

$$1) \quad a) \quad Y \sim N_n(\mu, \sigma^2 I_n)$$

$$\Rightarrow AY \sim N_n(A\mu, A\sigma^2 I_n A^T)$$

$$BY \sim N_n(B\mu, B\sigma^2 I_n B^T)$$

$$\begin{aligned} \text{Cov}(AY, BY) &= E[(AY - A\mu)(BY - B\mu)^T] \\ &= E[A(Y - \mu)[B(Y - \mu)]^T] \\ &= E[A(Y - \mu)(Y - \mu)^T B^T] \\ &= A E[(Y - \mu)(Y - \mu)^T] B^T \\ &= A \sigma^2 I_n B^T = \sigma^2 (AB^T) = \sigma^2 \cdot 0 = 0 \end{aligned}$$

$\Rightarrow AY$ and BY are independent

$$b) \quad AY \text{ and } BY \text{ are independent} \Rightarrow (AY)(BY)^T = 0$$

We have $(Y^T AY)^T (Y^T BY)$

$$= Y^T A^T Y Y^T B Y$$

$$= Y^T A Y (BY)^T Y \quad (\text{as } A^T = A, B^T = B \text{ by assumption})$$

$$= Y^T (0) Y = 0$$

$\Rightarrow Y^T AY$ and $Y^T BY$ are also independent

2) a) With $p = 3$, $y = X\beta + \epsilon$

$$\Rightarrow \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \epsilon_1 \\ \vdots \\ \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \epsilon_n \end{bmatrix}$$
$$= \begin{bmatrix} 1_n & x_{n1} & x_{n2} \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \epsilon$$

(i) With $H_0: \beta_2 = 0$

$$\Rightarrow X_{n \times 3} = \begin{bmatrix} 1_n & X_{11} & X_{12} \\ \vdots & \vdots & \vdots \\ X_{n1} & X_{n2} \end{bmatrix} = \begin{bmatrix} 1_n & X_1 & X_2 \end{bmatrix}_{n \times 3}$$

$$\text{And } L_{n \times 2} = \begin{bmatrix} 1_n & X_{11} \\ \vdots & \vdots \\ X_{n1} & X_{n2} \end{bmatrix} = \begin{bmatrix} 1_n & X_1 \end{bmatrix}_{n \times 2}$$

• Prove $\text{Im}(L) \subset \text{Im}(X)$

Let $a \in \mathbb{R}^n$ is a vector such that $a \in \text{Im}(L)$

$$\Rightarrow a = L_{n \times 2} v_{2 \times 1} \text{ for some } v \in \mathbb{R}^2$$

$$\Rightarrow a = \begin{bmatrix} 1_n & X_1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} v_1 + v_2 x_{11} \\ \vdots \\ v_1 + v_2 x_{n1} \end{bmatrix}$$

$$= \begin{bmatrix} v_1 + v_2 x_{11} + 0 x_{12} \\ \vdots \\ v_1 + v_2 x_{n1} + 0 x_{n2} \end{bmatrix} = \begin{bmatrix} 1_n & X_1 & X_2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ 0 \end{bmatrix}$$

$$\Rightarrow a \text{ can be written as } a = X \begin{bmatrix} v_1 \\ v_2 \\ 0 \end{bmatrix} \text{ for } \begin{bmatrix} v_1 \\ v_2 \\ 0 \end{bmatrix} \in \mathbb{R}^3$$

$$\Rightarrow a \in \text{Im}(X) \Rightarrow \text{Im}(L) \subset \text{Im}(X)$$

(ii) With $H_0: \beta_1 + \beta_2 = 0$, $X_{n \times 3} = \begin{bmatrix} 1_n & X_1 & X_2 \end{bmatrix}$

$$\Rightarrow H_0: \beta_1 = -\beta_2$$

$$\Rightarrow \text{Our model becomes } y_i = \beta_0 + \beta_1 (X_{i1} - X_{i2}) + \epsilon$$

$$\Rightarrow L_{n \times 2} = \begin{bmatrix} 1_n & (X_1 - X_2) \end{bmatrix}_{n \times 2}$$

Let $b \in \mathbb{R}^n$ is a vector in $\text{Im}(L)$

$$\Rightarrow b = L v \text{ for some } v \in \mathbb{R}^2$$

$$\Rightarrow b = \begin{bmatrix} 1_n & (X_1 - X_2) \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 1 & (x_{11} - x_{12}) \\ \vdots & \vdots \\ 1 & (x_{n1} - x_{n2}) \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$= \begin{bmatrix} v_1 + v_2 (x_{11} - x_{12}) \\ \vdots \\ v_1 + v_2 (x_{n1} - x_{n2}) \end{bmatrix} = \begin{bmatrix} v_1 + v_2 x_{11} - v_2 x_{12} \\ \vdots \\ v_1 + v_2 x_{n1} - v_2 x_{n2} \end{bmatrix}$$

$$= \begin{bmatrix} 1_n & X_1 & X_2 \end{bmatrix}_{n \times 3} \begin{bmatrix} v_1 \\ v_2 \\ -v_2 \end{bmatrix}_{3 \times 1} = X \begin{bmatrix} v_1 \\ v_2 \\ -v_2 \end{bmatrix}$$

$$\Rightarrow b \in \text{Im}(X) \Rightarrow \text{Im}(L) \subset \text{Im}(X)$$

$$b) SSE_x = Y^T(I - H_x)Y$$

$$SSE_L = Y^T(I - H_L)Y$$

$$\Rightarrow SSE_L - SSE_x = Y^T(I - H_L)Y - Y^T(I - H_x)Y \\ = Y^TY - Y^TH_LY - Y^TY + Y^TH_xY = Y^TH_xY - Y^TH_LY \\ = Y^T(H_x - H_L)Y$$

$$\Rightarrow f = \frac{(SSE_L - SSE_x)/(p-s)}{SSE_x/(n-p)} = \frac{Y^T(H_x - H_L)Y/(p-s)}{Y^T(I - H_x)Y/(n-p)}$$

c) • Symmetric

$$(H_x - H_L)^T = H_x^T - H_L^T$$

But H_x and H_L are symmetric themselves $\Rightarrow H_x = H_x^T, H_L = H_L^T$

$$\Rightarrow (H_x - H_L)^T = H_x - H_L$$

• Idempotent

$$(H_x - H_L)^2 = (H_x - H_L)(H_x - H_L) = H_x^2 - H_xH_L - H_LH_x + H_L^2 \\ = H_x - H_xH_L - H_LH_x + H_L$$

(because H_x, H_L are idempotent themselves $\Rightarrow H_x = H_x^2, H_L = H_L^2$)

Also $\text{Im}(L) \subset \text{Im}(X) \Rightarrow H_xH_L = H_LH_x = H_L$

$$\Rightarrow (H_x - H_L)^2 = H_x - H_L - H_L + H_L = H_x - H_L$$

$\Rightarrow H_x - H_L$ is idempotent

$$\bullet (I_n - H_x)(H_x - H_L) = I_nH_x - I_nH_L - H_x^2 + H_xH_L \\ = H_x - H_L - H_x + H_L = (H_x - H_x) + (H_L - H_L) = 0$$

d) From question 4) if $AB^T = 0 \Rightarrow AY$ & BY are independent

$$\text{From part c), } (I_n - H_x)(H_x - H_L)^T = 0$$

$\Rightarrow (I_n - H_x)Y$ and $(H_x - H_L)Y$ are independent

Also, $Y^T(H_x - H_L)Y$ is a quadratic form $\sim \chi^2$ with

$$df = \text{tr}(H_x - H_L) = \text{tr}(H_x) - \text{tr}(H_L) = p - s$$

$Y^T(I_n - H_x)Y$ is also $\sim \chi^2$ with $df = \text{tr}(I_n) - \text{tr}(H_x)$

$$\Rightarrow f = \frac{Y^T(H_x - H_L)Y/(p-s)}{Y^T(I - H_x)Y/(n-p)} \sim \frac{\chi_{(p-s)}^2/(p-s)}{\chi_{(n-p)}^2/(n-p)}$$

$$\Rightarrow f \sim F_{(p-s), (n-p)}$$

3) a) We want $\widehat{\bar{y}} = \bar{y}$ or $n^{-1} \mathbf{1}_n^T \hat{\mathbf{y}} = n^{-1} \mathbf{1}_n^T \mathbf{y}$

We have $\mathbf{1}_n \in \text{Im}(X) \Rightarrow \mathbf{1}_n = X\mathbf{v}$ for $\mathbf{v} \in \mathbb{R}^n$
 $(n \times 1) \quad n \times n \quad n \times 1$

$$\mathbf{1}_n^T \hat{\mathbf{y}} = \mathbf{v}^T X^T \hat{\mathbf{y}} = \mathbf{v}^T X^T (X(X^T X)^{-1} X^T \mathbf{y}) \quad (\text{because } \hat{\mathbf{y}} = H\mathbf{y})$$

$$= \mathbf{v}^T (X^T X (X^T X)^{-1}) X^T \mathbf{y} = \mathbf{v}^T \mathbf{I} X^T \mathbf{y} = \mathbf{v}^T X^T \mathbf{y}$$

$$\mathbf{1}_n^T \mathbf{y} = \mathbf{v}^T X^T \mathbf{y}$$

$$\Rightarrow \mathbf{1}_n^T \hat{\mathbf{y}} = \mathbf{1}_n^T \mathbf{y} \Rightarrow n^{-1} \mathbf{1}_n^T \hat{\mathbf{y}} = n^{-1} \mathbf{1}_n^T \mathbf{y} \Rightarrow \widehat{\bar{y}} = \bar{y}$$

b) We have $\sum (\hat{y}_i - \bar{y})^2 = SSR = \mathbf{y}^T (H - \frac{1}{n} J) \mathbf{y}$
 $\sum (y_i - \bar{y})^2 = SSTO = \mathbf{y}^T (\mathbf{I}_n - \frac{1}{n} J) \mathbf{y}$

with $J = \mathbf{1}_n \mathbf{1}_n^T$

as well as $\sum (\hat{y}_i - y_i)^2 = SSE = \mathbf{y}^T (\mathbf{I} - H) \mathbf{y}$

and $SSTO = SSR + SSE$

let $(\hat{y}_i - \widehat{\bar{y}}) = (\hat{y}_i - \bar{y}) = A$ (as $\widehat{\bar{y}} = \bar{y}$ in 3a)

$(y_i - \bar{y}) = B$

$$\Rightarrow \sum_{i=1}^n (\hat{y}_i - \widehat{\bar{y}})(y_i - \bar{y}) = \sum_{i=1}^n AB$$

$$= \sum_{i=1}^n (A^2 + B^2 - A^2 - B^2 - 2AB)$$

Let's look at $(A - B)^2 = A^2 - 2AB + B^2$

$$\Rightarrow AB = \frac{A^2 + B^2 - (A - B)^2}{2}$$

$$\Rightarrow \sum_{i=1}^n AB = \frac{1}{2} (\sum A^2 + \sum B^2 - \sum (A - B)^2)$$

$$= \frac{1}{2} (\sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \bar{y})^2 - \sum (\hat{y}_i - \bar{y} - y_i + \bar{y})^2)$$

$$= \frac{1}{2} (SSR + SSTO - SSE)$$

$$= \frac{1}{2} (SSR + SSR + SSE - SSE) = \frac{1}{2} \cdot 2 \cdot SSR$$

$$= SSR = \mathbf{y}^T (H - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) \mathbf{y}$$

We have $SSR \geq 0$, $\sqrt{\sum (\hat{y}_i - \widehat{\bar{y}})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2} = \sqrt{SSR} \cdot \sqrt{SSTO} > 0$

(because $SSR \geq 0$, $SSTO \geq 0$, and the term is in denominator)

$$\Rightarrow r_{\hat{y}, y} = \frac{SSR}{\sqrt{SSR} \sqrt{SSTO}} \geq 0 \Rightarrow \text{it will never be } < 0$$

c) $r_{\hat{y}, y}^2 = \left(\frac{SSR}{\sqrt{SSR} \sqrt{SSTO}} \right)^2 = \frac{SSR^2}{SSR \cdot SSTO} = \frac{SSR}{SSTO}$
 $= \frac{SSTO - SSE}{SSTO} = 1 - \frac{SSE}{SSTO} = R^2$

hw6

Giang Vu

10/1/2020

HOMEWORK 6

4)

a)

```
#read data
hw6_dt <- read.delim("/Users/giangvu/Desktop/STAT 2131 - Applied Stat Methods 1/HW/hw6/steam_text-1.txt")
#regress steam (Y) onto fat (X1) and glycerine (X2)
hw6_md <- lm(steam ~ fat + glycerine, data = hw6_dt)
#i)
hw6_sm<-summary(hw6_md) #mean model
hw6_anova <- summary(aov(hw6_md)) #variance model
hw6_sm
```

```
##
## Call:
## lm(formula = steam ~ fat + glycerine, data = hw6_dt)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.7977 -1.0015 -0.4424  1.0575  3.2397
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.625      2.247   2.058  0.0516 .
## fat            1.728      1.168   1.480  0.1529
## glycerine     -6.628      7.578  -0.875  0.3912
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.546 on 22 degrees of freedom
## Multiple R-squared:  0.1755, Adjusted R-squared:  0.1005
## F-statistic: 2.341 on 2 and 22 DF,  p-value: 0.1197
```

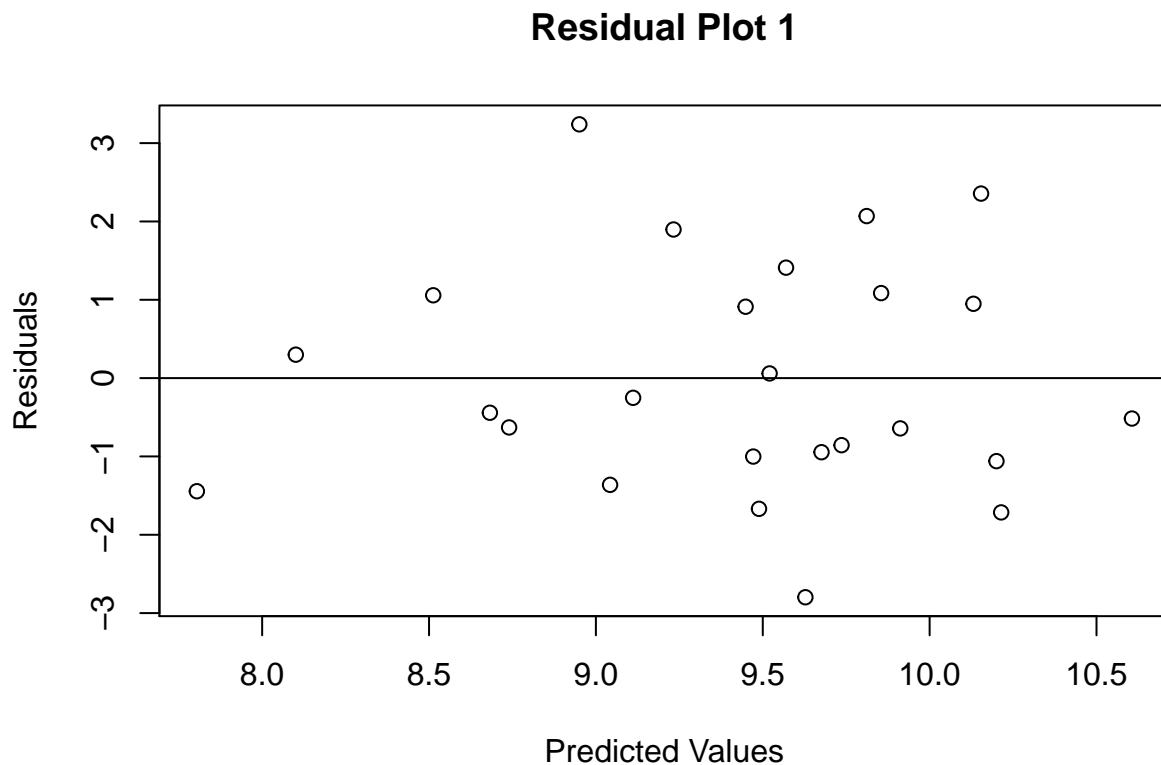
```
hw6_anova
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## fat            1   9.37   9.370   3.918 0.0604 .
```

```
## glycerine      1      1.83      1.829      0.765 0.3912
## Residuals     22     52.62      2.392
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

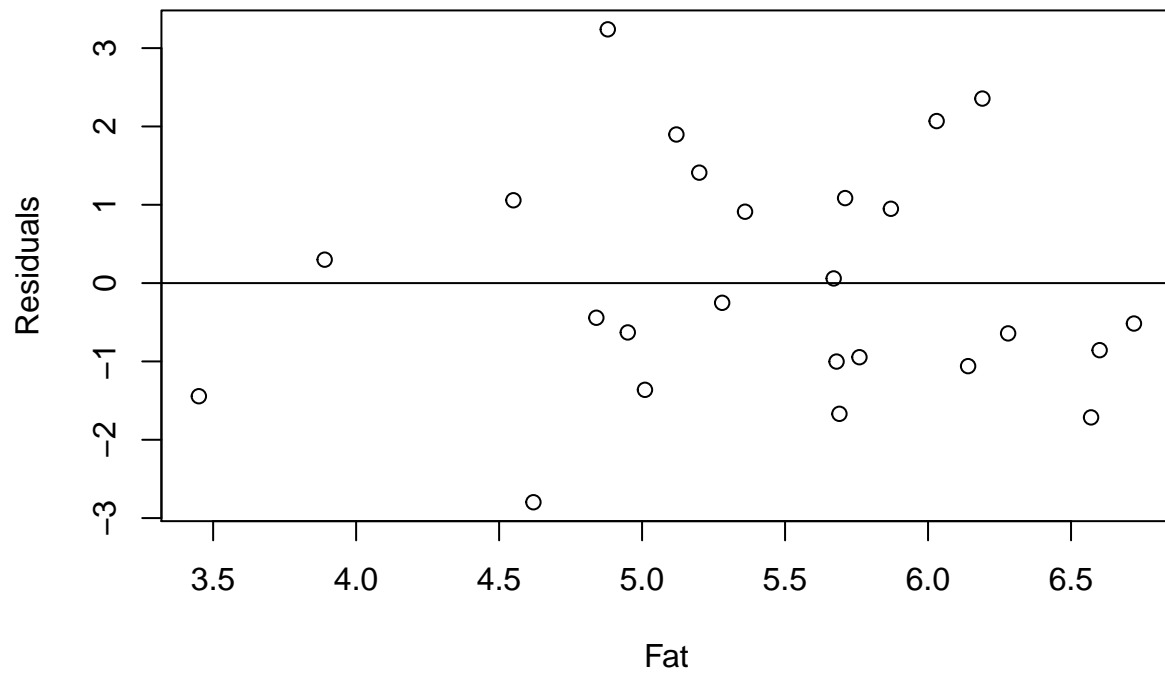
Interpreting the coefficients: As the value for variable fat increases by 1 unit, the value for steam is expected to increase by 1.728 units. As the value for variable glycerin increases by 1 unit, the value for steam is expected to decrease by -6.628 units. However, the coefficients estimates for both variables fat and glycerin in this model are not statistically significant given $\alpha = 0.05$

```
#ii) plot residual as fcn of  $\hat{Y}$ , fat, glycerine
plot(predict(hw6_md), resid(hw6_md),
      ylab = "Residuals", xlab = "Predicted Values", main = "Residual Plot 1")
abline(0,0)
```

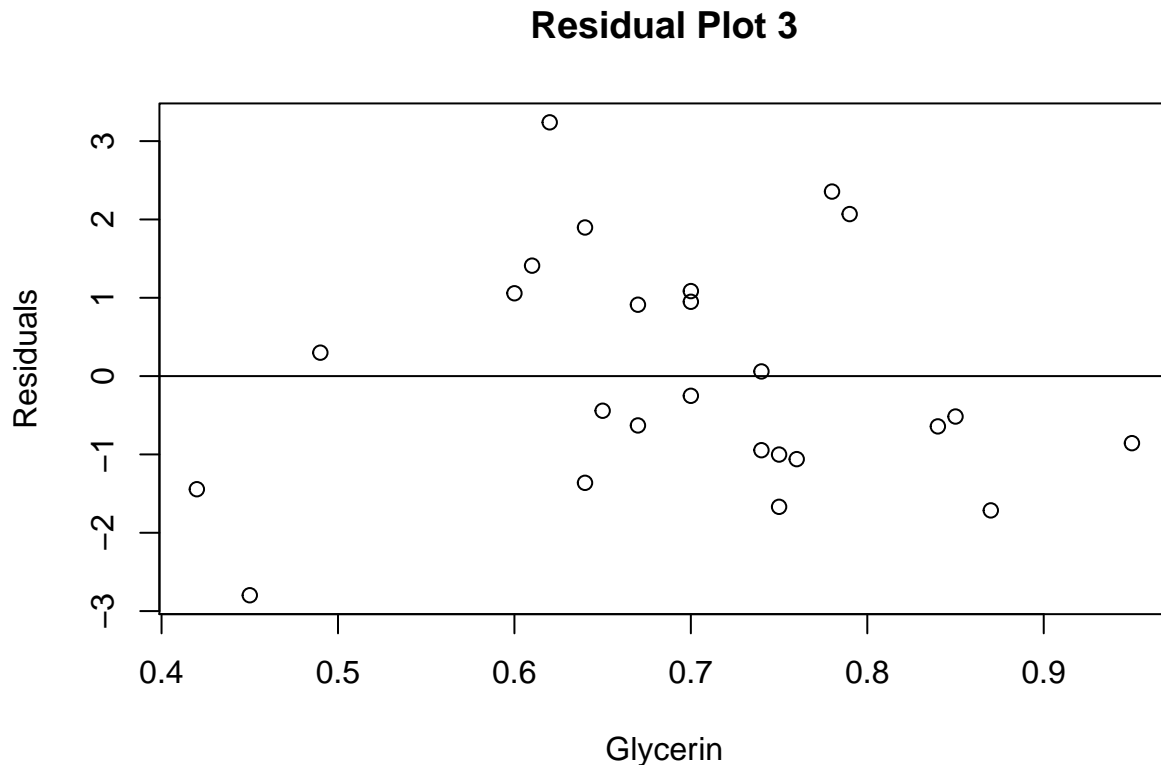


```
plot(hw6_dt$fat, resid(hw6_md),
      ylab = "Residuals", xlab = "Fat", main = "Residual Plot 2")
abline(0,0)
```

Residual Plot 2



```
plot(hw6_dt$glycerine, resid(hw6_md),  
      ylab = "Residuals", xlab = "Glycerin", main = "Residual Plot 3")  
abline(0,0)
```

Looking at the 3 plots, I don't see any relationship between residuals and glycerin, fat, or predicted Y.

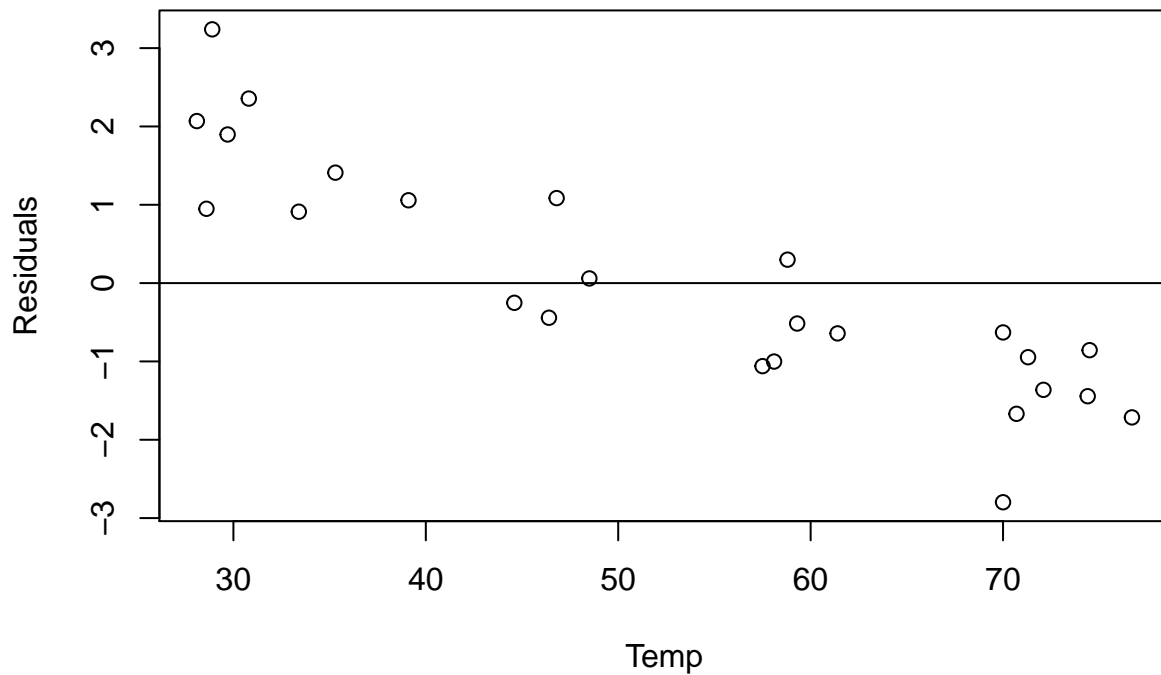
```
#iii)
linearHypothesis(hw6_md,c("fat=0","glycerine=0")) #F-test
```

```
## Linear hypothesis test
##
## Hypothesis:
## fat = 0
## glycerine = 0
##
## Model 1: restricted model
## Model 2: steam ~ fat + glycerine
##
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      24 63.816
## 2      22 52.617  2    11.199 2.3413 0.1197
```

p-value for F test is $0.1197 > 0.05 \Rightarrow$ fail to reject H_0 , also I noticed that R square of this model is really small too (0.1755) We can conclude that the two variables fat and glycerin cannot fully explain the variability of Y (steam). Their coefficients are not statistically significant.

```
#iv) temp against residual
plot(hw6_dt$temp, resid(hw6_md),
     ylab = "Residuals", xlab = "Temp", main = "Residual Plot 4")
abline(0,0)
```

Residual Plot 4



There's a decreasing pattern, so we can say that residuals can be explained by temp. This is probably what makes our original model a poor model.

b)

```
#i) regress steam (Y) onto fat (X1), glycerine (X2) and temp (X3)
hw6_md2 <- lm(steam ~ fat + glycerine + temp, data = hw6_dt)
summary(hw6_md2)
```

```
##
## Call:
## lm(formula = steam ~ fat + glycerine + temp, data = hw6_dt)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2348 -0.4116  0.1240  0.3744  1.2979
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.514814   1.062969   8.951 1.30e-08 ***
## fat           0.713592   0.502297   1.421   0.17
## glycerine     0.330497   3.267694   0.101   0.92
## temp        -0.079928   0.007884 -10.138 1.52e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.652 on 21 degrees of freedom
```

```
## Multiple R-squared:  0.8601, Adjusted R-squared:  0.8401
## F-statistic: 43.04 on 3 and 21 DF,  p-value: 3.794e-09
```

```
linearHypothesis(hw6_md2,c("fat=0","glycerine=0"))
```

```
## Linear hypothesis test
##
## Hypothesis:
## fat = 0
## glycerine = 0
##
## Model 1: restricted model
## Model 2: steam ~ fat + glycerine + temp
##
##      Res.Df    RSS Df Sum of Sq      F    Pr(>F)
## 1         23 18.223
## 2         21  8.927  2    9.2964 10.934 0.0005569 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

p value for F test = 0.0005569 < 0.05 => reject H0 We can say that fat and glycerine do have some impact on steam now that temp is included in the model. However, with a new variable temp added into the model and alpha = 0.05, the coefficient estimates for fat and glycerine are still not statistically significant (based on lm() results), but the coefficient estimate for temp is statistically significant.

- ii) Because of the inclusion of variable X3 (temp). Our original model is underfitting because we leave out an important variable - temp, so we saw that a lot of residuals are now accounted for by temp. With temp included, we have a better model, higher R squared.