1. $Y \sim X$, $n$, $X_i$ fixed

   $\hat{\beta_0}, \hat{\beta_1}$ LS estimates.

   $\hat{\epsilon_i} = Y_i - (\hat{\beta_0} + \hat{\beta_1} X_i)$

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

   $\qquad \rightarrow \epsilon_i$ do not depend on $X$

   $E[\epsilon_i] = 0 \; \forall i, \; Var[\epsilon_i] = \sigma^2, \; \epsilon's \; iid \; \Big\}$ Could assume

a) For a particular $i \in 1...n$,

   LS estimates unbiased
   $\swarrow$

   $E[\hat{\epsilon_i}] = E[Y_i - \hat{\beta_0} - \hat{\beta_1} X_i] = E[Y_i] - \beta_0 - \beta_1 X_i$

   ⓘⓘ

   $\qquad = \beta_0 + \beta_1 X_i + 0 - \hat{\beta_0} - \hat{\beta_1} X_i \; - \text{assuming moment assumptions}$ $\rightarrow$

   $\qquad = 0$ if we assume $E[\epsilon_i] = 0$

   $\qquad\qquad\qquad$ and w/ LS estimates.

b) For some $i, j$, $i \neq j$,

   $Cov(\hat{\epsilon_i}, \hat{\epsilon_j})$

   $\quad = E\Big[(\hat{\epsilon_i} - E[\hat{\epsilon_i}])(\hat{\epsilon_j} - E[\hat{\epsilon_j}])\Big]$

   ⌐ From the previous part a):

   └ $= E[\hat{\epsilon_i} \cdot \hat{\epsilon_j}]$ under moment assumptions

   └→ Since $\epsilon_i, \epsilon_j$ also independent under moment assumptions:

   $\qquad$ fixed
   $\hat{\epsilon_i} = Y_i - \hat{Y_i} = \beta_0 + \beta_1 X_i + \epsilon_i - \hat{\beta_0} - \hat{\beta_1} X_i = \overline{(\beta_0 - \hat{\beta_0}) + (\beta_1 - \hat{\beta_1}) X_i} + \epsilon_i$

   $\qquad$ fixed
   $\hat{\epsilon_j} = Y_j - \hat{Y_j} = \beta_0 + \beta_1 X_j + \epsilon_j - \hat{\beta_0} - \hat{\beta_1} X_j = \overline{(\beta_0 - \hat{\beta_0}) + (\beta_1 - \hat{\beta_1}) X_j} + \epsilon_j$

   Then the functions of independent r.v.'s are also independent.

   so $\hat{\epsilon_i}, \hat{\epsilon_j}$ are also indep. under moment assumptions.

   so $Cov(\hat{\epsilon_i}, \hat{\epsilon_j}) = 0$ if $\nearrow$

   $\Rightarrow$ ⓘⓘ

c)     90% PI

   $\varepsilon_i$ NOT normal.



   PI assumes normal distribution so coverage should
   be way worse than "usual" since $\varepsilon_i \cancel{\in}$ Normal.

   overall , (ii)


d)   Bootstrapping the residual will NOT reproduce the issues.
     Specifically, this method by definition chooses $\varepsilon_i$ assuming
     a normal distribution, so it will ultimately "cover up"
     the lack of normality of Y.

                (ii)

2. $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$   SLR , $X_i$'s fixed

a) Reduce data set into $\{(X_i, Y_i): x_{0.05} \le X_i \le x_{0.95}\} = M$

aka middle 90% of the data (X values)

i) $\hat{\beta}_{1,new} = \dfrac{\sum\limits_{i \in M}(X_i - \bar{X})(Y_i - \bar{Y})}{\sum\limits_{i \in M}(X_i - \bar{X})^2}$   vs.   $\hat{\beta}_{1,old} = \dfrac{\sum\limits_{i \in whole}(X_i - \bar{X})(Y_i - \bar{Y})}{\sum\limits_{i \in whole}(X_i - \bar{X})^2}$

Note that New dataset $M$ is of size $0.9n \leftarrow n = $ original dataset size

$E[\hat{\beta}_{1,old}] = \beta_1$ is unbiased, we know.

We expect that $\hat{\beta}_{1,new}$ will still be an unbiased estimator for $\beta_1$ because the new dataset only excludes high leverage points based on large X values, Since $\beta_1 = \dfrac{S_{xy}}{S_{xx}}$, the magnitude of change will be similar for both the numerator and denominator, so the ratio $S_{xy} : S_{xx}$ should stay the same, so $\hat{\beta}_{1,new}$ will be unbiased.

ii) $Var[\hat{\beta}_{1,old}] = \dfrac{1}{\sum\limits_{i=1 \cdots n}(X_i - \bar{X})^2}$ , $Var[\hat{\beta}_{1,new}] = \dfrac{1}{\sum\limits_{i=0.05n}^{0.95n}(X_i - \bar{X})^2}$

Note that $\sum\limits_{i=0.05n}^{0.95n}(X_i - \bar{X})^2 = \sum\limits_{i=1}^{n}(X_i - \bar{X})^2 - \underbrace{\sum\limits_{i=1}^{0.05n}(X_i - \bar{X})^2}_{large} - \underbrace{\sum\limits_{i=0.95n}^{n}(X_i - \bar{X})^2}_{large}$

Clearly, by creating the new dataset, we are removing those points that were unusually large or tiny (aka unusually far from $\bar{X}$). So,

$\sum\limits_{i=0.05n}^{0.95n}(X_i - \bar{X})^2 < \sum\limits_{i=1}^{n}(X_i - \bar{X})^2$ by definition. Thus, we expect

$Var[\hat{\beta}_{1,new}]$ to be larger than $Var[\hat{\beta}_{1,old}]$ from LS, since the denominator is much smaller than before. Also, a smaller multiple of $n$ in the denom contributes to variance being larger.

b) Now, keep all data points. Instead, truncate the values:

$$\tilde{x}_i = \begin{cases} x_{0.05} & , \ x_i < 0.05 \\ x_i & , \ x_{0.05} < x_i < x_{0.95} \\ x_{0.95} & , \ x_i > 0.95 \end{cases}$$

↗ but not assigning zeros

Note that this is similar to using Least Trimmed Squares instead of Least Squares. This method will for sure be more robust to the impact of outliers.

i) Now, we assign the same $x, y$ values of $x_{0.05} \to y_{0.05}$ and $x_{0.95} \to y_{0.95}$ to all points outside the middle 90% of points and refit the model with all $n$ of these points

In $\hat{\beta}_{1_{new}} = \frac{S_{xy_{new}}}{S_{xx_{new}}}$ vs. $\hat{\beta}_{1_{old}} = \frac{S_{xy_{old}}}{S_{xx_{old}}}$,

The new model is now: $\tilde{y}_i = \begin{cases} \beta_0 + \beta_1 \tilde{x}_i + \epsilon_i & \text{if } \tilde{x}_i = x_i \\ y_{0.05} & \text{if } \tilde{x}_i = x_{0.05} \\ y_{0.95} & \text{if } \tilde{x}_i = x_{0.95} \end{cases}$
of $\hat{y}$

So, $\hat{\beta}_{1_{new}}$ will now be biased

ii) $Var[\hat{\beta}_{1_{new}}]$ should be lower than $\hat{\beta}_{1_{now}}$

c) Now remove pts by Y values, not X.

i) $\hat{\beta}_{1,new}$ will be biased now since $E\left[\tilde{\beta}_{1,new}\right] = E\left[\dfrac{S_{xy}}{S_{xx}}\right]$

will be affected by change in $\sum\limits_{0.05r}^{0.95r}(y_i - \bar{y})^2$

3. weight ~ temp + hum + fert , n=37
   ↳ all centered + scaled.

a) Residual standard error $\hat{\sigma} = \sqrt{\dfrac{\sum(y_i - \hat{y})^2 \;\&\; RSS_{full}}{df = 33}}$

$= \sqrt{\dfrac{RSS}{33}}$ ← $RSS = y^T(I-H)\cdot y$

$H = X(X^TX)^{-1}X^T$

$F = \dfrac{RSS_{partial}\,/\,3}{RSS\,full\,/\,33}_{\;11} = 9.311 = \dfrac{11\cdot RSS_{part}}{RSS\,full} =$

b) Adj $R^2 = 1 - (1-R^2)\dfrac{n-1}{n-p-1}$

$= 1 - (1 - 0.4584)\cdot \dfrac{36}{32}$

$= 1 - 0.6093 = \boxed{0.3907}$

c) The standard errors for hum and temp being quite high while for fert and intercept being low is likely due to collinearity between hum and temp. Since these covariates are very dependent on each other, they in turn reduce each others' significance in this full model, (both have high p-values) while fert is relatively independent of them so its variance estimate is definitely lower.

d). $\frac{\overset{n=37}{\hat{\sigma}^2}}{\hat{\sigma}_{12}} \sim \frac{\chi^2_{33}}{\chi^2_{16}} \sim \boxed{F_{33,16}}$

$\Rightarrow$
$n=20$

e)