LAST NAME:_____FIRST NAME:_____

STUDENT NUMBER:_____

ENROLLED IN: (circle one)          STA 302          STA 1001

INSTRUCTIONS:
• Time: 90 minutes
• Aids allowed: calculator.
• A table of values from the $t$ distribution is on the last page (page 8).
• Total points: 50

---

**Some formulae:**

$$b_1 = \frac{\sum(X_i - \overline{X})(Y_i - \overline{Y})}{\sum(X_i - \overline{X})^2} = \frac{\sum X_i Y_i - n\overline{XY}}{\sum X_i^2 - n\overline{X}^2} \qquad b_0 = \overline{Y} - b_1\overline{X}$$

$$\text{Var}(b_1) = \frac{\sigma^2}{\sum(X_i - \overline{X})^2} \qquad \text{Var}(b_0) = \sigma^2\left(\frac{1}{n} + \frac{\overline{X}^2}{\sum(X_i - \overline{X})^2}\right)$$

$$\text{Cov}(b_0, b_1) = -\frac{\sigma^2\overline{X}}{\sum(X_i - \overline{X})^2} \qquad \text{SSTO} = \sum(Y_i - \overline{Y})^2$$

$$\text{SSE} = \sum(Y_i - \hat{Y}_i)^2 \qquad \text{SSR} = b_1^2\sum(X_i - \overline{X})^2 = \sum(\hat{Y}_i - \overline{Y})^2$$

$$\sigma^2\{\hat{Y}_h\} = \text{Var}(\hat{Y}_h) = \sigma^2\left(\frac{1}{n} + \frac{(X_h - \overline{X})^2}{\sum(X_i - \overline{X})^2}\right) \quad \sigma^2\{\text{pred}\} = \text{Var}(Y_h - \hat{Y}_h) = \sigma^2\left(1 + \frac{1}{n} + \frac{(X_h - \overline{X})^2}{\sum(X_i - \overline{X})^2}\right)$$

$$r = \frac{\sum(X_i - \overline{X})(Y_i - \overline{Y})}{\sqrt{\sum(X_i - \overline{X})^2\sum(Y_i - \overline{Y})^2}} \qquad \text{Working-Hotelling coefficient: } W = \sqrt{2\,F_{2,n-2;1-\alpha}}$$

---

| 1 | 2 | 3 | 4abc | 4def | 4ghi |
|---|---|---|------|------|------|
|   |   |   |      |      |      |

1. (10 points) A simple linear regression model is fit on $n$ observed data points.

   (a) What is the difference between $\beta_1$ and $b_1$?

   (b) What does it mean if $R^2 = 1$?

   (c) In lecture we showed $\sum_{i=1}^{n} e_i = 0$ and $\sum_{i=1}^{n} e_i X_i = 0$. Show that $\sum_{i=1}^{n} e_i \hat{Y}_i = 0$. (You may use the results shown in class if they are helpful.)

   (d) Explain why the result in (c) implies that the residuals and predicted values are uncorrelated and why this is useful.

2

2. (8 points) In order to carry out linear regression analyses, in addition to the assumption that a linear model is appropriate for the data, we have made the following assumptions:

- the expectation of the random errors is zero
- the variance of the errors is constant
- the errors are uncorrelated
- the errors are normally distributed

Assume that the independent variable is not random.

(a) Which of these additional assumptions are necessary to show that $b_1$ is unbiased for $\beta_1$?

(b) Derive the formula for the variance of $b_1$ and state which of the additional assumptions are necessary for the derivation.

3. (7 points) In lecture we have considered the Snow Gauge example. In this experiment, scientists measured the number of gamma rays (the `gain`) that make it through 10 samples of each of 9 densities of polystyrene. We fit a simple regression model with the logarithm of gain (`loggain`) as the dependent variable and `density` as the independent variable to these 90 points. A scientist argues that, since 10 samples were measured at each density, taking the mean of `loggain` at that density will result in a better estimate and the regression should then be run using the 9 resulting points. Will the least squares estimates of the slope and intercept change? Will the estimate of the error variance change? If there is a change, say whether it is larger or smaller. Justify your answers.

4. (25 points) The data analysed in this question are from a random sample of records of esales of homes in 1993 in the U.S. city of Albuquerque. The data collected include many variables about the homes sold, but we will only consider how well the size of the home (in square feet of usable floor space, variable name: `sqft`) can be used to predict the selling price (in hundreds of dollars, variable name: `price`) of the home.

Some output from SAS is given below. Note that some numbers have been replaced by letters.

```
                       Descriptive Statistics
                                        Uncorrected                    Standard
Variable            Sum            Mean           SS      Variance     Deviation
Intercept     116.00000        1.00000    116.00000             0             0
sqft             189751     1635.78448    337777165        238134     487.98988
price            123045     1060.73276    147252397        145518     381.46781

                       Analysis of Variance
                                   Sum of          Mean
Source                    DF      Squares        Square     F Value     Pr > F
Model                    (A)     13229494           (B)      430.28        (C)
Error                    114          (D)         30746
Corrected Total          (E)     16734535

              Root MSE                 (F)     R-Square      0.7906
              Dependent Mean     1060.73276     Adj R-Sq      0.7887
              Coeff Var            16.53058

                       Parameter Estimates
                       Parameter      Standard
    Variable     DF     Estimate         Error     t Value     Pr > |t|
    Intercept     1     -76.20835      57.17689       -1.33       0.1852
    sqft          1       0.69504       0.03351         (G)       <.0001
```

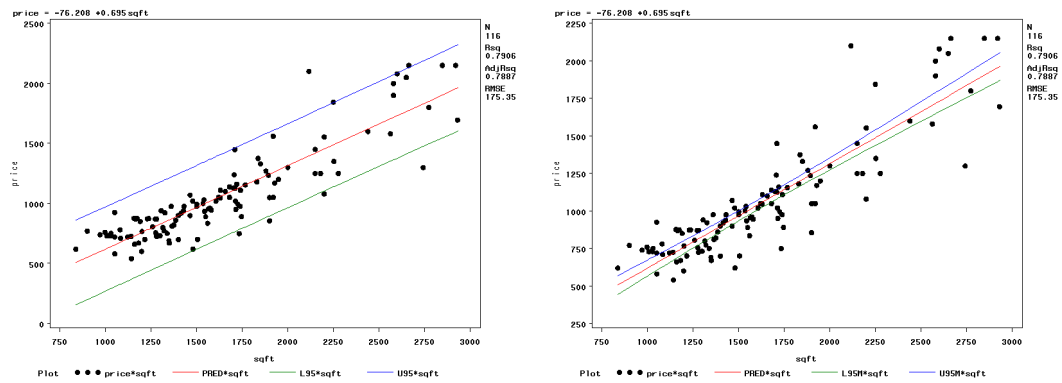(a) Find the 7 missing values (A through G) in the SAS output.

(b) How many houses are in the sample?

(c) Is the intercept statistically significantly different from 0? Justify your answer. Explain the meaning of the intercept for a real estate agent.

(d) A house with 2000 square feet of usable space came on the market (under the same market conditions as the houses used in this analysis). Predict its selling price.

(e) What is the standard error of the prediction in part (d)?
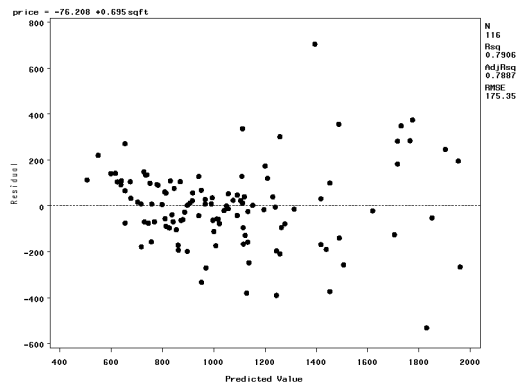
(f) Plots of the data including the regression line and 95% confidence intervals for the mean of Y and 95% prediction intervals for Y are given below.



i. Which plot is which? How do you know?

ii. For the plot on the right, show how to calculate the value on the lowest curve corresponding to $X = 1500$. In your answer include the numeric value.

(g) A plot of the residuals versus predicted values is below.



price = -76.208 +0.695 sqft

N
116
Rsq
0.7906
AdjRsq
0.7887
RMSE
175.35

Describe any problems you see in the residual plot. If the plot shows that any assumptions are being violated indicate which.

(h) A student hired by the real estate board to analyse these data argues that we should consider correlation rather than regression since the predictor variable is random. Respond to this comment.

(i) Several of the homes in the random sample used in this analysis were from a new housing development. Why should this be considered in carrying out the analyses?