

HOMEWORK 6

1. $Y \sim N(\mu, \sigma^2 I_n)$, $\mu \in \mathbb{R}^n$

$A, B \in \mathbb{R}^{n \times n}$

@ Show that if $AB^T = 0_{n \times n}$, AY and BY are ind.

For Y, A, B above,

$$AY \sim N(A\mu_y, A\sigma^2 I_n A^T), \quad BY \sim N(B\mu_y, B\sigma^2 I_n B^T)$$

(Proof: $E(AY) = AE(Y) = A\mu_y$

$$\text{Var}(AY) = E[(AY - A\mu_y)(AY - A\mu_y)^T]$$

$$= E[(AY - A\mu_y)(Y - \mu_y)^T A^T]$$

$$= AE[(Y - \mu_y)(Y - \mu_y)^T]A^T$$

$$= A\sigma^2 I_n A^T, \text{ and same for } BY)$$

If Y is as above, and $AB^T = 0$, the joint distribution of AY & BY is bivariate N .

To show independence of AY, BY , show that

$$\text{Cov}(AY, BY) = 0$$

$$\hookrightarrow \text{Cov}(AY, BY) = E[(AY - A\mu_y)(BY - B\mu_y)^T]$$

$$= E[A(Y - \mu_y)(Y - \mu_y)^T B^T]$$

$$= AE[(Y - \mu_y)(Y - \mu_y)^T]B^T = A\sigma^2 B^T = AB^T\sigma^2 = 0$$

So, AY and BY are independent

@ Suppose A, B are symmetric / idempotent, $AB = 0_{n \times n}$.

Show $Y^T A Y, Y^T B Y$ are independent.

Say $q_1 = Y^T A Y, q_2 = Y^T B Y$ and $AB = A^T T^T B = 0$

$$\Rightarrow T^T A T T^T B T = 0, \quad C = T^T A T \text{ and } D = T^T B T$$

Making $CK = (T^T A T)(T^T B T) = T^T A B T = T^T T = 0$, so $CK = KC$

$$\text{And } Q^T C Q = \begin{bmatrix} E_1 & 0 \\ 0 & 0 \end{bmatrix}, \quad Q^T D Q = \begin{bmatrix} 0 & 0 \\ 0 & E_2 \end{bmatrix}$$

Say $V = Q^T T^{-1} Y$, $E(V) = Q^T T^{-1} \mu$, ($V \rightarrow$ vector of standard normals)

$$\text{Var}(V) = Q^T T^{-1} \Sigma T^{-1} Q = I$$

so $Y = T Q V$ and $Y^T = V^T Q^T T^T$

$$q_1 = Y^T A Y = V^T Q^T T^T A T Q V = V^T Q^T T^T (T^{-1} C T^{-1}) T Q V$$

$$= V^T Q^T C Q V = V_1^T E_1 V_1, \text{ and } q_2 = V_2^T E_2 V_2$$

where V_1 is first half, V_2 is second half of V ,
so q_1 & q_2 are independent

2. $X \in \mathbb{R}^{n \times p}$, $p \leq n$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = Y = X\beta + E, \quad E \sim N_n(0_n, \sigma^2 I_n)$$

for some non-random $\beta \in \mathbb{R}^p$

$L \in \mathbb{R}^{n \times s}$ (full-rank design matrix), $s < p$

$$\text{im}(L) \subset \text{im}(X)$$

$$E(Y) = Ly \text{ for } y \in \mathbb{R}^s$$

Test $H_0: E(Y) = Ly \text{ for } y \in \mathbb{R}^s$

@ $p=3$, $E(Y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$

Find design matrices X, L , show that $\text{im}(L) \subset \text{im}(X)$

For $p=3$: $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$, 3×1 matrix $E = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \text{ and } X = \begin{bmatrix} x_1 & x_2 \\ x_1 & x_2 \\ \vdots & \vdots \\ x_1 & x_2 \end{bmatrix}_{n \times 2}$$

(i) $H_0: \beta_2 = 0$

$$L = \begin{bmatrix} 1 & x_1 \\ 1 & x_1 \\ \vdots & \vdots \\ 1 & x_1 \end{bmatrix}_{n \times 2}$$

x_i the same

(essentially removing a term)

For v , a 2×1 vector

and $a \in \text{im}(L)$

$$\Rightarrow a = [1_n \ x_1] v \rightarrow a = [1_n \ x_1] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$= \begin{bmatrix} v_1 + x_{11}v_2 + x_{12} \cdot 0 \\ \vdots \\ v_1 + x_{n1}v_2 + x_{n2} \cdot 0 \end{bmatrix} = [1 \ x_1 \ x_2] \begin{bmatrix} v_1 \\ v_2 \\ 0 \end{bmatrix}$$

$$\Rightarrow a \in \text{im}(X)$$

(ii) $H_0: \beta_1 + \beta_2 = 0$

$$L = \begin{bmatrix} 1 & x_{11} + x_{12} \\ 1 & x_{11} + x_{12} \\ \vdots & \vdots \\ 1 & x_{n1} + x_{n2} \end{bmatrix}_{n \times 1}$$

same x

(terms cancel)

$$Y = \beta_0 + \beta_2(x_{11} - x_{12}) + \epsilon_i$$

$$V = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}_{2 \times 1} \rightarrow a \in \text{im}(L)$$

$$a = [1_n \ x_1 \ x_2] \begin{bmatrix} v_1 \\ v_2 \\ v_2 \end{bmatrix}$$

$$= \begin{bmatrix} v_1 + v_2(x_{11} + x_{12}) & 0 \\ \vdots & \vdots \\ v_1 + v_2(x_{n1} + x_{n2}) & 0 \end{bmatrix} \rightarrow a \in \text{im}(X)$$

2, continued

$$\textcircled{b} H_X = X(X^T X)^{-1} X^T, \quad H_L = L(L^T L)^{-1} L^T$$

SSE_X = SSE when design matrix is XSSE_L = SSE when design matrix is L

$$F\text{-stat: } f = \frac{(SSE_L - SSE_X)/(p-s)}{SSE_X/(n-p)}$$

show that we can write f as:

$$f = \frac{Y^T(H_X - H_L)Y/(p-s)}{Y^T(I_n - H_X)Y/(n-p)}$$

$$\text{Denominator: } SSE_X = \sum \epsilon_i^2 = \sum (y_i - \hat{y}_i)^2$$

$$\begin{aligned} &= e^T e = (Y - XB)^T(Y - XB) = Y^T Y - 2B^T X^T Y + B^T X^T X B \\ &= Y^T Y - 2B^T X^T Y + B^T X^T X (X^T X)^{-1} X^T Y \quad (\text{sub for } B) \\ &= Y^T Y - 2B^T X^T Y + B^T I X^T Y = Y^T Y - Y^T X (X^T X)^{-1} X^T Y \\ &= Y^T (I - H_X) Y \end{aligned}$$

and rank(I - H) = n - p (degrees of freedom)

$$\text{so denominator: } SSE_X/(n-p) = Y^T(I - H_X)Y/(n-p)$$

Numerator:

$$SSE_X = Y^T(I - H_X)Y, \text{ so } SSE_L = Y^T(I - H_L)Y$$

$$\text{expand: } Y^T(I - H_L)Y = Y^T Y - Y^T L(L^T L)^{-1} L^T Y$$

$$\begin{aligned} SSE_L - SSE_X &= (Y^T Y - Y^T L(L^T L)^{-1} L^T Y) - (Y^T Y - Y^T X(X^T X)^{-1} X^T Y) \\ &= -Y^T L(L^T L)^{-1} L^T Y + Y^T X(X^T X)^{-1} X^T Y \\ &= -Y^T H_L Y + Y^T H_X Y = Y^T H_X Y - Y^T H_L Y = Y^T (H_X - H_L) Y \end{aligned}$$

and rank(H_X - H_L) = (n - s) - (n - p) = -s + p = p - s

$$\text{so } f = \frac{(SSE_L - SSE_X)/(p-s)}{SSE_X/(n-p)} = \frac{Y^T(H_X - H_L)Y/(p-s)}{Y^T(I_n - H_X)Y/(n-p)}$$

⊛ Show that (H_X - H_L) is symmetric and that (I_n - H_X)(H_X - H_L) = 0_{n×n}

(H_X - H_L) symmetric: (H_X - H_L)^T = (H_X - H_L)

$$\begin{aligned} (H_X - H_L)^T &= H_X^T - H_L^T = [(X(X^T X)^{-1} X^T)^T] - [(L(L^T L)^{-1} L^T)^T] \\ &= X[(X^T X)^{-1}]^T X^T - L[(L^T L)^{-1}]^T L^T = X[(X^T X)^T]^{-1} X^T - L[(L^T L)^T]^{-1} L^T \\ &= X(X^T X)^{-1} X^T - L(L^T L)^{-1} L^T = H_X - H_L \end{aligned}$$

next page

20, continued

Show that $(H_X - H_L)$ is idempotent

$$\begin{aligned}
 (H_X - H_L) &= (H_X - H_L)^2 = (H_X - H_L)(H_X - H_L) \\
 &= (X(X^T X)^{-1} X^T - L(L^T L)^{-1} L^T)(X(X^T X)^{-1} X^T - L(L^T L)^{-1} L^T) \\
 &= H_X^2 - X(X^T X)^{-1} X^T L(L^T L)^{-1} L^T - L(L^T L)^{-1} L^T X(X^T X)^{-1} X^T + H_L^2 \\
 &= H_X - H_X H_L - H_L H_X + H_L \quad \text{because } \text{im}(L) \subseteq \text{im}(X) \\
 &= H_X - 2H_L + H_L = H_X - H_L \quad \text{by properties of } H
 \end{aligned}$$

Have
idempotent,
proof in
HW5

Show that $(I_n - H_X)(H_X - H_L) = 0_{n \times n}$

$$\begin{aligned}
 (I_n - H_X)(H_X - H_L) &= I_n H_X - I_n H_L - H_X H_X + H_X H_L \\
 &= H_X - H_L - H_X + H_L = 0_{n \times n}
 \end{aligned}$$

① Use ① to prove $f \sim F(p-s, n-p)$ when
 $H_0: E(Y) = L\gamma$, $\gamma \in \mathbb{R}^s$ is true

From ①, $AB^T = 0$ means AY and BY are independent, and from the results in

$$\textcircled{2}, (I_n - H_X)(H_X - H_L)^T = (I_n - H_X)(H_X - H_L) = 0$$

so the numerator part & denominator part: $(H_X - H_L)Y$ & $(I_n - H_X)Y$ are independent (under H_0).

Additionally, $Y^T(H_X - H_L)Y$ is quadratic form,

χ^2 distributed with $df = \text{tr}(H_X - H_L)$

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

$$Y^T(I_n - H_X)Y \rightarrow \text{tr}(I_n - H_X) = \text{tr}(I_n) - \text{tr}(H_X) = n - p$$

Which tells us that under H_0 ,

$$f \sim F_{p-s, n-p} = \frac{Y^T(H_X - H_L)Y / p-s}{Y^T(I_n - H_X)Y / n-p}$$

3. $Y \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times p}$ full-rank design matrix and

$$R^2 = 1 - \frac{SSE}{SST_0} \quad (\text{regress } Y \text{ onto } X) \quad (= \frac{SSR}{SST_0})$$

$$\hat{Y} = \frac{1}{n} 1_n^T \hat{Y}, \quad \bar{Y} = \frac{1}{n} 1_n^T Y, \quad \text{and}$$

$$R^2_{YX} = \frac{\sum (Y_i - \hat{Y})(Y_i - \bar{Y})}{\sqrt{\sum (Y_i - \hat{Y})^2} \sqrt{\sum (Y_i - \bar{Y})^2}}$$

@ Use assumption $1_n \in \text{im}(X)$ to show that $\hat{Y} = \bar{Y}$

$$1_n = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \text{ of } n \text{ rows, } \in \text{im}(X)$$

$$\text{Show that } \hat{Y} = \bar{Y} \rightarrow \frac{1}{n} 1_n^T \hat{Y} = \frac{1}{n} 1_n^T Y$$

$$\frac{1}{n} 1_n^T \hat{Y} \Rightarrow \frac{1}{n} \text{ cancels from both sides}$$

$$\rightarrow 1_n^T \hat{Y} \Rightarrow 1_n^T \in \text{im}(X), \text{ so } 1_n^T = X^T v \text{ for some vector } v$$

$$\rightarrow (X^T v)^T \hat{Y} = v^T X^T \hat{Y} = v^T X^T X (X^T X)^{-1} X^T Y \quad (\hat{Y} = HY)$$

$$= v^T I X^T Y = 1_n^T Y$$

next page \rightarrow

3, continued

- ⑥ Use ③ to show that $\sum (\hat{y}_i - \bar{y})(y_i - \bar{y}) = Y^T (H - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) Y$ will $r^2_{y,x}$ ever be smaller than 0? why?

From ③, $\hat{y} = \hat{y}$, so we have

$$[\sum (\hat{y}_i - \bar{y})(y_i - \bar{y})]^2 = Y^T (H - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) Y$$

we have $J = \mathbf{1}_n \mathbf{1}_n^T$

$$\sum (\hat{y}_i - \bar{y})(y_i - \bar{y}) \cdot \sum (\hat{y}_i - \bar{y})(y_i - \bar{y}) = (SSR)(SST_0)$$

$$= [Y^T (H - \frac{1}{n} J) Y] [Y^T (I - \frac{1}{n} J) Y]$$

$$= Y^T (H - \frac{1}{n} J) Y Y^T (I - \frac{1}{n} J) Y$$

$$= Y^T (H - \frac{1}{n} J) (I - \frac{1}{n} J) Y$$

$$= Y^T (H I - \frac{1}{n} H J - \frac{1}{n} J I + (\frac{1}{n} J)^2) Y$$

$$= Y^T (H - \frac{1}{n} J - \frac{1}{n} J + \frac{1}{n} J) Y$$

$$HJ = J \text{ and } (\frac{1}{n} J)^2 = \frac{1}{n} J$$

$$= Y^T (H - \frac{1}{n} J) Y$$

$r^2_{y,x}$ will always be ≥ 0 because it's the square root of SSR (always ≥ 0), and denominator is squared, so also, always ≥ 0 . So, $r^2_{y,x}$ always ≥ 0

- ⑦ Show that $R^2 = r^2_{y,x}$

$$r^2_{y,x} = \frac{\sum (\hat{y}_i - \bar{y})(y_i - \bar{y})}{\sqrt{\sum (\hat{y}_i - \bar{y})^2} \sqrt{\sum (y_i - \bar{y})^2}} \text{ and from ③, } \hat{y} = \bar{y}$$

$$\text{so } (r^2_{y,x})^2 = \frac{(\sum (\hat{y}_i - \bar{y})(y_i - \bar{y}))^2}{\sum (\hat{y}_i - \bar{y})^2 \sum (y_i - \bar{y})^2}$$

$$\text{and from ③, } \sum (\hat{y}_i - \bar{y})(y_i - \bar{y}) = Y^T (H - \frac{1}{n} J) Y$$

$$\rightarrow = [Y^T (H - \frac{1}{n} J) Y]^2$$

$$= [Y^T (H - \frac{1}{n} J) Y] [Y^T (I - \frac{1}{n} J) Y]$$

$$= \frac{Y^T (H - \frac{1}{n} J) Y}{Y^T (I - \frac{1}{n} J) Y}$$

$$\Rightarrow \frac{SSR}{SST_0}$$

← sub matrix form, SSR $\hat{=}$ SST₀

$$= 1 - \frac{SSE}{SST_0} = R^2$$

2131 HW6

Orly Olbum

Question 4

(a) Suppose you regress steam (Y) onto fat (X1) and glycerine (X2).

(i) Write down the model you are assuming when performing this regression (i.e. what is the mean and variance model). Provide an interpretation for the coefficients in the mean model.

```
model = lm(steam ~ fat + glycerine, data = steam)
summary(model)
```

```
##
## Call:
## lm(formula = steam ~ fat + glycerine, data = steam)
##
## 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
```

```
summary(aov(model))
```

```
##              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
```

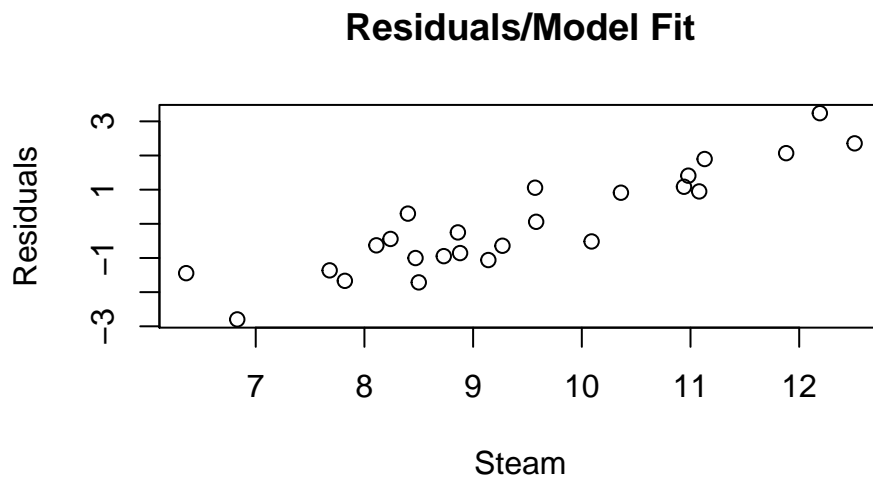
For a first order linear regression model, we assume: - the mean model is linear in both X variables - the regression function is a plane - the regression function is linear in both X variables - the association between Y and one of the X variables does not depend on the other X variable - we assume independent residuals (and normally distributed)

The fitted linear model above shows an equation of: $\text{steam} = 4.625 + 1.728\text{xfat} + -6.628\text{xglycerine}$, meaning for every unit increase in steam, there is a 1.728 unit increase in fat and a 6.628 unit decrease in glycerine. At $\alpha = 0.05$, neither of these variables prove to be significant as linear predictors for steam output. This leaves us with an underfit model, because we are leaving out variable that could explain the variation in steam.

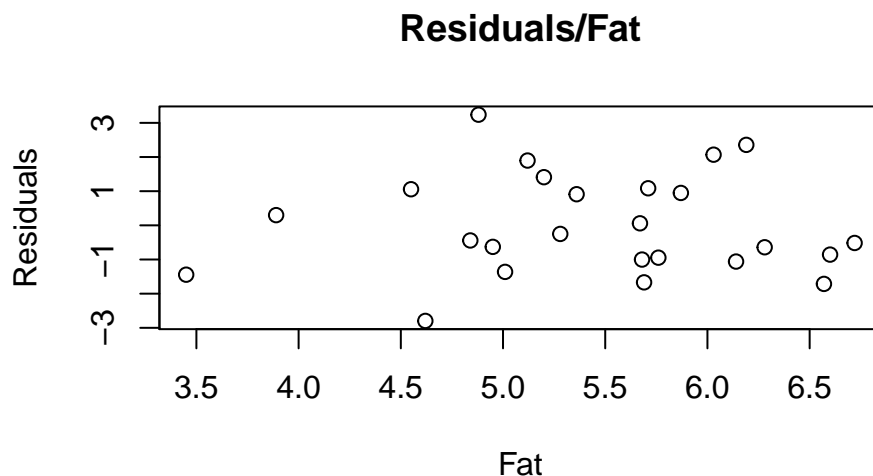
The variance model (ANOVA output) also does not show significance at $\alpha = 0.05$, but it does show us that there is more variation due to the fat variable than the glycerine variable in steam.

(ii) In separate plots, plot \hat{e} as a function of \hat{Y} , fat and glycerine. Do you see any evidence that the mean or variance model is incorrect?

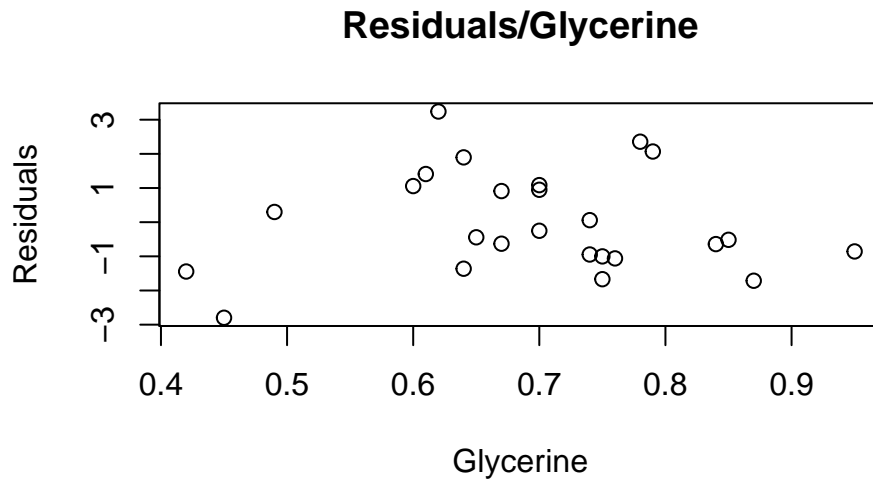
```
model.res = resid(model)
plot(steam$steam, model.res, ylab = "Residuals", xlab = "Steam", main="Residuals/Model Fit")
```



```
plot(steam$fat, model.res, ylab = "Residuals", xlab = "Fat", main="Residuals/Fat")
```




```
plot(steam$glycerine, model.res, ylab = "Residuals", xlab = "Glycerine", main="Residuals/Glycerine")
```



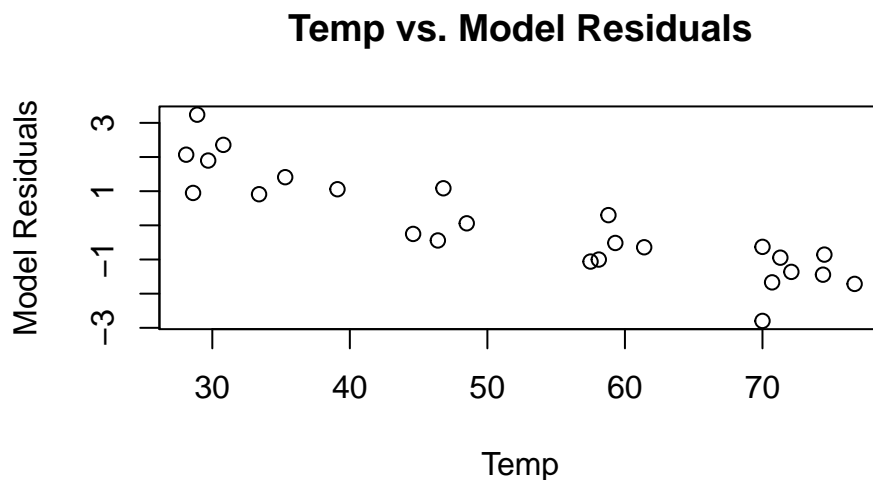
The first plot shows a strong positive association between steam and the model residuals. Because the R^2 of the model is very low, this relationship indicates a potentially poor model. If the R^2 was higher, the dependent variable's variation would be explained more by the independent variables. This is not the case here. The second plot and third plots show no (obvious) association between the variables and the model residuals, indicating a correct model.

(iii) Consider the null hypothesis that the coefficients for both fat and glycerine are 0. At a significance level of $\alpha = 0.05$, what do you conclude about these coefficients?

At $\alpha = 0.05$, we fail to reject the hypothesis of coefficients being 0 and conclude that we do not have evidence to support fat and glycerine being significant predictors for steam.

(iv) Plot the variable "temp" against the residuals from this regression. What can you conclude from this plot?

```
plot(steam$temp, model.res, ylab = "Model Residuals", xlab = "Temp", main = "Temp vs. Model Residuals")
```



We see in this plot that there is a negative association between temp and the model residuals.

(b) Now regress steam (Y) onto fat (X1), glycerine (X2) and temp (X3).

```
model2 = lm(steam ~ fat + glycerine + temp, data = steam)
summary(model2)
```

```
##
## Call:
## lm(formula = steam ~ fat + glycerine + temp, data = steam)
##
## 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
```

```
summary(aov(model2))
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## fat           1   9.37    9.37  22.042 0.000124 ***
## glycerine     1   1.83    1.83   4.304 0.050512 .
## temp          1  43.69   43.69 102.776 1.52e-09 ***
## Residuals    21   8.93    0.43
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(i) Consider the null hypothesis that the coefficients for both fat and glycerine are 0. At a significance level of $\alpha = 0.05$, what do you conclude about these coefficients?

From the mean model, we can see that both fat and glycerine still are not significant predictors of steam, but the variance model shows that with temp in the mix as another independent variable, fat is a significant predictor. Glycerine remains insignificant, at $\alpha = 0.05$.

(ii) Why are the P values from this test so much smaller than those from part (a)?

The R^2 of this model has risen by a lot, indicating that the independent variables are now accounting for much more of the variation in the dependent variable (steam) than in the prior model. Because of this, the p-values will be lower, because they are dependent on the F-values, which in turn depends on the mean squared error (and therefore variance). We see that in both the mean and variance models, temp is a significant predictor for steam. A plot of the residuals for this model against steam will show no obvious association, telling us that the variation in steam is explained by the model itself (more so than the first model).

```
model2.res = resid(model2)
plot(steam$steam, model2.res, ylab = "Residuals", xlab = "Steam", main="Residuals/Model 2 Fit")
```

