# Solution for HW2, MATH3805

1. (a)
$$Y = X\beta + \epsilon$$

$$X = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix}^T, Y = \begin{pmatrix} y_1 & y_2 & \cdots & y_n \end{pmatrix}^T$$

Then we can write the derivatives of the sum of squared errors and set it equal to 0. Then the normal equation is

$$X^\top X \hat{\beta} - X^\top Y = 0$$

(b)

(c) $$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \left(X^T X\right)^{-1} X^T Y$$

$$\left(X^T X\right)^{-1} X^T Y = \frac{1}{n\sum x_i^2 - \left(\sum x_i\right)^2} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

So we have
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$
$$\hat{\beta}_0 = \bar{y} - \bar{x}\hat{\beta}_1$$

where $\bar{y} = \sum_{i=1}^n y_i$ and $\bar{x} = \sum_{i=1}^n x_i$ That is

$$\hat{\beta}_0 = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n\sum x_i^2 - \left(\sum x_i\right)^2}$$
$$\hat{\beta}_1 = \frac{n\sum x_i y_i - \sum x_i \sum y_i}{n\sum x_i^2 - \left(\sum x_i\right)^2}$$

(d)

(e)
$$\begin{bmatrix} \text{var}\left(\hat{\beta}_0\right) & \text{cov}\left(\hat{\beta}_0, \hat{\beta}_1\right) \\ \text{cov}\left(\hat{\beta}_0, \hat{\beta}_1\right) & \text{var}\left(\hat{\beta}_1\right) \end{bmatrix} = \sigma^2 \left(X^T X\right)^{-1}$$

$$\text{var}\left(\hat{\beta}_0\right) = \frac{\sigma^2 \sum x_i^2}{n \sum x_i^2 - \left(\sum x_i\right)^2}$$

$$\text{var}\left(\hat{\beta}_1\right) = \frac{n\sigma^2}{n \sum x_i^2 - \left(\sum x_i\right)^2}$$

2. Exercise B.16

   (a) $Y = (4, 3, 3, 1, -1)^T, X = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \end{pmatrix}^T$

   (b) $X^T X = \begin{pmatrix} 5 & 0 \\ 0 & 10 \end{pmatrix}, X^T Y = \begin{pmatrix} 10 \\ -12 \end{pmatrix}$

   (c) $\hat{\beta} = \left(X^T X\right)^{-1} X^T Y = \begin{pmatrix} 2 \\ -1.2 \end{pmatrix}$

   (d) $\hat{y} = 2 - 1.2x$

3. Exercise B.19
$$H_0 : \beta_1 = 0 \quad \text{vs} \quad H_1 : \beta_1 \neq 0$$
$$s = 0.7303$$
$$t = \frac{\hat{\beta}_1}{s_{\hat{\beta}_1}} = \frac{\hat{\beta}_1}{s\sqrt{c_{11}}} = -5.1962$$

   p-value: 0.01385. or Reject region: $|t| > t_{\frac{\alpha}{2},3}$

   Reject $H_0, \beta_1 \neq 0$

4. Exercise B.22
$$E(y|x = 1) = 0.8$$
$$\text{var}(E[y|x = 1]) = s^2 x_p^T \left(X^T X\right)^{-1} x_p = 0.16$$

   90% confidence interval for mean is $(-0.1413, 1.7413)$

   $E(y|x = 1)$ will fall into the confidence interval with probability 90%.

   Exercise B.23

   $\text{var}(y|x = 1) = s^2 + s^2 x_p^T \left(X^T X\right)^{-1} x_p = 0.6933$

   90% prediction interval is $(-1.1593, 2.7593)$

   $y$ will fall into the prediction interval with probability 90%.

5. Exercise B.30

(a)

$$y = \begin{pmatrix} 5.2 \\ 0.3 \\ -1.2 \\ 2.2 \\ 6.2 \\ 5 \\ -0.1 \\ -1.1 \\ 2.0 \\ 6.1 \end{pmatrix}, X = \begin{pmatrix} 1 & -2 & 2 \\ 1 & -1 & -1 \\ 1 & 0 & -2 \\ 1 & 1 & -1 \\ 1 & 2 & 2 \\ 1 & -2 & 2 \\ 1 & -1 & -1 \\ 1 & 0 & -2 \\ 1 & 1 & -1 \\ 1 & 2 & 2 \end{pmatrix}$$

(b) $\hat{\beta} = \left(X^T X\right)^{-1} X^T Y = (2.46, 0.41, 1.6143)^T$

$y = 2.46 + 0.41 x_1 + 1.6143 x_2$

(c) $SSE = \left(Y - X\hat{\beta}\right)^T \left(Y - X\hat{\beta}\right) = 2.4363$

$s^2 = SSE/(7) = 0.3480$

(d)

$$F-value = \frac{\left[(y - \bar{y})^T (y - \bar{y}) - SSE\right]/2}{SSE/7} = \frac{(78.7640 - 2.4363)/2}{2.4363/7} = 109.6534$$

$$p - value = 1 - F_{(2,7)}(109.6661) = 0.0000.\text{or}$$
$$F - value > F_{(2,7)}^{-1}(0.95) = 4.74$$

Reject $H_0$. The model contribute information for predicting y.

(e)

$$R^2 = \frac{SSR}{SSR + SSE} = 0.969$$

96.9% of y's variation in samples can be explained by the model.

(f)

$$\text{var}(\hat{\beta}) = s^2 \left(X^T X\right)^{-1} = \begin{pmatrix} 0.0348 & 0 & 0 \\ 0 & 0.0174 & 0 \\ 0 & 0 & 0.0124 \end{pmatrix}$$

$$t - value = \frac{\hat{\beta}}{s_{\hat{\beta}}} = \frac{0.41}{\sqrt{0.0174}} = 3.1082$$

p-value: 0.01713. or $|t-value| > 2.3646$

Reject $H_0, \beta_1 \neq 0$

The practical implication is the extrusion pressure will effect the strength of the new plastic.

3

(g)
$$\hat{y} = x\hat{\beta} = 4.8686$$
$$\text{Var}(E(\hat{y})) = s^2 x^T \left(X^T X\right)^{-1} x = 0.1541$$

90% confidence interval for mean is
$$\left(\hat{y} \pm t_{0.05,7}\sqrt{\text{Var}(E(\hat{y}))}\right) = (4.1248, 5.6124)$$

(h)
$$\text{Var}(\hat{y}) = s^2 + s^2 x^T \left(X^T X\right)^{-1} x = 0.5022$$

90% prediction interval is
$$\left(\hat{y} \pm t_{0.05,7}\sqrt{\text{Var}(\hat{y})}\right) = (3.5260, 6.2111)$$

6. Exercise 4.11

(a)
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

(b)
$$\hat{y} = 21087.951 + 108.451x_1 + 557.910x_2 - 340.166x_3 + 85.681x_4$$

(c) Holding the value of the remaining variables fixed, the mean change in y for every 1-unit increase in $x_1$ is 108.451. For $x_2$, the mean change is 557.910. For $x_3$, the mean change is -340.166. For $x_4$, the mean change is 85.681.

(d) SPSS output:
T-statistic for $\beta_1$ : 1.222 and p-value: $0.236 > \alpha$. we can not reject $H_0$. So $x_2$ is a useless predictor.

(e) $R^2 = 0.912$ $R_a^2 = 0.894$. $R^2$ represents that 91.2% variation of $y$ could be explained. And $R_a^2$ represents that 89.4% variation in $y$ can be explained when considering the sample size and the number of parameters. Like $R^2$, adjusted $R^2$ also evaluates how many percent of the variation in $y$ can be explained by the multiple regression model. However, unlike $R^2$, adjusted $R^2$ takes into account (adjusted for) both the sample size and the number of parameters such that a model of more parameter will have a heavy penalty so that adjusted $R^2$ cannot be forced to 1 by simply adding more and more parameters. Thus the $R_a^2$ will be preferred as it takes the sample size and the number of parameters into account.

(f) F-value = 51.720 p-value = 0.0000.
reject $H_0$, at least one should not be 0.

**Model: MODEL1**
**Dependent Variable: RFEWIDTH**

| Number of Observations Read | 25 |
|---|---|
| Number of Observations Used | 25 |

| Analysis of Variance | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 4 | 49163 | 12291 | 51.72 | <.0001 |
| Error | 20 | 4752.76913 | 237.63846 | | |
| Corrected Total | 24 | 53915 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 15.41553 | R-Square | 0.9118 |
| Dependent Mean | 88.32000 | Adj R-Sq | 0.8942 |
| Coeff Var | 17.45417 | | |

| Parameter Estimates | | | | | |
|---|---|---|---|---|---|
| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
| Intercept | 1 | 21088 | 18553 | 1.14 | 0.2691 |
| REDSHIFT | 1 | 108.45084 | 88.73979 | 1.22 | 0.2359 |
| LINEFLUX | 1 | 557.90980 | 315.99021 | 1.77 | 0.0927 |
| LUMINOSITY | 1 | -340.16553 | 320.76260 | -1.06 | 0.3016 |
| AB1450 | 1 | 85.68102 | 6.27334 | 13.66 | <.0001 |

Figure 1: SAS output without intercept for Exercise 4.11

7. Exercise 4.13

   (a) The first order model is:

   $$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$$

   where RPM is $x_1$, CPRATIO is $x_2$, INLETTEMP is $x_3$, EXHTEMP is $x_4$, AIRFLOW is $x_5$, HEATRATE is $y$.

   (b)

   $$\hat{y} = 13614 + 0.08879 x_1 + 0.3519 x_2 - 9.2009 x_3 + 14.3939 x_4 - 0.8480 x_5$$

   (c) $\beta_0$ represent the $y$-intercept of the line and $\beta_1$ represent the slope. Holding the value of the remaining variables fixed, the mean change in $y$ for every 1-unit increase in RPM $x$ is 0.08879. For CPRATIO ($x_2$), the mean change is 0.3519. For INLETTEMP ($x_3$), the mean change is 9.2009. For EXHTEMP ($x_4$), the mean change is 14.3939. For AIRFLOW ($x_5$), the mean change is -0.8480.

   (d) $s = 458.8284$. It values the variation of $y$, is an estimator of $\sigma$. mean $\pm 2s$ provide a rough confidence interval.

   (e) The adjusted $R^2$ is 0.9172. $R_a^2$ represents that 91.72% variation in $y$ can be explained when considering the sample size and the number of parameters. Like $R^2$, adjusted $R^2$ also evaluates how many percent of the variation in $y$ can be explained by the multiple regression model. However, unlike $R^2$, adjusted $R^2$ takes into account (adjusted for) both the sample size and the number of parameters such that a model of more parameter will have a heavy penalty so that adjusted $R^2$ cannot be forced to 1 by simply adding more and more parameters. Thus the $R_a^2$ will be preferred as it takes the sample size and the number of parameters into account.

   (f) F-value is 147.30, and P-value $< 0.0001$. So the overall model is useful.

## The SAS System

### The REG Procedure
### Model: MODEL1
### Dependent Variable: HEATRATE

| Number of Observations Read | 67 |
|---|---|
| Number of Observations Used | 67 |

| Analysis of Variance | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 5 | 155055273 | 31011055 | 147.30 | <.0001 |
| Error | 61 | 12841935 | 210524 | | |
| Corrected Total | 66 | 167897208 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 458.82843 | R-Square | 0.9235 |
| Dependent Mean | 11066 | Adj R-Sq | 0.9172 |
| Coeff Var | 4.14613 | | |

| Parameter Estimates | | | | | |
|---|---|---|---|---|---|
| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
| Intercept | 1 | 13614 | 870.01294 | 15.65 | <.0001 |
| RPM | 1 | 0.08879 | 0.01391 | 6.38 | <.0001 |
| CPRATIO | 1 | 0.35190 | 29.55568 | 0.01 | 0.9905 |
| INLETTEMP | 1 | -9.20087 | 1.49920 | -6.14 | <.0001 |
| EXHTEMP | 1 | 14.39385 | 3.46095 | 4.16 | 0.0001 |
| AIRFLOW | 1 | -0.84796 | 0.44211 | -1.92 | 0.0598 |

Figure 2: SAS output without intercept for Exercise 4.13

Exercise 4.24

(a) Under the condition RPM=7500, CPRATIO=13.5, INLETTEMP=1000, EXHTEMP =525, AIRFLOW=10.0, $y$ will appear in the interval $(11599.6, 13665.5)$ with probability 95%.

(b) Under the condition RPM=7500, CPRATIO=13.5, INLETTEMP=1000, EXHTEMP =525, AIRFLOW=10.0, $E(y)$ will appear in the interval $(12157.9, 13107.1)$ with probability 95%.

(c) Yes. The confidence interval only considers variance of $X\beta$, but the prediction interval should consider the sum of two variance of $X\beta$ and residual.

Exercise 4.32

(a) The linear order model is:

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 \beta_5 x_5 + \beta_6 x_2 x_5 + \beta_7 x_3 x_5$$

where RPM is $x_1$, INLETTEMP is $x_2$, EXHTEMP is $x_3$, CPRA-TIO is $x_4$, AIRFLOW is $x_5$.

(b)

$$\hat{y} = 13646 - 0.04560 x_1 - 12.6752 x_2 + 23.0025 x_3 - 3.0227 x_4 + 1.2882 x_5$$
$$+ 0.0162 x_2 x_5 - 0.0414 x_3 x_5$$

(c)
$$H_0 : \beta_6 = 0 \quad vs \quad H_1 : \beta_6 \neq 0$$

The t-statistic for $x_2 x_5$ is 4.40 and p-value is less than 0.0001. So inlet temperature and air flow rate interact is useful to explain heat rate.

(d)
$$H_0 : \beta_7 = 0. \quad vs \quad H_1 : \beta_7 \neq 0$$

The t-statistic for $x_3 x_5$ is $-3.77$ and p-value is less than 0.0004. So exhaust temperature and air flow rate interact is useful to explain heat rate.

(e) part linear relationship between heat rate $y$ and temperature (both inlet and exhaust) depends on air flow rate.

## The SAS System

### The REG Procedure
### Model: MODEL1
### Dependent Variable: HEATRATE

| Number of Observations Read | 67 |
|---|---|
| Number of Observations Used | 67 |

### Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 7 | 158234406 | 22604915 | 138.02 | <.0001 |
| Error | 59 | 9662802 | 163776 | | |
| Corrected Total | 66 | 167897208 | | | |

| Root MSE | 404.69286 | R-Square | 0.9424 |
|---|---|---|---|
| Dependent Mean | 11066 | Adj R-Sq | 0.9356 |
| Coeff Var | 3.65694 | | |

### Parameter Estimates

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 13646 | 1068.17448 | 12.77 | <.0001 |
| RPM | 1 | 0.04599 | 0.01602 | 2.87 | 0.0057 |
| INLETTEMP | 1 | -12.67517 | 1.54155 | -8.22 | <.0001 |
| EXHTEMP | 1 | 23.00252 | 3.76778 | 6.11 | <.0001 |
| CPRATIO | 1 | -3.02265 | 26.41853 | -0.11 | 0.9093 |
| AIRFLOW | 1 | 1.28815 | 3.56266 | 0.36 | 0.7190 |
| IA | 1 | 0.01615 | 0.00367 | 4.40 | <.0001 |
| EA | 1 | -0.04143 | 0.01098 | -3.77 | 0.0004 |

Figure 3: SAS output without intercept for Exercise 4.32

9

Exercise 4.64

(a)
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2$$

where PRM is $x_1$, CPR is $x_2$

(b) $H_0$: $\beta_4 = \beta_5 = 0$ vs $H_1$: at least one are unequal to zero.

(c) Reduced:
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

Complete:
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2$$

(d)
$$SSE_R = 25310639, SSE_C = 19370350, MSE_C = 317547$$

(e) $F = \frac{(SSE_R - SSE_C)/2}{SSE_C /61} = 9.3534$

(f) $F^{-1}(0.9|2, 61) = 2.3917$.

Thus the rejection region was $F > 2.3917$.

(g) the curvature terms in the complete second-order model are useful.

**The SAS System**

**The REG Procedure**
**Model: MODEL1**
**Dependent Variable: HEATRATE**

| Number of Observations Read | 67 |
|---|---|
| Number of Observations Used | 67 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 3 | 142586570 | 47528857 | 118.30 | <.0001 |
| Error | 63 | 25310639 | 401756 | | |
| Corrected Total | 66 | 167897208 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 633.84239 | R-Square | 0.8492 |
| Dependent Mean | 11066 | Adj R-Sq | 0.8421 |
| Coeff Var | 5.72761 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 12065 | 418.52997 | 28.83 | <.0001 |
| RPM | 1 | 0.16969 | 0.03467 | 4.89 | <.0001 |
| CPRATIO | 1 | -146.06557 | 26.65913 | -5.48 | <.0001 |
| RC | 1 | -0.00242 | 0.00312 | -0.78 | 0.4401 |

Figure 4: SAS output without intercept for Exercise 4.64

## The SAS System

### The REG Procedure
### Model: MODEL1
### Dependent Variable: HEATRATE

| Number of Observations Read | 67 |
|---|---|
| Number of Observations Used | 67 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 5 | 148526859 | 29705372 | 93.55 | <.0001 |
| Error | 61 | 19370350 | 317547 | | |
| Corrected Total | 66 | 167897208 | | | |

| Root MSE | 563.51284 | R-Square | 0.8846 |
|---|---|---|---|
| Dependent Mean | 11066 | Adj R-Sq | 0.8752 |
| Coeff Var | 5.09209 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 15583 | 1142.85985 | 13.63 | <.0001 |
| RPM | 1 | 0.07823 | 0.11044 | 0.71 | 0.4814 |
| CPRATIO | 1 | -523.13391 | 103.37571 | -5.06 | <.0001 |
| RPMSQ | 1 | -1.80598E-7 | 0.00000197 | -0.09 | 0.9272 |
| CPRSQ | 1 | 8.84007 | 2.16320 | 4.09 | 0.0001 |
| RC | 1 | 0.00445 | 0.00558 | 0.80 | 0.4282 |

Figure 5: SAS output without intercept for Exercise 4.64

8. Exercise 4.26

   (a) F-value= 226.35 and p-value< 0.001, the overall model is useful.

   (b) t-value is $-3.09$. Under $\alpha = 0.05$, this variable is significant.

   (c) Given $x_2 = 1$, $\hat{y} = 0.044 + 0.269x_1$

   (d) Given $x_2 = 7$, $\hat{y} = 0.308 - 0.673x_1$

   (e) In part c, $y$ is positive related to $x$. when $x$ increase, $y$ will increase. For part $d$, $y$ is negative related to $x$.
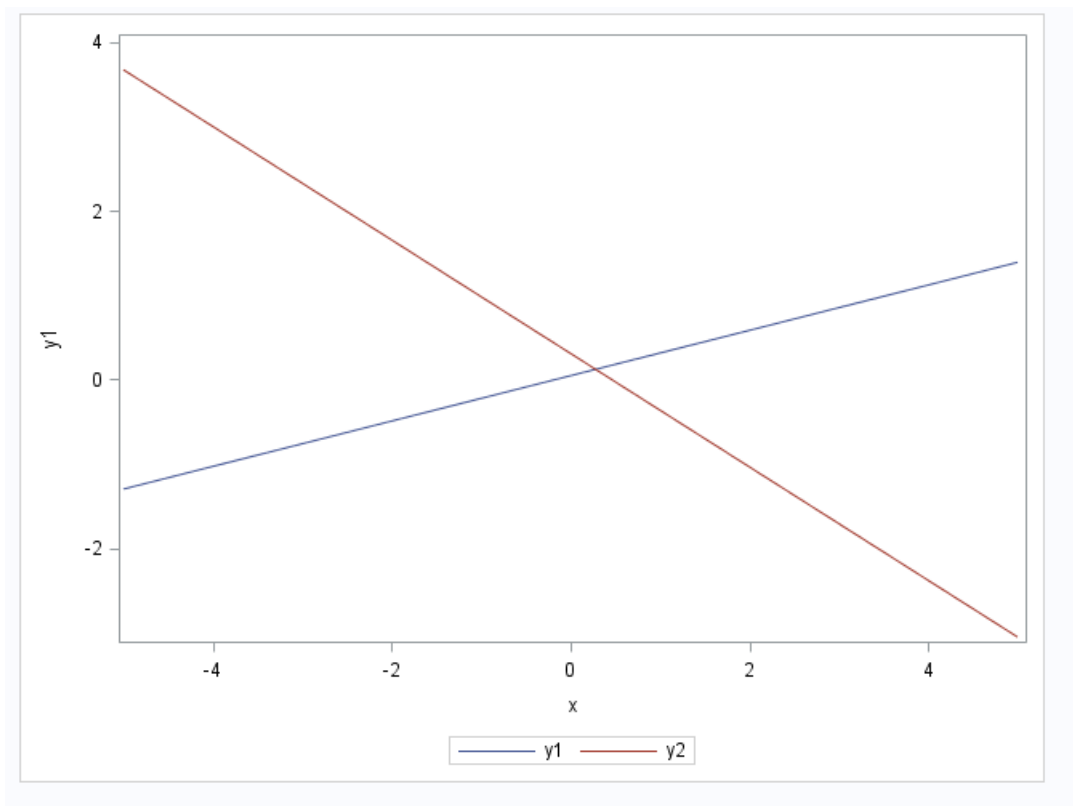


Figure 6: SAS output without intercept for Exercise 4.26

9. Exercise 4.37

   (a) not exact linear relationship.

   (b)
   $$H_0 : \beta_2 = 0.vs.H_1 : \beta_2 \neq 0$$

   t-value is 2.69 and p-value is 0.031.

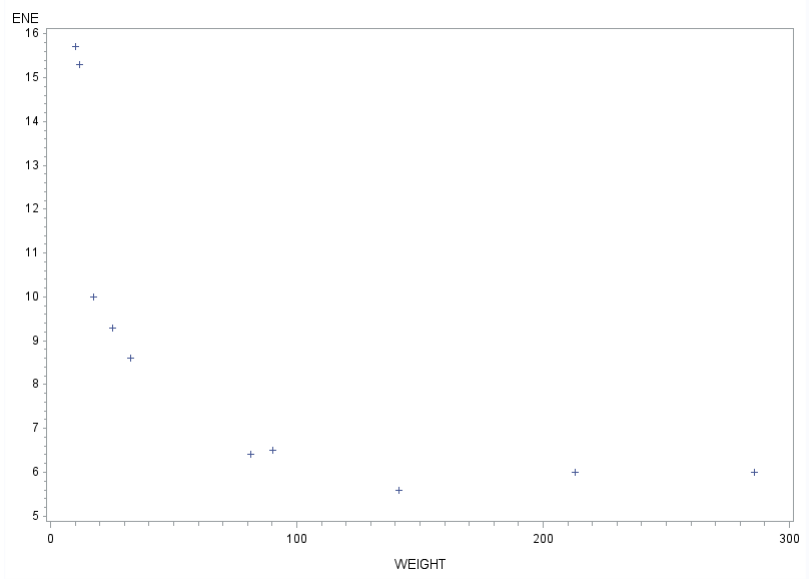   Under $\alpha = 0.10$, the quadratic variable is significant.



Figure 7: Scatter plot for Exercise 4.37

## The SAS System

### The REG Procedure
### Model: MODEL1
### Dependent Variable: ENE

| Number of Observations Read | 10 |
|---|---|
| Number of Observations Used | 10 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 2 | 94.65852 | 47.32926 | 9.83 | 0.0093 |
| Error | 7 | 33.70548 | 4.81507 | | |
| Corrected Total | 9 | 128.36400 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 2.19433 | R-Square | 0.7374 |
| Dependent Mean | 8.94000 | Adj R-Sq | 0.6624 |
| Coeff Var | 24.54504 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 13.71274 | 1.30625 | 10.50 | <.0001 |
| WEIGHT | 1 | -0.10184 | 0.02881 | -3.53 | 0.0095 |
| WEIGHTSQ | 1 | 0.00027348 | 0.00010160 | 2.69 | 0.0310 |

Figure 8: SAS output without intercept for Exercise 4.37

15

10. Exercise 4.59

   (a)
   $$\hat{y} = 80.22 + 156.5x_1 + 272.84x_2 + 760.1x_1x_2 - 42.3x_1^2 + 47x_1^2x_2$$

   (b)
   $$H_0 : \beta_1 = \ldots = \beta_5 = 0 \quad vs \quad H_1 : \exists i, \beta_i \neq 0$$

   F-value is 417.05, and p-value $< 0.0001$. So the overall model is useful.

   (c) There is no enough evidence to indicate that $y$ is curvilinearly related to $x_1$. We should compare with a reduced model:
   $$E(y) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_1x_2$$

   Exercise 4.69

   (a) In the null model, we consider the curvilinearly relationship in two ways: 1. $x_1$ is curvilinearly related to $y$ directly; 2 $x_1$ is curvilinearly related to $y$ based on $x_2$.

   (b) $H_0$: $\beta_4 = \beta_5 = 0$ vs $H_1$: at least one is not 0.

   (c) The curvilinearly variables of $x_1$ are useless.

   (d) $F = \frac{(SSE_R - SSE_C)/2}{SSE_C/(n-6)} = \frac{(89171 - 88819)/2}{88819/30} = 0.0594$

11. Exercise 4.88

   (a)
   $$\hat{y} = 10.6590 - 0.28161x_1 + 0.00267x_1^2$$

   (b) $R_a^2 = 0.8770$(1 point ) the percentage of variation of sample with penalty of degree can be explained by the model.

   (c) $s = 4.5486$, estimate $\sigma$, reflect the variation of y.

   (d)
   $$H_0 : \beta_1 = \beta_2 = 0.vs.H_1 : \exists i, \beta_i \neq 0$$

   F-value is 33.08 , and P -value is 0.0003. So the overall model is useful.

   (e)
   $$H_0 : \beta_2 = 0.vs.H_1 : \beta_2 \neq 0$$

   t-value is 2.13 and P-value is $0.0706 > 0.05$. The evidence is not enough to conclude that the percentage improvement y increase more quickly for more costly fleet modifications than for less costly fleet modifications.

(f)
$$H_0 : \beta_3 = \ldots = \beta_5 = 0.vs.H_1 : \exists 3 \le i \le 5, \beta_i \ne 0$$

$$F = \frac{(SSE_R - SSE_C)/3}{\text{SSE}_C/(n-6)} = 0.3301, F^{-1}(0.95|3,4) = 6.5914$$

The type of base $x_2$ is useless.

**The SAS System**

**The REG Procedure**
**Model: MODEL1**
**Dependent Variable: PERCENT**

| Number of Observations Read | 10 |
|---|---|
| Number of Observations Used | 10 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 2 | 1368.77501 | 684.38750 | 33.08 | 0.0003 |
| Error | 7 | 144.82499 | 20.68928 | | |
| Corrected Total | 9 | 1513.60000 | | | |

| Root MSE | 4.54855 | R-Square | 0.9043 |
|---|---|---|---|
| Dependent Mean | 17.20000 | Adj R-Sq | 0.8770 |
| Coeff Var | 26.44504 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 10.65904 | 14.55009 | 0.73 | 0.4876 |
| COST | 1 | -0.28161 | 0.28088 | -1.00 | 0.3494 |
| COSTSQ | 1 | 0.00267 | 0.00125 | 2.13 | 0.0706 |

Figure 9: SAS output without intercept for Exercise 4.88

## The SAS System

### The REG Procedure
### Model: MODEL1
### Dependent Variable: PERCENT

| Number of Observations Read | 10 |
|---|---|
| Number of Observations Used | 10 |

| | | Analysis of Variance | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 5 | 1397.51481 | 279.50296 | 9.63 | 0.0238 |
| Error | 4 | 116.08519 | 29.02130 | | |
| Corrected Total | 9 | 1513.60000 | | | |

| Root MSE | 5.38714 | R-Square | 0.9233 |
|---|---|---|---|
| Dependent Mean | 17.20000 | Adj R-Sq | 0.8274 |
| Coeff Var | 31.32059 | | |

| | | Parameter Estimates | | | |
|---|---|---|---|---|---|
| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| |
| Intercept | 1 | 2.03190 | 22.64502 | 0.09 | 0.9328 |
| COST | 1 | -0.10364 | 0.48513 | -0.21 | 0.8413 |
| BASE | 1 | 49.76686 | 52.36930 | 0.95 | 0.3958 |
| COSTSQ | 1 | 0.00189 | 0.00234 | 0.81 | 0.4643 |
| CB | 1 | -0.87476 | 0.93575 | -0.93 | 0.4028 |
| COSTSQBASE | 1 | 0.00353 | 0.00398 | 0.89 | 0.4246 |

Figure 10: SAS output without intercept for Exercise 4.88