
Problem Set 1
Introduction to R
E63 Big Data Analytics

Harvard Extension School
Andrew Caide

September 4, 2017
Problem Set 1

Contents

1	Problem Set Questions	2
1.1	Problem 1	2
1.2	Problem 2	2
1.3	Problem 3	3
1.4	Problem 4	3
1.5	Problem 5	3
1.6	Problem 6	4
1.7	Document History	4
2	Solutions	5
2.1	Problem 1	6
2.2	Problem 2	9
2.3	Problem 3	9

Chapter 1

Problem Set Questions

1.1 Problem 1

Binomial distribution describes coin tosses with potentially doctored or altered coins. Value of p is the probability that head comes on top. If both the head and the tail have the same probability, $p = 0.5$. If the coin is doctored or altered, p could be larger or smaller. Plot on three separate graphs the binomial distribution for $p = 0.3$, $p = 0.5$ and $p = 0.8$ for the total number of trials $n = 60$ as a function of k , the number of successful (head up) trials. Subsequently, place all three curves on the same graph. For each value of p , determine 1st Quartile, median, mean, standard deviation and the 3rd Quartile. Present those values as a vertical box plot with the probability p on the horizontal axis. (15%)

1.2 Problem 2

Problem 2. Finish the plot of the correlation between waiting times and durations of Old Faithful data. Recreate the scatter plot of waiting vs. duration times. As we mentioned in class, the best linear assessment in the sense of the least squares fit of a relationship (proportionality) between two or many variables can be achieved with R function `lm()`. `lm` stands for the linear model. The first argument of `lm()` is called formula accepts a model which starts with the response variable, waiting in our case, followed by a tilde (symbol `~`, read as "is modeled as") followed by the (so called Wilkinson-Rogers) model on the right. In our case we simply assume that waiting time is proportional to the duration time and that `model` reads: `formula = waiting ~ duration`. The second argument of function `lm()` is called data and, in our case, will take value `faithful`, the data set containing our data. Store the result of function `lm()` in a variable. The name of that variable is not essential. Call it `model`. Print

the variable. The first component of that variable is the intercept of calculated line with the vertical axis (waiting, here) and the second is the slope of the line. Convince yourself that line with those parameters will truly lie on your graph. Function `abline()` adds a line to the previously created graph. Next, pass the variable model to the function `abline()`. Make that line somewhat thicker and blue. Use `help(functionName)` to find details about invocations of both `lm()` and `abline()` functions. (20%)

1.3 Problem 3

Calculate the covariance matrix of the faithful data. Determine the eigenvalues and eigenvectors of that matrix. Demonstrate that two eigenvectors are mutually orthogonal. Demonstrate that the eigenvector with the larger eigenvalue is parallel with line discovered by `lm()` function in the previous problem. (15%)

1.4 Problem 4

You noticed that eruptions clearly fall into two categories, short and long. Let us say that short eruptions are all which have duration shorter than 3.1 minute. Add a new column to data frame `faithful` called `type`, which would have value `"short"` for all short eruptions and value `"long"` for all long eruptions. Next use `boxplot()` function to provide your readers with some basic statistical measures for waiting. In a separate plot present the box plot for duration times. Please note that `boxplot()` function also accepts as its first argument a formula such as `waiting ~ type`, where `waiting` is the numeric vector of data values to be split in groups according to the grouping variable `type`. The second argument of function `boxplot()` is called `data`, which in our case will take the name of our dataset, i.e. `faithful`. Find a way to add meaningful legends to your graphs. Subsequently, present both boxplots on one graph. (20%)

1.5 Problem 5

Create a matrix with 40 columns and 100 rows. Populate each column with random variable of the uniform distribution type. Make those distributions symmetric around zero. Let the distribution for each column appear like the one on slide 92 of the lecture note, except centered around zero. Present two distributions contained in any two randomly selected columns of your matrix on two separate plots. Convince yourself that generated distributions are (close to) uniform. (15%)

1.6 Problem 6

Start with your matrix from problem 5. Add yet another column to that matrix and populate that column with the sum of original 40 columns. Create a histogram of values in the new column showing that the distribution starts to resemble the Gaussian curve. Add a true, calculated, Gaussian curve to that diagram with the parameters you expect from the sum of 40 random variables of uniform distribution with values between -1 and 2. (15%)

1.7 Document History

- September 2 2017: first draft document

Chapter 2

Solutions

2.1 Problem 1

This question explores binomial distributions with different weights (or probabilities). A set of results that are binomially distributed have only two results, like a coin flip (heads or tails).

Firstly we explore what a weighted distribution looks like, with three weights: 0.3 (20% lower success rate), 0.5 (equally weighted), and 0.8 (30% higher success rate).

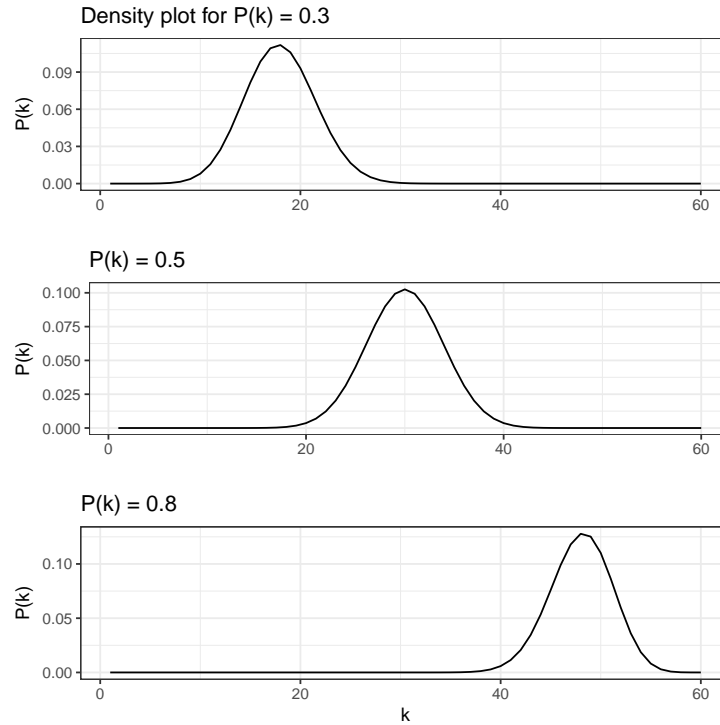


Figure 2.1: Question 1, Individual Density Plots of Binomial distributions of three trials where $n = 60$.

To confirm the distributions in Figure 2.1, we can recall the mean of a binomial distribution:

$$\mu_k = np$$

where μ_k is the mean (on k axis), with probability p , and n number of trials. Furthermore binomial distributions have a variance:

$$\sigma_k^2 = np(1 - p)$$

where σ_k^2 is the standard deviation. It should be noted that as the probability increases the variance tightens as the probability increases. An eyeball test indicate our plots were spot on: $\mu_k\{0.3, 0.5, 0.8\} = 60 * \{0.3, 0.5, 0.8\} = \{18, 30, 48\}$.

Before moving on, I would like to use histograms for the next plot because coinflip measurements are discrete: there's no way somebody could have obtained 15.3 heads + 0.7 tails. The plot will have 60 bins, one per possible outcome ($n = 60$). It should also be noted that the probability values are going to appear greater - they should. The (continuous) density plots weigh the probability on all space, including invalid results. Now that the results have made discrete the likelihood of outcomes should look a lot more reasonable.

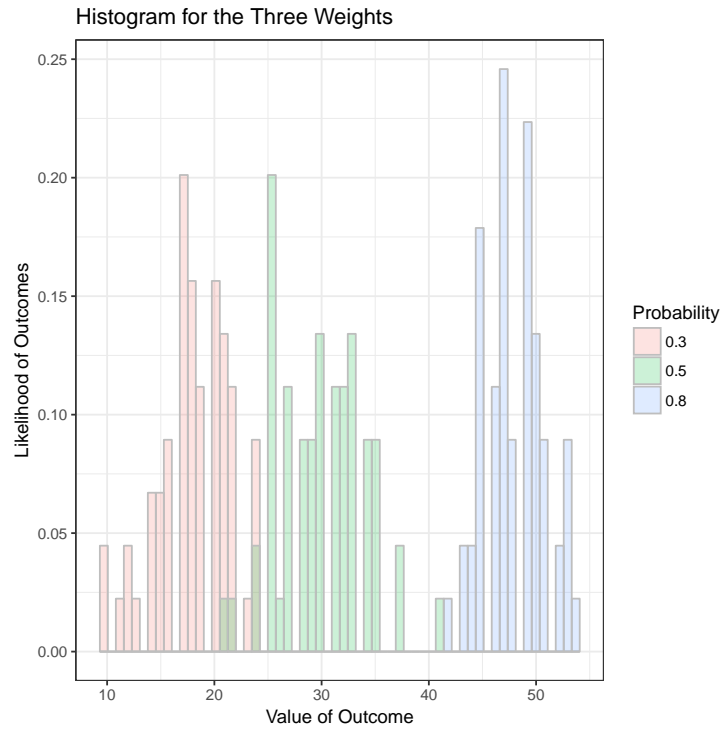


Figure 2.2: Question 1, Binomial distributions of three trials where $n = 60$.

The summary of the binomial statistics are outlined in the following table, and boxplots on figure 2.3 further illustrate the median, 25% and 75% quartiles, and the max and min.

Table 2.1: Summary of Question 1 Data. Length of data set, 25 Percent quartile, Mean, Median, Standard Deviation, and 75 Percent quartile.

Probability	LowerQuartile	Median	Mean	SD	UpperQuartile
0.3	16	18	18	3	21
0.5	27	30	30	4	33
0.8	46	48	48	3	50

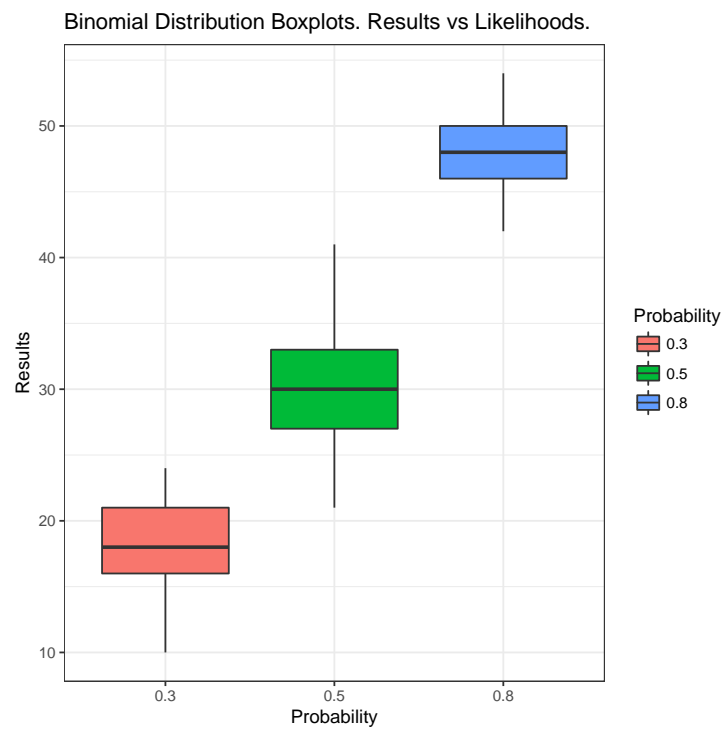


Figure 2.3: Question 1, Box Plots for the Binomial Distributions. Statistics outlined in aforementioned table.

2.2 Problem 2

This question investigates the relationship of two measures in the 'faithful' built in data-set. This data set comes from the Old Faithful Geyser in Yellowstone National Park, USA. Two measures will be examined in this data set: the time between geyser eruptions (in minutes), and the duration of the eruptions (in minutes). First we look at the data on a scatter-plot and see if we can apply a linear fit with a reasonable correlation.

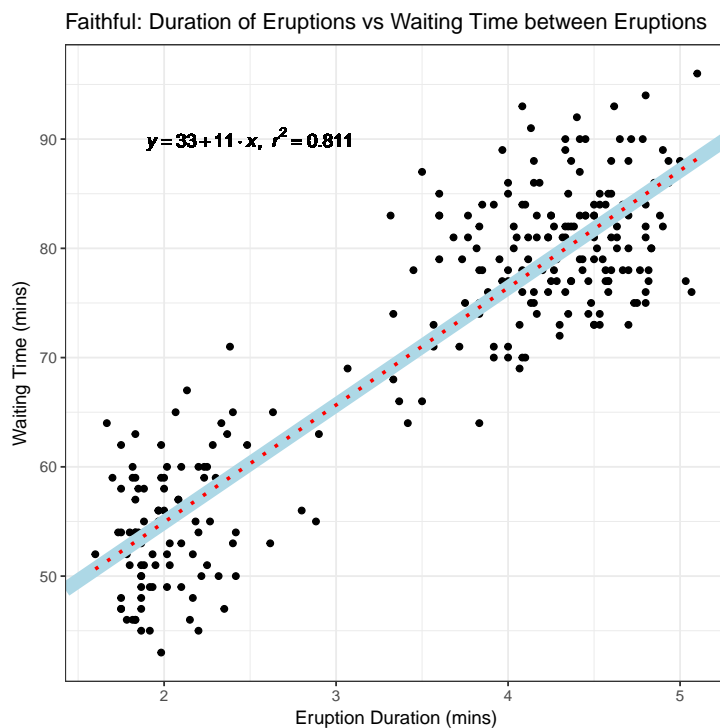


Figure 2.4: Question 2, Linear Fit on Faithful Data

The linear fit constructed from the `lm()` function is shown in blue. To verify, I allowed ggplot to apply its fit (without additional input), indicated by the red dots. The correlation coefficient is weak at $r^2 = 0.811$.

2.3 Problem 3

This problem revisits the faithful data. In the previous question we found ourselves unsatisfied with the correlation coefficient from the linear plot, but the plot certainly indicates two populations of events. Let's take a more rigorous

look.

First let's observe the covariance of the faithful data.

```
> df <- faithful
> covariance <- cov(df)
> covariance # See table
```

Table 2.2: Covariance result.

eruptions	waiting
1	14
14	185

The positive covariance between eruption durations and waiting times indicates a positive relationship between the two variables, which was already made fairly obvious. It is safe to ignore the diagonals in this matrix. Note: a covariance of 0 should indicate no relationship, and the magnitude of the covariance indicates a stronger relationship. Therefore a covariance of -130 indicates a tight (or relatively stronger than our 13.98), but negatively trending relationship.

From this we can calculate the eigen values and eigen vectors.

```
> eig.values <- eig$values
eigen.values

[1] 185.8818239    0.2442167

> eig.vectors <- eig$vectors
> eig.vectors
```

```
      [,1]      [,2]
[1,] -0.9885959 -0.1505924
[2,]  0.1505924 -0.9885959
```

We can also demonstrate orthogonality.

```
> eig.vectors * t(eig.vectors)
```

```
      [,1] [,2]
[1,]    1    0
[2,]    0    1
```