widehat{E[Y|X]}) = Var(\hat{\beta}_0 + X\hat{\beta}_1)$$

$$\begin{bmatrix} a & b \end{bmatrix} = \begin{bmatrix} 1 & X \end{bmatrix} \text{ and } \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix}$$

# Confidence Intervals

$$var(\widehat{E[Y|X]}) = Var(\hat{\beta}_0 + X\hat{\beta}_1)$$

$$\begin{bmatrix} a & b \end{bmatrix} = \begin{bmatrix} 1 & X \end{bmatrix} \text{ and } \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix}$$

Define:

$$\hat{\Sigma} = \begin{bmatrix} \sigma_{\hat{\beta}_0}^2 & \sigma_{\hat{\beta}_0, \hat{\beta}_1}^2 \\ \sigma_{\hat{\beta}_0, \hat{\beta}_1}^2 & \sigma_{\hat{\beta}_1}^2 \end{bmatrix}$$

# Confidence Intervals

$$\widehat{var}(E[\widehat{Y|X}]) = Var(\hat{\beta}_0 + X\hat{\beta}_1)$$

$$\begin{bmatrix} a & b \end{bmatrix} = \begin{bmatrix} 1 & X \end{bmatrix} \text{ and } \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix}$$

Define:

$$\hat{\Sigma} = \begin{bmatrix} \sigma_{\hat{\beta}_0}^2 & \sigma_{\hat{\beta}_0, \hat{\beta}_1}^2 \\ \sigma_{\hat{\beta}_0, \hat{\beta}_1}^2 & \sigma_{\hat{\beta}_1}^2 \end{bmatrix}$$

$$\widehat{var}(E[\widehat{Y|X}]) = Var(\hat{\beta}_0 + X\hat{\beta}_1) = \begin{bmatrix} 1 & X \end{bmatrix} \hat{\Sigma} \begin{bmatrix} 1 \\ X \end{bmatrix}$$

# Calculations in R

We can use matrix multiplication to efficiently calculate intervals for multiple observations.

In R, we do this by using `%*%`

Let  $X$  = design matrix = 
$$\begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$$

$\widehat{E[Y|X]} = X \%*\% \text{coef}(\text{modelname})$  gives us predicted values for all rows of the  $X$  matrix.

# Calculations in R

We can use matrix multiplication to efficiently calculate intervals for multiple observations.

In R, we do this by using `%*%`

Let  $X$  = design matrix = 
$$\begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$$

$\widehat{E[Y|X]} = X \%*\% \text{coef}(\text{modelname})$  gives us predicted values for all rows of the  $X$  matrix.

$$\widehat{E[Y|X]} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 + X_1 \hat{\beta}_1 \\ \hat{\beta}_0 + X_2 \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_0 + X_n \hat{\beta}_1 \end{bmatrix}$$

# Calculations in R

$$\widehat{var}(E[\widehat{Y}|X]) = X \%*\% \text{vcov}(\text{modelname}) \%*\% t(X)$$

- `vcov(model)` =  $\hat{\Sigma}$  = estimated variance-covariance matrix of  $\hat{\beta}$  (works for `lm`, `gls`, maybe others)
- We use `t` to get a transpose of a matrix

# Calculations in R

$$\widehat{var}(E[\widehat{Y}|X]) = X \%*\% \text{vcov}(\text{modelname}) \%*\% t(X)$$

- `vcov(model)` =  $\hat{\Sigma}$  = estimated variance-covariance matrix of  $\hat{\beta}$  (works for `lm`, `gls`, maybe others)
- We use `t` to get a transpose of a matrix

We end up with a matrix that looks something like:

$$\begin{bmatrix} var(\hat{Y}_1) & cov(\hat{Y}_1, \hat{Y}_2) & \cdots & cov(\hat{Y}_1, \hat{Y}_n) \\ cov(\hat{Y}_2, \hat{Y}_1) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & cov(\hat{Y}_{n-1}, \hat{Y}_n) \\ cov(\hat{Y}_n, \hat{Y}_1) & \cdots & cov(\hat{Y}_n, \hat{Y}_{n-1}) & var(\hat{Y}_n) \end{bmatrix}$$

# Calculations in R

$$\widehat{var}(E[\widehat{Y}|X]) = X \%*\% \text{vcov}(\text{modelname}) \%*\% t(X)$$

- `vcov(model)` =  $\hat{\Sigma}$  = estimated variance-covariance matrix of  $\hat{\beta}$  (works for `lm`, `gls`, maybe others)
- We use `t` to get a transpose of a matrix

We end up with a matrix that looks something like:

$$\begin{bmatrix} var(\hat{Y}_1) & cov(\hat{Y}_1, \hat{Y}_2) & \cdots & cov(\hat{Y}_1, \hat{Y}_n) \\ cov(\hat{Y}_2, \hat{Y}_1) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & cov(\hat{Y}_{n-1}, \hat{Y}_n) \\ cov(\hat{Y}_n, \hat{Y}_1) & \cdots & cov(\hat{Y}_n, \hat{Y}_{n-1}) & var(\hat{Y}_n) \end{bmatrix}$$

Pull off the diagonal elements (the variances) - see the textbook for code.

# Prediction Intervals

$$var(\hat{Y}_i|X_i) = var(\hat{\beta}_0 + X\hat{\beta}_1 + \epsilon_i)$$

- $var(\epsilon_i) = var(Y_i|X_i) = \sigma_i^2$  and estimated by  $\hat{\sigma}_i^2$ .
- In many cases,  $\hat{\sigma}_i^2$  is independent of  $\begin{bmatrix} \hat{\beta}_0 & \hat{\beta}_1 \end{bmatrix}$
- This implies  $cov(\hat{\sigma}_i^2, \hat{\beta}_0) = cov(\hat{\sigma}_i^2, \hat{\beta}_1) = 0$

# Prediction Intervals

$$\text{var}(\hat{Y}_i|X_i) = \text{var}(\hat{\beta}_0 + X\hat{\beta}_1 + \epsilon_i)$$

- $\text{var}(\epsilon_i) = \text{var}(Y_i|X_i) = \sigma_i^2$  and estimated by  $\hat{\sigma}_i^2$ .
- In many cases,  $\hat{\sigma}_i^2$  is independent of  $\begin{bmatrix} \hat{\beta}_0 & \hat{\beta}_1 \end{bmatrix}$
- This implies  $\text{cov}(\hat{\sigma}_i^2, \hat{\beta}_0) = \text{cov}(\hat{\sigma}_i^2, \hat{\beta}_1) = 0$

So, to construct a prediction interval, we approximate  $\text{var}(\hat{Y}_i|X_i)$  with:

$$\text{var}(\hat{Y}_i|X_i) \approx \widehat{\text{var}}(\hat{\beta}_0 + X\hat{\beta}_1) + \hat{\sigma}_i^2 = X\hat{\Sigma}X' + \hat{\sigma}_i^2.$$

## Additional Notes

Note, these estimates are approximate in that:

- They rely on asymptotic normality (central limit theorem)  
[think difference between  $t$  and  $z$ ]
- They ignore uncertainty in the variance parameters

# Temporal or Spatial Correlation

$$\begin{aligned}y_i &= \beta_0 + \beta_1 x_i + \epsilon_i \\ \epsilon_i &\sim N(0, \Omega)\end{aligned}$$

- Time series:  $\text{cor}(\epsilon_i, \epsilon_j) = \rho^{|t_i - t_j|}$

# Temporal or Spatial Correlation

$$\begin{aligned}y_i &= \beta_0 + \beta_1 x_i + \epsilon_i \\ \epsilon_i &\sim N(0, \Omega)\end{aligned}$$

- Time series:  $\text{cor}(\epsilon_i, \epsilon_j) = \rho^{|t_i - t_j|}$
- Spatial data:  $\text{cor}(\epsilon_i, \epsilon_j)$  depends on distance between points.

# Temporal or Spatial Correlation

$$\begin{aligned}y_i &= \beta_0 + \beta_1 x_i + \epsilon_i \\ \epsilon_i &\sim N(0, \Omega)\end{aligned}$$

- Time series:  $\text{cor}(\epsilon_i, \epsilon_j) = \rho^{|t_i - t_j|}$
- Spatial data:  $\text{cor}(\epsilon_i, \epsilon_j)$  depends on distance between points.

If these interest you, I highly recommend taking Brian Aukema's class.