

Applied Regression Analysis

Homework 1

Mitchell Matheny

June 14, 2018

For problems 6 through 10. I have typed out the short version of the work with solutions. The hand written work is attached at the end of this document.

1 Exercise 1.1

The study used in this exercise looks at a study of teenage gambling in Britain. Make a numerical and graphical summary of the data, commenting on what you notice.

When initially looking at the summary of the data. One notices that the number male and female entries are treated as numerical data. We therefore need to convert them into what is called a factor. This is accomplished by using the following command:

```
teengamb$sex <- factor(teengamb$sex)
```

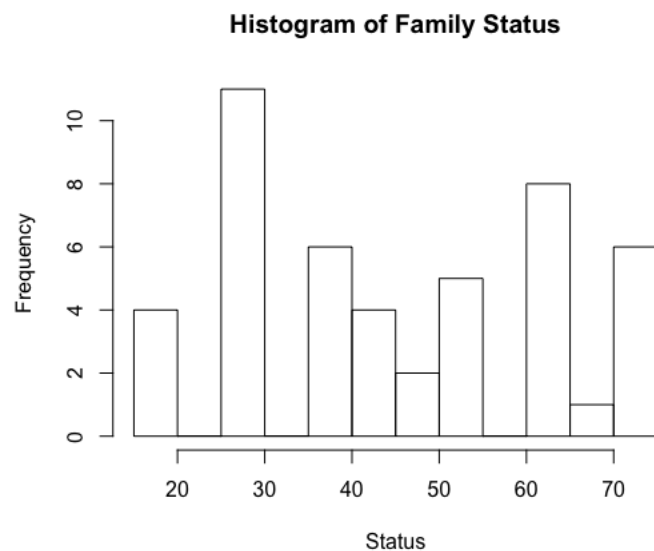
After doing so the summary looks like this:

```
> summary(teengamb)
      sex      status      income      verbal      gamble
Male :28   Min.   :18.00   Min.    : 0.600   Min.    : 1.00   Min.    : 0.0
Female:19   1st Qu.:28.00   1st Qu.: 2.000   1st Qu.: 6.00   1st Qu.: 1.1
              Median :43.00   Median : 3.250   Median : 7.00   Median : 6.0
              Mean    :45.23   Mean     : 4.642   Mean     : 6.66   Mean    :19.3
              3rd Qu.:61.50   3rd Qu.: 6.210   3rd Qu.: 8.00   3rd Qu.:19.4
              Max.    :75.00   Max.     :15.000   Max.     :10.00   Max.    :156.0
```

We can now easily see how many males and females were included in this study. This will also make it much easier to differentiate the sex on the various plots to follow.

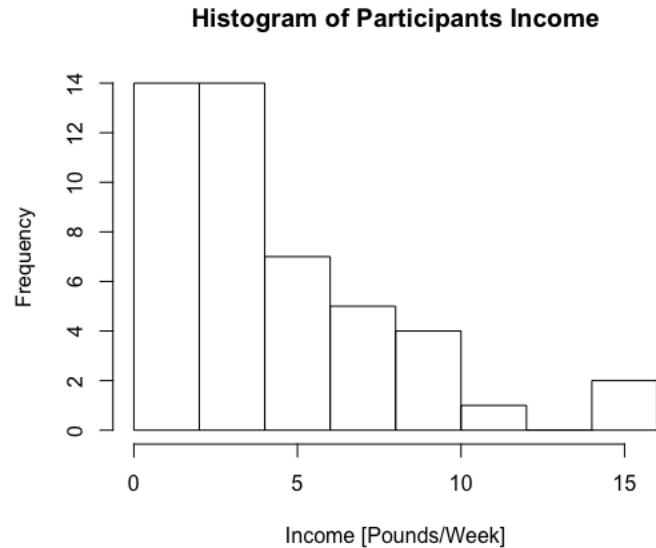
Comments on Numerical Summary and Univariate Plots:

Status The numerical values for family status seem a bit arbitrary. The documentation doesn't really specify how these numbers were calculated, this would take a bit more research into the actual study. With numerical values ranging from 18 to 75, both the mean and median are above the middle of the range of 40. Indicating that the societal status of the parents of the teenagers in this study are on the higher end of the spectrum. We can see this more clearly by looking at a histogram of the data:



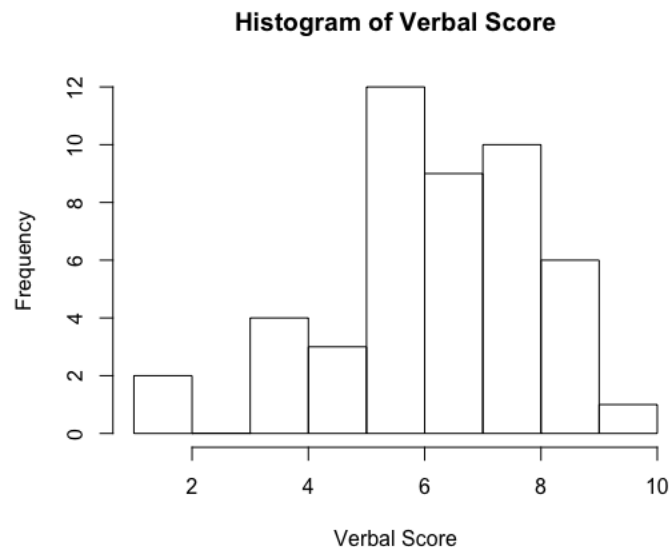
As we can see, the majority of the participants have a status above 40.

Income: The numerical values for income are given in pounds per week. The average income is 4.6 pounds per week, which is equivalent to only \$6.16 U.S. This seems like a very low value compared to the average of 500 pounds per week (\$670 U.S), but considering that these are teenagers, the values seem acceptable. Looking at the histogram:



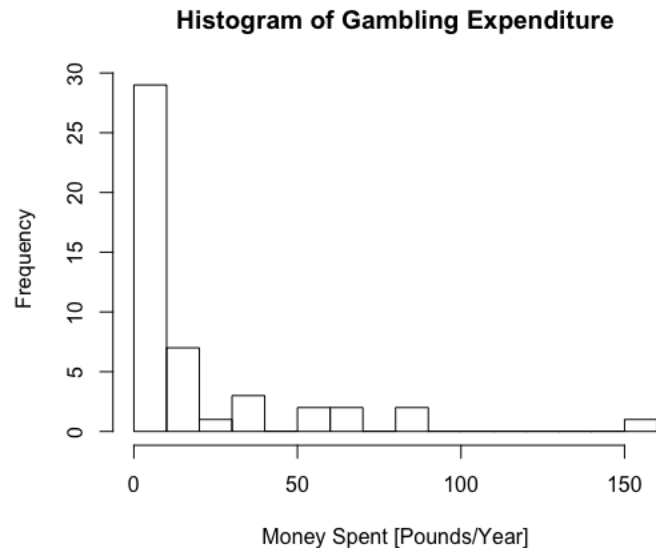
We can see the the 60% of the participants make less than 5 pounds per week (\$6.69 U.S), and the average is screwed by the 2 participants making 15 pounds per week (\$20.00 U.S)

Verbal: The verbal tests were scored on a scale of 0 to 12. Looking at the summary, we can see that the average value of 6.6 is just above the 50% correct mark, possibly indicating that the scores are screwed towards the higher end of the scale. Looking at the histogram:



We can clearly see that the majority of the participants scored 50% or better.

Gamble: The participants reported the amount of money the spent on gambling in pounds per year. Looking at the summary, we can see that the values have a large range. The minimum being 0 pounds per year spent to 156 pounds spent per year. Due to the large maximum value the mean is shifted to a higher value of 19.3. We can see however that the majority of the participants spend little to no money on gambling as the median value is 6. Looking at the histogram, we can see that this is that case.



Multivariate Plots

For the Multivariate plots, I will only make comments on plots that i notice trends on.

When looking at all of the plots relating a variable to sex, the only visually noticeable trend would be between the amount gambled per year: Men spent more per year than women. There is also a slight trend between the amount of money earned per week and the amount gambled per year. There appears to be a slight upward linear trend, indicating that the more earned per week, the more likely the participant was to spend money on gambling. This however is a speculation, and would need further analysis to confirm. The last trend i noticed was there is a slight linear trend when comparing verbal scores and parental status. The participants with a higher family status tended to score higher on the verbal test. Further analysis is needed to determine the true relationship.

Plots are located after problems 6-10

2 Problems 6-10

Problem 6:

Suppose that Y is a random variable with pdf, $f(y) = \lambda e^{-\lambda y}$ Determining the mean and variance of Y.
The mean is just the expectaion value.

$$E[y] = \int y \lambda e^{-\lambda y} dy \quad (1)$$

Solving via integration by parts, the mean is therefore:

$$E[y] = \frac{-e^{-\lambda y}(\lambda y + 1)}{\lambda} \quad (2)$$

To find the Variance:

$$Var[y] = E[y^2] - E[y]^2 \quad (3)$$

$$E[y^2] = \int y^2 \lambda e^{-\lambda y} dy = \frac{e^{-\lambda y}(\lambda^2 y^2 + 2\lambda y + 2)}{\lambda^2}$$

Therefore the variance:

$$Var[y] = \frac{e^{-\lambda y}(\lambda^2 y^2 + 2\lambda y + 2)}{\lambda^2} + \frac{-e^{-2\lambda y}(\lambda y + 1)^2}{\lambda^2} \quad (4)$$

Problem 7:

Determining the mean and variance of $4Y + 2$.

Using the following property:

$$E[aY + b] = aE[Y] + b \quad (5)$$

$$E[4Y + 2] = 4E[y] + 2$$

Therefore the mean is:

$$E[4Y + 2] = \frac{-4e^{-\lambda y}(\lambda y + 1)}{\lambda} + 2 \quad (6)$$

Using a similar argument for the variance:

$$Var[aY + b] = a^2 Var[Y] = 16 Var[Y] \quad (7)$$

$$Var[4Y + 2] = \frac{16e^{-\lambda y}(\lambda^2 y^2 + 2\lambda y + 2)}{\lambda^2} + \frac{-16e^{-2\lambda y}(\lambda y + 1)^2}{\lambda^2}$$

Problem 8:

Suppose that random variables X and Y have a joint PDF given by $f(x, y) = 6(1 - y); 0 \leq x \leq y \leq 1$. Determine $f_x, f_y, E(X), E(Y), E(XY)$, and $cov(X, Y)$.

$$f_x = \int_x^1 (6 - 6y) dy$$

$$f_x = 3x^2 - 6x + 3$$

$$f_y = \int_0^y (6 - 6y) dx$$

$$f_y = -6y^2 + 6y$$

$$E[x] = \int_0^1 x f_x dx = \int_0^1 (3x^3 - 6x^2 + 3x) dx$$

$$E[x] = \frac{1}{4}$$

(8)

$$E[y] = \int_0^1 x f_y dy = \int_0^1 y(-6y^2 + 6y) dy$$

$$E[y] = \frac{1}{2}$$

$$E[xy] = \int_0^1 \int_0^1 xy f(x, y) dx dy = \int_0^1 \int_0^1 xy(6 - 6y) dx dy$$

$$E[xy] = 1$$

$$Cov[xy] = E[xy] - E[x]E[y]$$

$$Cov[xy] = \frac{7}{8}$$

Problem 9:

Suppose that $y = (Y_1, Y_2)^T$ is a bivariate normal random vector with mean and variance;

$$E[y] = \begin{pmatrix} 1 & 2 \end{pmatrix}^T \quad (9)$$

$$Var[y] = \begin{pmatrix} 1 & 0.25 \\ 0.25 & 2 \end{pmatrix}$$

What is the distribution, mean, and variance of $a^T y$, where $a = \begin{pmatrix} 1 & 3 \end{pmatrix}^T$

$$\begin{aligned} f(y) &= a^T y = \begin{pmatrix} 1 & 3 \end{pmatrix} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \\ f(y) &= Y_1 + 3Y_2 \end{aligned} \quad (10)$$

For the mean:

$$\begin{aligned} E[a^T y] &= a^T E[y] \\ E[a^T y] &= \begin{pmatrix} 1 & 3 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} \end{aligned} \quad (11)$$

$$E[a^T y] = 7$$

For the variance:

$$\begin{aligned} Var[a^T y] &= a Var[y] a^T \\ Var[a^T y] &= \begin{pmatrix} 1 & 3 \end{pmatrix} \begin{pmatrix} 1 & 0.25 \\ 0.25 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \end{pmatrix} \end{aligned} \quad (12)$$

$$Var[a^T y] = 20.5$$

Problem 10:

Using the same y as in problem 4, determine the distribution, mean, and variance of $z = Ay$, where

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \quad (13)$$

$$z = Ay = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \quad (14)$$

$$z = \begin{pmatrix} 2Y_1 + Y_2 \\ Y_1 + 2Y_2 \end{pmatrix}$$

For the mean:

$$\begin{aligned} E[AY] &= AE[Y] \\ E[AY] &= \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} \end{aligned} \quad (15)$$

$$E[AY] = \begin{pmatrix} 4 \\ 3 \end{pmatrix}$$

For the Variance:

$$\begin{aligned} Var[AY] &= AVar[Y]A^T \\ Var[AY] &= \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 0.25 \\ 0.25 & 2 \end{pmatrix} \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \end{aligned} \quad (16)$$

$$Var[AY] = \begin{pmatrix} 3.5 & 2.5 \\ -0.5 & 6 \end{pmatrix}$$

