

# Solutions to the Review 1

①  $\pi = 0.6$  calls resolved within 5 minutes = 5

ⓐ  $\underbrace{SS \dots S}_{20} = 0.6^{20}$  or  $P(x=20) = \binom{20}{20} 0.6^{20} 0.4^0$   
 $= \boxed{0.6^{20}}$

ⓑ  $\underbrace{FF \dots F}_{20} = 0.4^{20}$  or  $P(x=0) = \binom{20}{0} 0.6^0 0.4^{20}$   
 $= \boxed{0.4^{20}}$

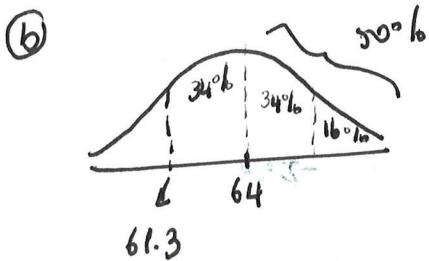
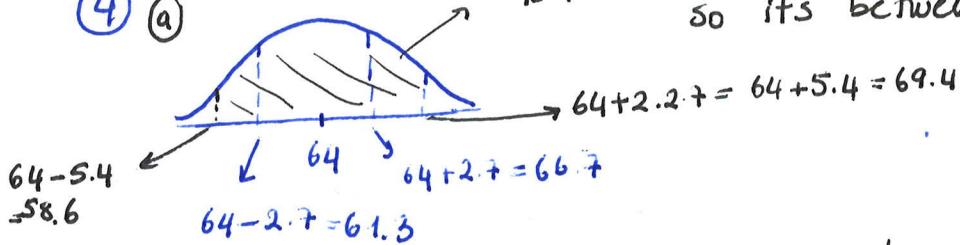
ⓒ  $P(x \geq 1) = 1 - P(x=0) = \boxed{1 - 0.4^{20}}$

② Observational, because ...

③ ⓑ  $4316 + 8986 + 1182 = 14,484$

ⓑ No, this poll clearly consists of volunteers who decided to participate. Note that we have no guarantee that these responses came from 14,484 distinct people; some may have voted more than once.

④ ⓑ 95% is within 2 stand. dev. away from the mean  
 $\therefore$  it's between 58.6 and 69.4



34% + 34% + 16% = 84% of the women are taller than 61.3.

⑤  $\frac{3.1 + 2.4 + 1.5 + x + 1.7}{5} = 2.3$

$$\begin{aligned} 8.7 + x &= 11.5 \\ x &= 2.8 \end{aligned}$$

ⓑ New data is  $\frac{x}{5}$   
 $3.1 \quad 2.4 \quad 1.5 \quad 2.8 \quad 5$        $\bar{x} = \frac{3.1 + 2.4 + 1.5 + 2.8 + 5}{5} = 2.96$

Mean is a sensitive measure of the center.

⑥  $x = \text{depth of flooding}$   
 $y = \text{flood damage}$

a)  $\hat{y} = 13.96 + 3.170x$

b) The plot with the line drawn in suggests that perhaps a simple linear regression model may not be appropriate. It suggests that there might be a curvilinear relationship may exist between depth and damage.

c) when  $x=6.5$   $\hat{y} = 13.96 + 3.170(6.5) = 34.565$

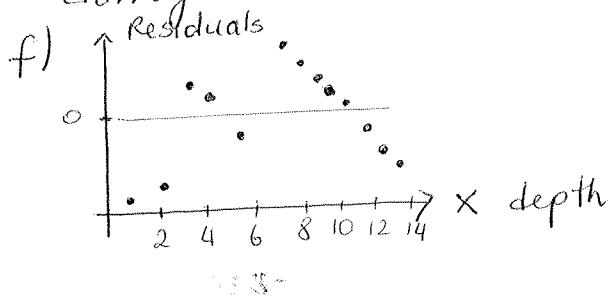
d) Since  $x=18$  is outside of the range of  $x$ -values for which data has been collected, we have no information concerning the relationship in the vicinity of  $x=18$ . So one would not want to use the least squares line to predict flood damage when depth of flooding is 18 feet.

e) We want to find first the mean and the standard deviation of the flood damage data.

$$\bar{x}_y = 36.15 \quad z \text{ score for } 26 = \frac{26 - 36.15}{13.31} = -0.763.$$

$$s_y = 13.31$$

So when flood damage is 26 you know that you are within one standard deviation of the mean of the flood damage data set.



Finding the residuals on minitab  
 Stat  $\rightarrow$  Regression  $\rightarrow$  Fitted Line Plot  
 $\rightarrow$  storage  $\checkmark$  residuals  
 $\checkmark$  fits.

This will create a new column with residuals and fits.  
 Then graph a scatterplot with  
 $y \rightarrow$  residuals  
 $x \rightarrow$  flood depth.

g) Stat  $\rightarrow$  Display Descriptive Statistics.

Variables: y

Mean	St Dev	Min	Q <sub>1</sub>	Median	Q <sub>3</sub>	Max
36.15	13.31	10	27	43	46.5	119

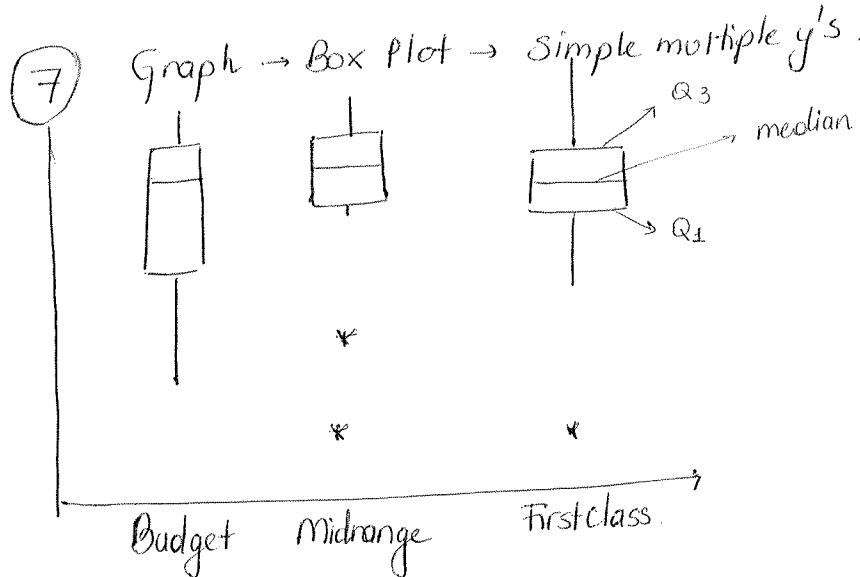
Five # summary consists of

min Q<sub>1</sub> median Q<sub>3</sub> max

h) Mean and median are very off from each other. This would be a data set that's skewed negatively that's why the mean is pulled

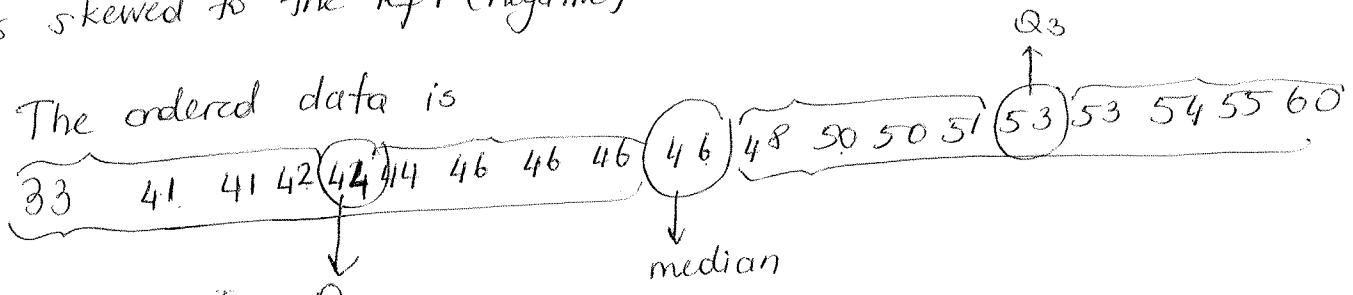
(3)

down towards the left tail. For center and spread I would use five number summary and also IQR.



The median franchise cost is about the same for each type of the hotel. The variability of franchise cost differs substantially for the three types. With the exception of outliers first class hotels vary very little in the cost of franchises while budget hotels vary a lot. Ignoring the outliers the distribution of franchise costs for first class hotels is quite symmetric, midrange hotels slightly skewed to the right (positive) while budget hotels distribution is skewed to the left (negative).

8) The ordered data is



a)  $IQR = 53 - 44 = 9$

$$1.5 \cdot 9 = 13.5$$

$$Q_3 + 13.5 = 53 + 13.5 = 66.5$$

mild outliers on the positive side  
no outliers on the negative side.

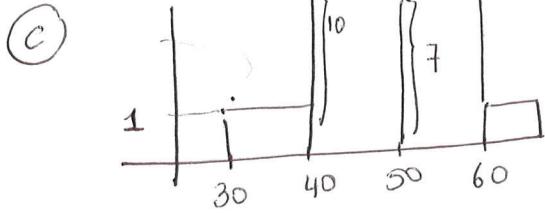
$$Q_1 - 13.5 = 44 - 13.5 = 30.5$$

//

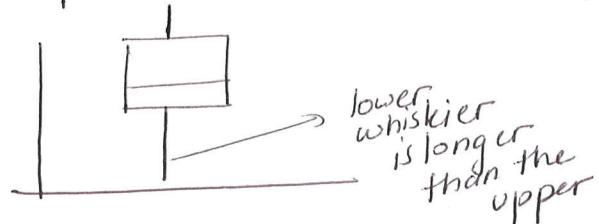
Hence there aren't any extreme outliers.

(b)

3	3
4	1 1 2 4 4 6 6 6 6 8
5	0 0 1 3 3 4 5
6	0



Boxplot



(d) The mean of this data is 47.53 and the median is 46. This suggests a little bit of positive skewness in the middle half of the data. And looking at the boxplot we see negative skewness in the extremes of the data.

- (10)
- (a) Correlation coefficient is between -1 and 1 and doesn't have a unit and is a numerical value between 2 numerical variables.
  - (b) Correlation coefficient is a numerical value between two numerical variables, however gender is a categorical variable.

- (11)
- |              | Mean  | Median | Q <sub>1</sub> | Q <sub>3</sub> |
|--------------|-------|--------|----------------|----------------|
| a) Pure tone | 106.2 | 72     | 38             | 155.5          |
| monkey call  | 176.6 | 141    | 91             | 205.5          |
- These describe very well that all the numerical measurements are higher for monkey call than for pure tones.

b) Call = 93.92 + 0.7783 Tone

c) r squared = .408 or 40.8%

40.8% of the variability in monkey calls can be explained by the linear relationship between pure tones and monkey calls.

- (1) The 3<sup>rd</sup> point (241, 485) has the largest residual (203.5045)

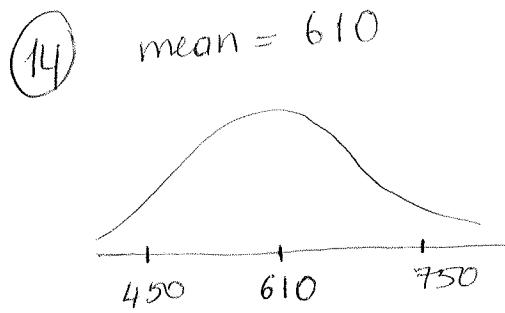
e) The first point  $(174, 500)$  is an outlier in the x-direction. (5)  
The value of  $r = .63$  q with this point but  $r$  drops to  
.479 when this point is removed. This point has a  
major impact on  $r$ .

f) No, the Jones range from 19 to 474 (the outlier) in  
the data set so this would go well beyond the range of  
variable data. The regression line should not be used  
for extrapolation. → 2 Pts.

- (12) • The different varieties have different lengths.  
Yellow is the shortest, red is next and bihai is the longest.  
The means, medians or any other measures of location  
clearly shows this fact.
- The variability as measured by IQR (height of each box)  
or by standard deviation is highest for red. The yellow  
has the lowest /smallest variability and bihai is in  
between.
  - The distribution of lengths for the yellow variety  
is symmetric, the other two distributions are  
skewed to the right.

(13) (Box-plots would be a good one to derive the above)  
Restricting the range of GPAs to perfect 4.0 will not  
allow us to explore the relationship with the explanatory  
variables.

b) Many possible answers. For example a simple random  
sample of students could be selected or a stratified  
sample to obtain equal numbers of students for each  
year could be obtained. Randomization should be used  
and association between GPA and other explanatory variables should  
be examined.



As a guess we can take

$$\frac{750 - 610}{140} \text{ as 2 st deviation}$$

$$\text{and then standard deviation} = \frac{140}{2} = 70.$$

(5)  $\frac{594 - 610}{70} = -0.228$

(15) (a) weight versus diameter  $r^2 = 94.4\%$   
Weight versus length  $r^2 = 72.3\%$

Diameter seems to give a better linear association  
than length.

(b) when we look at weight versus squared diameter  
we see that we get even a better linear association  
with  $r^2 = 99.7\%$

30 schools

2 schools is 6.6% of 30 schools.

We do know for normal distributions.

95% of the observations fall within  
2 standard deviation away from the  
mean

(6)