4.26 in book

1.

**Table 4E.6** Tensile Strength Data for Exercise 4.26

| Strength | Percentage Hardwood | Strength | Percentage Hardwood |
|----------|--------------------|----------|--------------------|
| 160 | 10 | 181 | 20 |
| 171 | 15 | 188 | 25 |
| 175 | 15 | 193 | 25 |
| 182 | 20 | 195 | 28 |
| 184 | 20 | 200 | 30 |

$x$ = strength ; $y$ = % Hardwood

a) $n = 10$

$\sum x_i = 160 + 171 + \ldots + 200 = 1829$

$\sum x_i^2 = 160^2 + 171^2 + \ldots + 200^2 = 335825$

$\sum y_i = 10 + 15 + \ldots + 30 = 208$

$\sum y_i^2 = 10^2 + 15^2 + \ldots + 30^2 = 4684$

$\sum x_i y_i = 160(10) + 171(15) + \ldots + 200(30) = 38715$

$S_{xy} = \sum x_i y_i - \dfrac{(\sum y_i)(\sum x_i)}{n} = 38715 - \dfrac{(208)(1829)}{10} = 671.8$

$S_{xx} = \sum x_i^2 - \dfrac{(\sum x_i)^2}{n} = 335825 - \dfrac{(1829)^2}{10} = 1300.9$

$\tilde{\beta}_1 = \dfrac{S_{xy}}{S_{xx}} = \dfrac{671.8}{1300.9} = 0.5164$ ;

$\bar{x} = \dfrac{1}{10}(\sum x_i) = \dfrac{1}{10}(1829) = 182.9$

$\bar{y} = \dfrac{1}{10}(\sum y_i) = \dfrac{1}{10}(208) = 20.8$

$\beta_0 = \bar{y} - \tilde{\beta}_1 \cdot \bar{x} = 20.8 - 0.5164(182.9) = -73.65$

$\boxed{\hat{y} = 0.516x - 73.7}$ where $\begin{cases} x = \text{Strength} \\ \hat{y} = \text{Percentage Hardwood} \end{cases}$

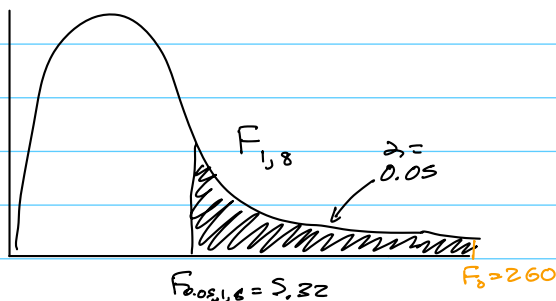b) ANOVA for linear regression

$H_0: \beta_1 = 0$

$H_a: \beta_1 \neq 0$

Test Statistic:

$SS_R = \sum(\hat{y}_i - \bar{y})^2 = (8.97 - 20.8)^2 + (14.65 - 20.8)^2 + \ldots + (29.63 - 20.8)^2 = 346.93$

$SS_E = \sum(y_i - \hat{y}_i)^2 = (10 - 8.97)^2 + \ldots + (30 - 29.63)^2 = 10.67$

$F_0 = \dfrac{SS_R/1}{SS_E/(n-2)} = \dfrac{346.43}{10.67/(8)} = 260$

Reject if: $F_0 > F_{0.05,1,8}$



$F_{1,8}$   $\alpha = 0.05$

$F_{0.05,1,8} = 5.32$   $F_0 = 260$

$260 > 5.32 \rightarrow$ Reject $H_0$

We have sufficient evidence to support the claim that there is a significant linear relationship between percentage hardwood ($y$) and strength ($x$).

c)

**Regression Equation**

Percentage Hardwood = -73.65 + 0.5164 Strength

**Coefficients**

| Term | Coef | SE Coef | 95% CI | T-Value | P-Value | VIF |
|------|------|---------|--------|---------|---------|-----|
| Constant | -73.65 | 5.87 | (-87.19, -60.12) | -12.55 | 0.000 | |
| Strength | 0.5164 | 0.0320 | (0.4426, 0.5903) | 16.12 | 0.000 | 1.00 |

**Model Summary**

| S | R-sq | R-sq(adj) | PRESS | R-sq(pred) | AICc | BIC |
|---|------|-----------|-------|------------|------|-----|
| 1.15513 | 97.01% | 96.64% | 17.4802 | 95.11% | 39.03 | 35.94 |

**Analysis of Variance**

| Source | DF | Seq SS | Contribution | Adj SS | Adj MS | F-Value | P-Value |
|--------|-----|--------|--------------|--------|--------|---------|---------|
| Regression | 1 | 346.93 | 97.01% | 346.93 | 346.925 | 260.00 | 0.000 |
| Strength | 1 | 346.93 | 97.01% | 346.93 | 346.925 | 260.00 | 0.000 |
| Error | 8 | 10.67 | 2.99% | 10.67 | 1.334 | | |
| Total | 9 | 357.60 | 100.00% | | | | |

95% CI on $\beta_1$ (slope)

(0.4426, 0.5903)

In repeated sampling, the interval calculated captures the the slope $\beta_1$ 95% of the time.

Also calculations from previous page (excel assisted) match the minitab output

d) $R^2 = \dfrac{SS_R}{SS_T} = \dfrac{SS_R}{SS_R + SS_E} = \dfrac{346.93}{346.93 + 10.67} = 97.0\%$

97% of the variability in the target variable is accounted for by the regression model.

2.

Data from
problem ①

**Table 4E.6** Tensile Strength Data for Exercise 4.26

| Strength | Percentage Hardwood | Strength | Percentage Hardwood |
|----------|---------------------|----------|---------------------|
| 160 | 10 | 181 | 20 |
| 171 | 15 | 188 | 25 |
| 175 | 15 | 193 | 25 |
| 182 | 20 | 195 | 28 |
| 184 | 20 | 200 | 30 |

$x =$ strength; $y = \%$ Hardwood

Residuals

$e_1 = 10 - 8.97 = 1.03$

$e_2 = 15 - 14.65 = .345$

$\vdots$

$e_{10} = 30 - 29.68 = 0.869$

a) Residual Normal Probability Plot



Probability Plot of Q1 residuals
Normal - 95% CI

Mean 2.842171E-15
StDev 1.089
N 10
AD 0.249
P-Value 0.665

Insufficient evidence to reject the claim that the data is not normal. As $p$-value $= 0.665 > 0.05 = \alpha$, we can support the normality assumption of the model.
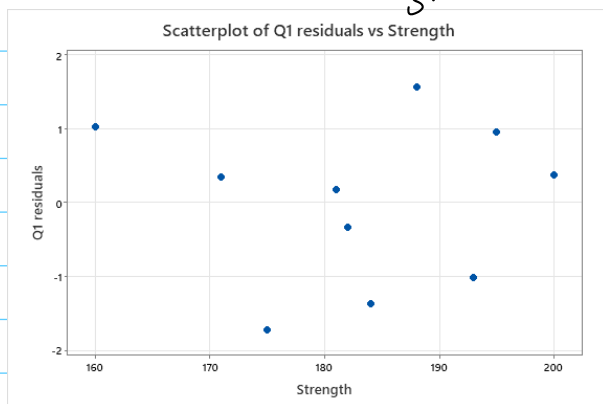
b) Residuals vs. Fitted Value



Versus Fits
(response is Percentage Hardwood)

The residuals are evenly distributed about $y = 0$. There is no clear pattern, and thus the equal variance assumption is valid here.

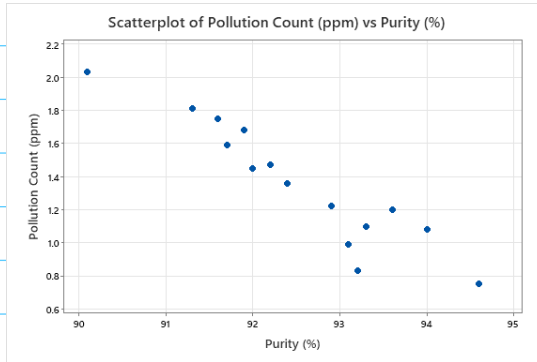If there was a pattern, I would apply a transformation

c) Residuals vs. Strength values (x)



Scatterplot of Q1 residuals vs Strength

The residuals show no patterns and are evenly distributed about 0. The equal variance assumption is again valid.

**3.**

| Purity (%) | 93.3 | 92.0 | 92.4 | 91.7 | 94.0 | 94.6 | 93.6 |
|---|---|---|---|---|---|---|---|
| Pollution count (ppm) | 1.10 | 1.45 | 1.36 | 1.59 | 1.08 | 0.75 | 1.20 |
| Purity (%) | 93.1 | 93.2 | 92.9 | 92.2 | 91.3 | 90.1 | 91.6 | 91.9 |
| Pollution count (ppm) | 0.99 | 0.83 | 1.22 | 1.47 | 1.81 | 2.03 | 1.75 | 1.68 |

a)



Scatterplot of Pollution Count (ppm) vs Purity (%)

Without a fitted line there is a clear negative linear relationship between Pollution Count (y) and Purity Percent (x). I'm convinced!

Regression:

$\hat{y} = 24.23 - .303x$

$\begin{cases} \hat{y} = \text{Pollution Count (ppm)} \\ x = \text{Purity Percent} \end{cases}$

b) ANOVA test for linear regression

$H_0: \beta_1 = 0$
$H_a: \beta_1 \neq 0$
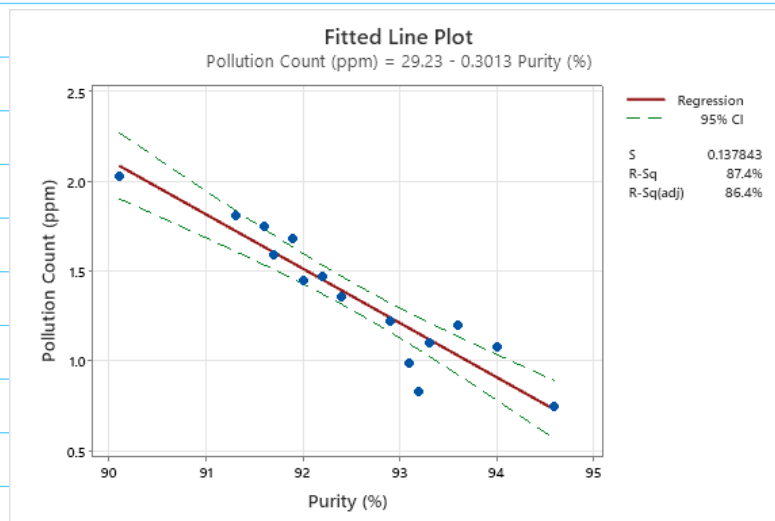
$\}$ Minitab p-value $\approx 0 < \alpha = 0.05$

p-value $< \alpha$, so we have significant evidence to reject the claim that there is no linear relationship. We can accept there is a significant linear relationship between Pollution Count (y) and Purity Percent (x).

c) 95% CI for: Purity % = $(-0.3698, -.2827)$ ← slope
Pollution Count (Constant) = $(22.89, 35.57)$ ← intercept

In repeated samples, the intervals calculated capture the true slope or intercept 95% of the time.

d)



Fitted Line Plot
Pollution Count (ppm) = 29.23 - 0.3013 Purity (%)

Regression
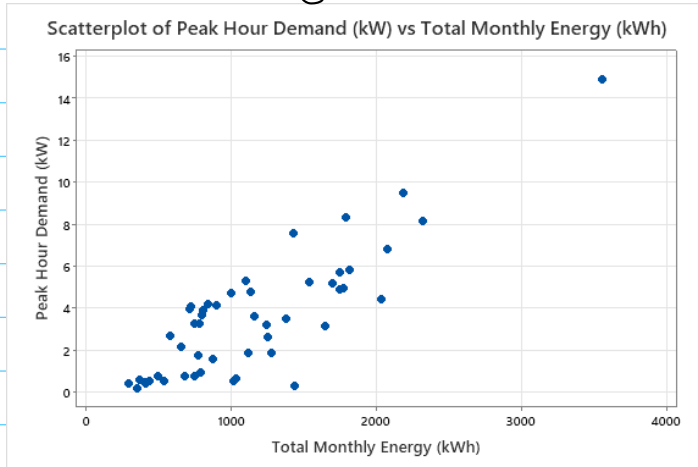95% CI

S        0.137843
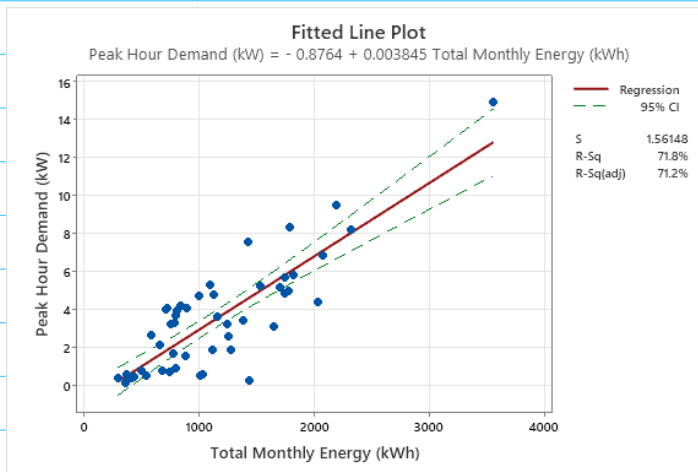R-Sq     87.4%
R-Sq(adj) 86.4%

In repeated experiments there is a 95% chance that the CI on the responses calculated contain the true responses.

4. use $\begin{cases} x = \text{Total Monthly Energy Used (kWh)} \\ y = \text{Peak Hour Demand (KW)} \end{cases}$

a) Scatter plot of $y$ vs $x$

Scatterplot of Peak Hour Demand (kW) vs Total Monthly Energy (kWh)

b) Linear Fit

Fitted Line Plot
Peak Hour Demand (kW) = - 0.8764 + 0.003845 Total Monthly Energy (kWh)

| S | 1.56148 |
| R-Sq | 71.8% |
| R-Sq(adj) | 71.2% |

$\hat{y} = -0.8764 + .003845x$

$\begin{cases} x = \text{Total Monthly Energy Used (kWh)} \\ y = \text{Peak Hour Demand (KW)} \end{cases}$

c) $H_0: \beta_1 = 0$ $\qquad\qquad$ $H_a: \beta_1 \neq 0$
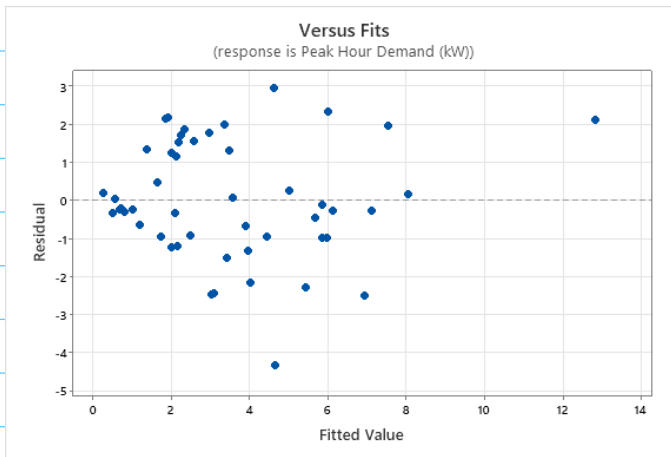
## Analysis of Variance

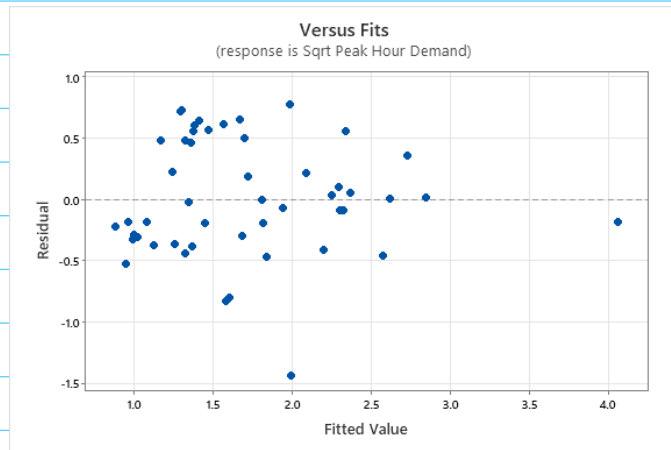| Source | DF | Seq SS | Contribution | Adj SS | Adj MS | F-Value | P-Value |
|---|---|---|---|---|---|---|---|
| Regression | 1 | 297.7 | 71.78% | 297.7 | 297.728 | 122.11 | 0.000 |
| Total Monthly Energy (kWh) | 1 | 297.7 | 71.78% | 297.7 | 297.728 | 122.11 | 0.000 |
| Error | 48 | 117.0 | 28.22% | 117.0 | 2.438 | | |
| Total | 49 | 414.8 | 100.00% | | | | |

p-value $\approx 0 < \alpha = 0.05$

p-value $< \alpha$, so we have significant evidence to reject the claim that there is no linear relationship. We can accept there is a significant linear relationship between the total Monthly Energy Used $(x)$, and the Peak Hour Demand $(y)$.

d)

**Versus Fits**
(response is Peak Hour Demand (kW))



There seems to be a feather like pattern, despite being roughly displaced about 0. At lower fitted values, residuals are closer together. Thus, the equality of variance may not be a good assumption and we should apply a data transformation to remedy this.

e)

**Versus Fits**
(response is Sqrt Peak Hour Demand)



Yes, the data is more evenly distributed about zero and there is not a pattern corresponding to the fitted value. With this transformation, the equality of variance assumption appears to be a valid assumption.

5. a) Problem 4.30

(i) ==Satisfaction $(y)$ = 131.10 − 1.290 · Age $(x_1)$== $\begin{cases} y = \text{Satisfaction} \\ x_1 = \text{Age} \end{cases}$
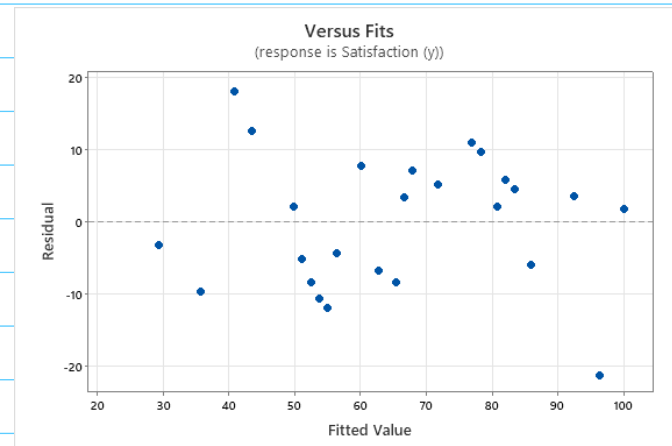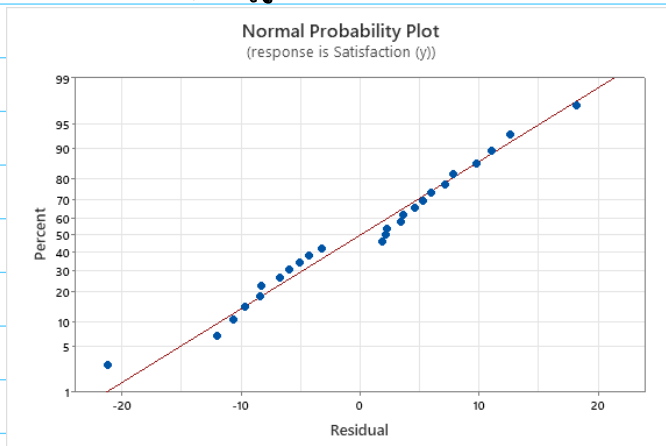
(ii) $H_0: \beta_1 = 0$      $H_a: \beta_1 \neq 0$

Minitab → p-value = 0 < α = 0.05

we have significant evidence to reject the claim that there is no linear relationship. We can accept there is a significant linear relationship between the Age of the patient $(x_1)$, and their Satisfaction $(y)$.

(iii) 81.24% of the variability is accounted for by the regressor variable age.

b) Problem 4.31



The normal probability plot of the residuals is linear with no outstanding variability, so the normality assumption is valid. The residuals vs. Fitted value is evenly distributed about ⓪ and no clear pattern is seen. Thus the equality of variance assumption is valid. The model seems to be a good fit.

c) Problem 4.32

(:) Satisfaction $(y) = 143.47 - 1.031 \cdot Age(x_1) - 0.556 \cdot Severity(x_2)$

$\left\{ \begin{array}{l} y = Satisfaction \qquad x_1 = Age \\ \qquad\qquad\qquad x_2 = Severity \end{array} \right.$

$H_0: B_1 = B_2 = 0 \qquad H_a: B_i \neq 0$ for at least one $i$

$Minitab \rightarrow p\text{-value} = 0 < \alpha = 0.05$

we have significant evidence to reject the claim that there is no linear relationship. We can accept there is a significant linear relationship between the Age of the patient $(x_1)$, Severity of their illness $(x_2)$, and their satisfaction $(y)$.
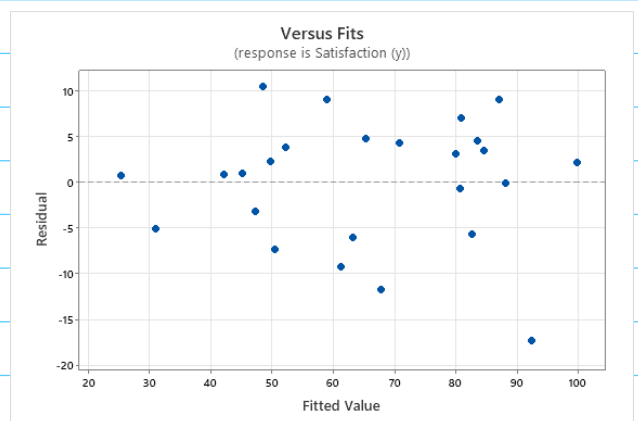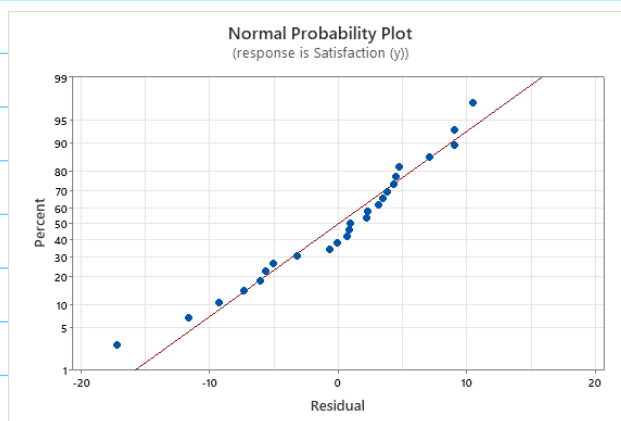
(:) $R^2(x_1) = 81.24\%$ ← Age
$R^2(x_2) = 60.42\%$ ← Severity

The portion of variability is significant $(> \alpha = 5\%)$ in both cases. Both regressors are needed.

(:::) $R^2_{mult.} = 86.72\%$ (adj)
$R^2_{single} = 80.43\%$ (adj)

The adjusted $R^2$ value, which accounts for overfitting effects, improved with the full fit expression. Thus we can conclude adding severity improved the model's quality.

d) Problem 4.33



The normal probability plot of the residuals is linear with no outstanding variability, so the normality assumption is valid. The residuals vs. Fitted value is evenly distributed about ① and no clear pattern is seen. Thus, the equality of variance assumption is valid. The model seems to be a good fit.