

A N S W E R S

STAT 2112: Mock Final

1. In an American metropolitan area, the list of the known inhabitants indicates that out of 700,000 people, 46,000 of them are 17 to 26 years old, 400,000 of them are 27 to 65 years old, 54,000 of them are 66 years old and older, the rest are too young to attend college. The demographic unit of the local municipality wants to report to the common council the number of people who are attending school in a higher education institution like a technical college or a university. The three strata which were defined above are sufficient for the purposes of the common council to allocate money to the local technical college. The precision which the chairman of the subcommittee of the common council has requested is $CV = 3\%$ for stratum I, $CV = 10\%$ for stratum II, and $CV = 20\%$ for the stratum III.

Based on the Census Bureau's American Community Survey for this metropolitan area, about 5,000 with a CV of 3% or a mean of $\frac{5000}{46000} = .108695$ with a standard deviation $\sqrt{\frac{5000(.03)}{46000}} = .05710$ of those who are 17 to 26 years old are attending or have attended college; 4,000 with a CV of 6% or a mean of $\frac{4000}{400000} = 1\%$ with a standard deviation of $\sqrt{\frac{4000(.06)}{400000}} = .02449$ of those who are in the 27 to 65 years old stratum of the population and who have attended college since age 27 ; and of the retired group, 150 with a CV of 10% or a mean of $\frac{150}{54000} = .27\%$ with a standard deviation of $\sqrt{\frac{150(.10)}{54000}} = .0166$ are attending college or have taken a college level courses since age 66.

From other surveys which this demographic unit has recently conducted, about 60% respond in the 17 to 26 years old group, about 40% respond in the 27 to 65 years old group, and about 80% respond in the 66 years old and older group.

- (a) Calculate the sampling sizes for each of the three strata.

ANSWER:

- i. For Stratum I

$$n_I = 505$$

- ii. For Stratum II

$$n_{II} = 1495$$

- iii. For Stratum III

$$n_{III} = 1107$$

- (b) In addition to discovering the proportion of people in each stratum who are attending or have attended college, the enumerator asked for the the respondent's monthly rent or mortgage payment or if no rent nor mortgage payment then the monthly real estate tax. The following descriptive statistics were obtained.

	Stratum I	Stratum II	Stratum III
n	500	1500	1110
Monthly Rent (\bar{x})	750	1750	500
$\frac{s}{\sqrt{n}}$	30	105	50
1 st quartile	690	1405	400
2 nd quartile	700	1800	510
3 rd quartile	910	2100	650
maximum	1500	4200	5000
minimum	0	150	250

i. Construct a box plot for each stratum.

ANSWER: See Figure 1

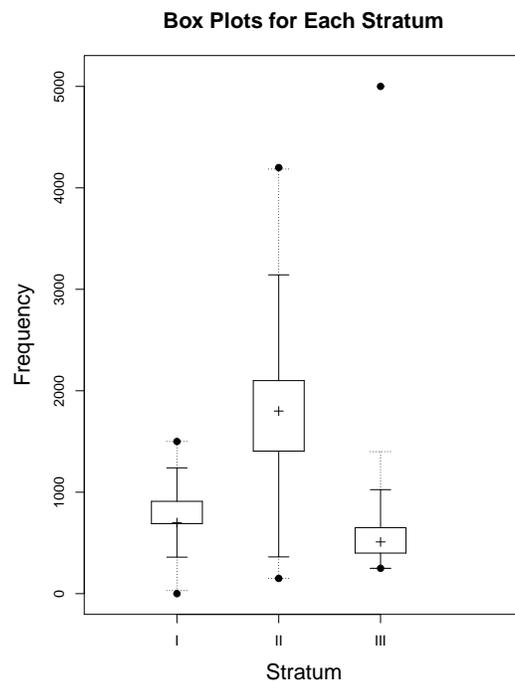


Figure 1:

ii. Calculate the 95% confidence interval of the average monthly rent for each stratum.

ANSWER:

A. For Stratum I:

The 95% confidence interval for stratum I is:(691,809)

B. For Stratum II:

The 95% confidence interval for stratum II is:(1544,1956)

C. For Stratum III:

The 95% confidence interval for stratum III is:(401,599)

iii. Test at $\alpha = .05$ the hypothesis that

$$H_0 : \mu_I = \mu_{III} \text{ vs } H_1 : \mu_I \neq \mu_{III}$$

ANSWER:

A. $\alpha = .05 \frac{\alpha}{2} = .025 \nu = n_1 + n_3 - 2 = 500 + 1110 - 2 = 1608$

B. $t_{1608,.025} = 1.96144$

C. Note that: $var(\bar{x}_I) = s_I^2 = \left(\frac{s_I}{\sqrt{n_I}}\right)^2 n_I = (30)^2 500 = 450000$ and $var(\bar{x}_{III}) = s_{III}^2 = \left(\frac{s_{III}}{\sqrt{n_{III}}}\right)^2 n_{III} = (50)^2 1110 = 2775000$

$$s_p = 1433.003$$

where $\nu = n_I + n_{III} - 2 = 1608$

$$T = \frac{\bar{x}_I - \bar{x}_{III}}{\sqrt{s_p^2 \sqrt{\frac{1}{n_I} + \frac{1}{n_{III}}}}} = 3.239115$$

D. Is $\|T\| = 3.239115 > 1.96144$? Yes

E. Therefore, we reject the null hypothesis.

iv. What assumptions need to be made about the survey?

A. Is the sample representative of the population?

B. Are the strata disjoint?

C. Did the respondents answer the questions independently of other respondents and of the enumerator?

D. Can it be assumed that the respondents treat the questions identically?

E. Does the concept of the survey and the questions make sense?

F. Is there normality in the data?

G. Does the Normal distribution fit the scales which are used in the questions?

2. Design of Experiment

(a) Write a 2×2 factorial design of experiment of your conception.

i. List the two factors and the two levels.

ANSWER: Use five tablespoons of honey or ten tablespoons of honey in making an apple pie. Use five tablespoons of brown sugar or ten tablespoons of brown sugar in making an apple pie.

ii. Explain your concept of the theory which underlies your proposed experiment and what you hope to find from the experiment.

ANSWER: Supposedly, the honey tastes better than sugar. Three judges will sample each of four apple pies. On a scale of 1 bland to 5 tasty, the testers will judge the pies.

iii. What assumptions need to be made about the survey?

ANSWER: We assume that the judges' evaluations are reproducible and identical, that the scale is discriminating enough, and that the recipe for the apple pies is closely followed, so that, all things being equal, the difference in taste is due to the honey and brown sugar.

iv. Be original and keep the design of the experiment simple.

ANSWER: I will use the factors and levels described above three times to conduct a 2×2 factorial design experiment replicated three times.

3. A pediatrician claims that the weight of the second child is more than the weight of a woman's first child regardless of the sex of the child at birth. He asserts that there exists a linear relationship between the weights of the first and second child. The following table contains thirteen observations in which the weight of the first child is the explanatory variable and the weight of the second child is the response variable.

Child	1	2	3	4	5	6	7	8	9	10	11	12	13
First	5.13	5.25	6.71	4.71	3.77	5.81	8.29	6.36	6.08	4.91	3.19	5.64	6.37
Second	6.43	5.83	7.78	5.27	4.45	6.61	9.51	7.01	8.49	5.62	4.04	6.10	7.25

$\sum_{i=1}^{13} x_i = 72.22$ $\sum_{i=1}^{13} x_i^2 = 421.8674$ $\sum_{i=1}^{13} x_i y_i = 491.9085$ $\sum_{i=1}^{13} y_i = 84.39$. After the line was fitted to the data, the sum of squared errors was computed to be: $SSE=2.839573$.

(a) Write the linear model: ANSWER: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ where $\epsilon_i \sim N(0, \sigma^2)$

(b) What are the two principal reasons for using linear models:

ANSWER:

i. To Interpolate

ii. To show an inherent relationship between x and y.

Suppose that the weight of your first child was 8 lbs 7 oz. What is the predicted weight of your second child? To that end, compute:

(c) $SS_{xx} =$

ANSWER: $SS_{xx} = 20.65752$

(d) $SS_{xy} =$

ANSWER: $SS_{xy} = 23.08959$

(e) $\hat{\beta}_0 =$

ANSWER: $\hat{\beta}_0 = 0.2821025$

(f) $\hat{\beta}_1 =$

ANSWER:

$\hat{\beta}_1 = 1.1177329$

(g) When $x_0=8.4375$, find $\hat{y} =$

ANSWER: $\hat{y} = 9.712974$

(h) The 95% for β_1 is: (0.87, 1.36). Is the model valid?

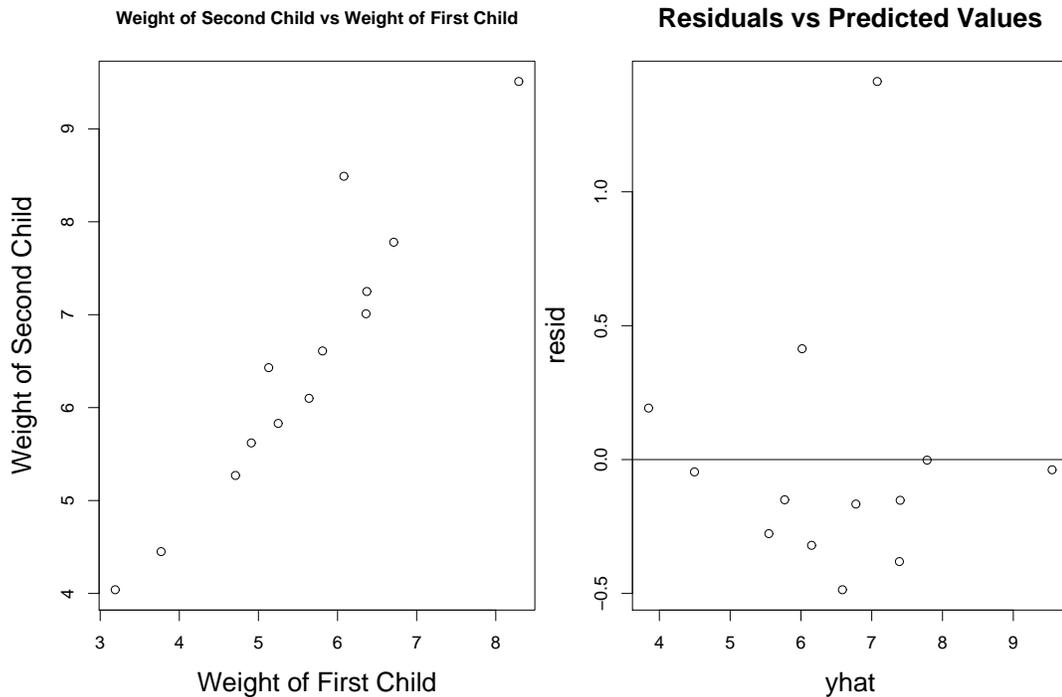
ANSWER:

- i. The theory seems reasonable since it is based on the pediatrician's many observations.
- ii. The plot of the data show a linear relationship between x and y.
- iii. The 95% CI about β_1 shows that $H_0 : \beta_1 = 0$ can be rejected in favor of the alternative hypothesis, $H_1 : \beta_1 \neq 0$, because $0 \notin (0.87, 1.36)$. We note based on the ANOVA table given below in Table 1 that $H_0 : \beta_0 = 0$ can be rejected.
- iv. The plot of residuals versus predicted values shows no pattern, but in the QQ plot shown in Figure 2 there appears to be a departure from normality.
All four criteria are answered in the affirmative; therefore, the model seems to be a valid one, though with some reservations about the three points which depart from normality. They have to be investigated.

(i) What is the 95% confidence interval of $E[\hat{y}]$, when $x_0=8.4375$. ANSWER: The 95% CI about y_0 is: (8.35,11.07)

Child	1	2	3	4	5	6	7	8	9	10	11	12	13
Predicted	6.01	6.15	7.78	5.54	4.49	6.77	9.54	7.39	7.07	5.77	3.84	6.58	7.40
Residual	0.410	-0.320	-0.002	-0.276	-0.045	-0.166	-0.038	-0.380	1.412	-0.150	0.192	-0.486	-0.152

(j) Plot the data and residuals versus predicted values on the following templates.



(k) Test the hypothesis $H_0 : \beta_1 = 0$ vs $H_1 : \beta_1 \neq 0$ at $\alpha = .05$ by using the F test statistic.

Table 1: Analysis of Variance Table

Source of Variation	df	Sum of Squares	Mean Sum of Squares	F statistic
Mean	1	547.8209		
Regression	1	25.808	25.808	99.97557
Residual Error	11	2.8396	.258143	
Total	13	576.4685		

ANSWER: $F_{1,11,.05} = 4.844336$

Because $F = 99.97557 > 4.844336 = F_{1,11,.05}$, then reject $H_0 : \beta_1 = 0$. We note that the p-value = 7.402023×10^{-7} .

(l) Interpret the QQ Plot: Test of Normality of the Residuals shown in Figure 2.

ANSWER: According to Figure 2, a departure of the residuals from normality appears on the right side of the graph. These three points correspond to children: 11, 2, 3, respectively. It is necessary to investigate the reasons why these three points are different from the rest, because it was assumed in the formulation of the linear model that $\epsilon_i \sim N(0, \sigma^2)$.

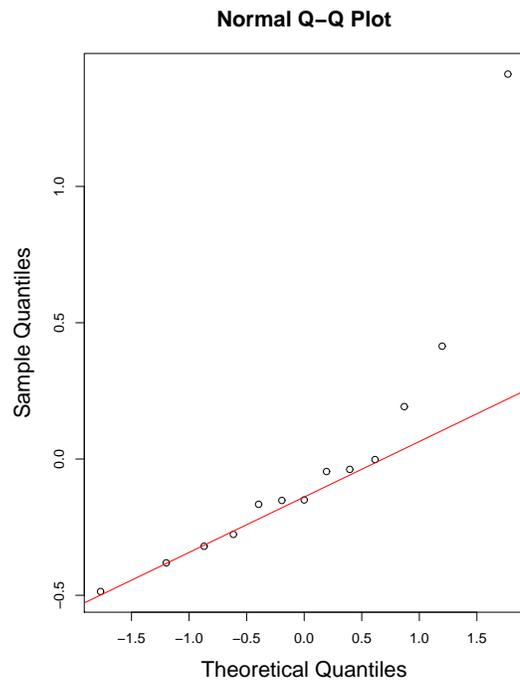


Figure 2: