# Biostat 250B HW1

## Elvis Cui

## January 15, 2021

## 1 Commentary

Recall the classical linear regression model:

$$\mathbf{y} = \mathbf{X}\beta + \epsilon$$
$$\epsilon \sim \mathcal{N}(0, \mathbf{W})$$

As practitioners, we would like to estimate $\beta$ based on observed data and address the goodness-of-fit based on some measurements such as **coefficient of determinants**.

There are 2 common ways to estimate $\beta$:

- Ordinary Least Squares (OLS): $\widehat{\beta}_{OLS} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$.

- Weighted Least Squares (WLS): $\widehat{\beta}_{WLS} = (\mathbf{X}^T\mathbf{W}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{W}^{-1}\mathbf{y}$.

For the second one, usually we write it in the **transformed** version, i.e.,

$$\widehat{\beta}_{WLS} = (\mathbf{X}_*^T\mathbf{X}_*)^{-1}\mathbf{X}_*^T\mathbf{y}_*$$

where

$$\mathbf{X}_* = \mathbf{W}^{-1/2}\mathbf{X},\ \mathbf{y}_* = \mathbf{W}^{-1/2}\mathbf{y}$$

Then, the coefficient of determination can be calcualted in three ways:

- 
$$R_{OLS}^2 = 1 - \frac{\|\mathbf{y} - \mathbf{X}\widehat{\beta}_{OLS}\|_2^2}{\mathbf{y}^{\mathbf{T}}\mathbf{y} - n\bar{y}^2}$$

- 
$$R_{WLS}^2 = 1 - \frac{\|\mathbf{y}_* - \mathbf{X}_*\widehat{\beta}_{WLS}\|_2^2}{\mathbf{y}_*^{\mathbf{T}}\mathbf{y}_* - n\bar{y}_*^2}$$

- 
$$\text{pseudo } R_{WLS}^2 = 1 - \frac{\|\mathbf{y} - \mathbf{X}\widehat{\beta}_{WLS}\|_2^2}{\mathbf{y}^{\mathbf{T}}\mathbf{y} - n\bar{y}^2}$$

The first two are more intuitive than the last one. While usually the second one ($R_{WLS}^2$) is higher than the first one ($R_{OLS}^2$), the third one (pseudo $R_{WLS}^2$) is always less than the first one. Therefore, "sole reliance on the **coefficient of determination** may fail to reveal important data characteristics and model inadequacies".