

Regression Assumptions

- linear : $E(Y|X) = X\beta$
 - Y - continuous random variable
 - e.g. $Y=0, 1$ range is unlimited
 - $Y = \text{pos. integers}$ never be integers
 - integer $\Rightarrow Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$
 - Variance is constant $\epsilon_i \sim N(0, \sigma^2)$ constant
 - No outlier & influential points
- polynomial reg ch6
• transformation ch7
- logistic reg
• Poisson reg
~~ch12~~
STAT 4006
- Weighted Least Square ch5
- ~~259~~

5.1 Weighted Least Square (WLS)

- Model

$$E(Y | X = x_i) = \beta' x_i$$

$$\text{Var}(Y | X = x_i) = \frac{\sigma^2}{w_i}$$

assumed known

- Alternative representation

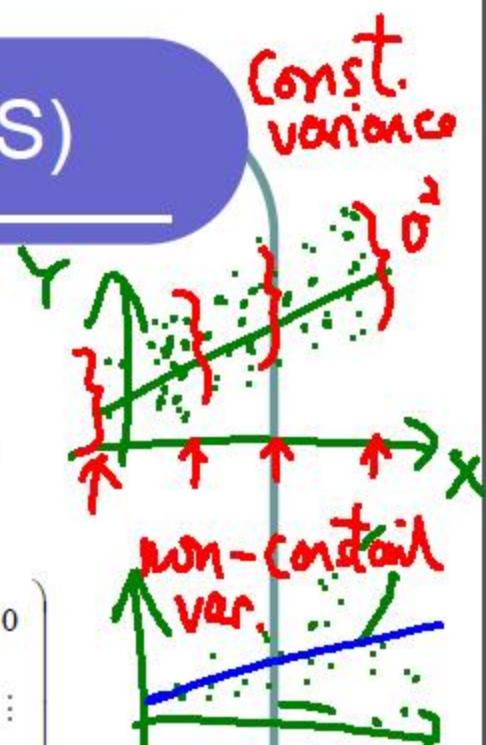
$$e_i \sim \text{id} \quad \text{var}(e_i) = \frac{\sigma^2}{w_i}$$

$$Y = X\beta + e,$$

↖
AXI

$$\text{Var}(e) = \sigma^2 W^{-1} = \sigma^2 \begin{pmatrix} \frac{1}{w_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{w_2} & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \cdots & 0 & \frac{1}{w_n} \end{pmatrix}$$

- Errors are independent but **not** identically distributed





5.1 Weighted Least Square (WLS)

- Model

$$E(Y | X = x_i) = \beta' x_i$$

$$\text{Var}(Y | X = x_i) = \frac{\sigma^2}{w_i}$$

assumed known

- Alternative representation

$$Y = X\beta + e,$$

$$\text{Var}(e) = \sigma^2 W^{-1} = \sigma^2 \begin{pmatrix} \frac{1}{w_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{w_2} & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \cdots & 0 & \frac{1}{w_n} \end{pmatrix}$$

- Errors are independent but **not** identically distributed

$$W = \begin{pmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w_n \end{pmatrix}$$

assume iid
that $y_i \sim \sigma^2$

	mean GPA G	mean IQ Q	class size
STAT 3003	3.3	130	40
STAT 3004	3.4	140	50
STAT 3005	3.5	150	60
STAT 3007	3.7	170	30
STAT 3008	3.999	199	112

$$G_1 = \frac{\sum_{i=1}^{40} y_i}{40} \quad \text{Var}(G_1) = \frac{\sigma^2}{40}$$

$$\text{Var}(Y|x) = \frac{\sigma^2}{W_i} \Rightarrow W_1 = 40$$

$$G_5 = \frac{\sum_{i=1}^{120} z_i}{120} \quad \text{Var}(G) = \frac{\sigma^2}{120}$$

$$\Rightarrow W_5 = 120$$

	$Y = \text{monthly rainfall}$	temp
Jan	$Z_1 + \dots + Z_{31}$.
Feb	$Z_1 + \dots + Z_{28}$.
March	$Z_1 + \dots + Z_{31}$.
April	.	.
	.	.

$Z_i = \text{daily rainfall}$
 $\sim \mathcal{N}(\mu, \sigma^2)$

$$Y_1 = Z_1 + \dots + Z_{31}$$

$$\text{Var}(Y_1) = 31\sigma^2$$

$$Y_2 = Z_1 + \dots + Z_{28}$$

$$\text{Var}(Y_2) = 28\sigma^2$$

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

$$\Rightarrow w_1 = \frac{1}{31}$$

$$w_2 = \frac{1}{28}$$

j :



5.1 Weighted Least Square (WLS)

- **Model**

$$E(Y | X = x_i) = \beta' x_i, \quad Var(Y | X = x_i) = \frac{\sigma^2}{w_i}$$

assumed known

- **Examples of known w_i**

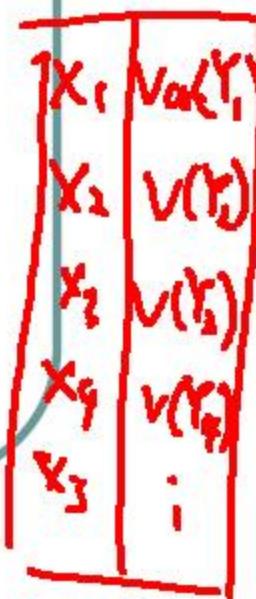
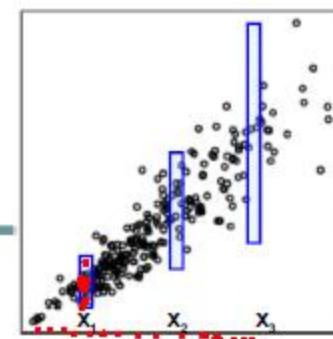
- The i -th Observations is an average of n_i variables

$$Y_i = \frac{z_{i1} + z_{i2} + \dots + z_{in_i}}{n_i}, \quad Var(Y_i | X = x_i) = \frac{Var(z_{i1} | X = x_i)}{n_i} = \frac{\sigma^2}{n_i} \Rightarrow w_i = n_i$$

- The i -th Observations is a sum of n_i variables

$$Y_i = z_{i1} + z_{i2} + \dots + z_{in_i}, \quad Var(Y_i | X = x_i) = n_i Var(z_{i1} | X = x_i) = n_i \sigma^2 \Rightarrow w_i = 1/n_i$$

- When sample size is large, estimate $Var(Y | X = x_i)$ by computing sample variance of the Y with X close to x_i .
- Subject knowledge...
- Guess from the scatterplot...





ordinary OLS : $RSS = \sum (y_i - \hat{y}_i)^2$

5.1 Estimators for the parameters

- Residual sum of square
 - Standardized by the variance of each observation

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

$$\begin{aligned}
 \text{RSS}(\beta) &= \sum \frac{(y_i - \hat{y}_i)^2}{\text{Var}(y_i)} = \frac{1}{\sigma^2} \sum w_i (y_i - \hat{y}_i)^2 \\
 &\propto \sum w_i (y_i - \hat{y}_i)^2 \\
 &= (y_1 - \hat{y}_1 \quad \dots \quad y_n - \hat{y}_n) \begin{pmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{pmatrix} \begin{pmatrix} y_1 - \hat{y}_1 \\ \vdots \\ y_n - \hat{y}_n \end{pmatrix} \\
 &= (Y - X\beta)' W (Y - X\beta)
 \end{aligned}$$



5.1 Estimators for the parameters

- Residual sum of square

$$\sqrt{W} = \begin{pmatrix} \sqrt{w_1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sqrt{w_n} \end{pmatrix}, \quad \sqrt{W} \times \sqrt{W} = W$$

$\sqrt{w} \curvearrowleft \curvearrowright \sqrt{w}$

$$\text{RSS}(\beta) = (Y - X\beta)' W (Y - X\beta) = (\sqrt{W} Y - \sqrt{W} X\beta)' (\sqrt{W} Y - \sqrt{W} X\beta)$$

Recall: In multiple linear regression,

$$\hat{\beta} = (X' X)^{-1} X' Y \quad \text{minimizes} \quad (Y - X\beta)' (Y - X\beta)$$

$$\hat{\beta}_w = (M' M)^{-1} M' Z \quad (Z - M\beta)' (Z - M\beta)$$

- Therefore, the WLS estimator is

$$\hat{\beta}_w = ((\sqrt{W} X)' \sqrt{W} X)^{-1} (\sqrt{W} X)' \sqrt{W} Y = (X' W X)^{-1} X' W Y$$

$$M = \sqrt{w} X$$

$$Z = \sqrt{Y} Y$$