

Examining the number of characters in an email

Your Name Here

Put the date here

Load packages and data set.

We will be using the *email50* data set which you have already downloaded and placed in your class folder. Read this data into R by filling in your path. If you don't recall where it is at, check lab 3 and simply copy the code used there to read it in.

```
email50 <- read.delim("_____/email50.txt", header=TRUE, sep="\t")
```

At this point click the Knit HTML button at the top and check to see if your document is created correctly.

Summary Statistics

Let's examine the number of characters `num_char` contained in an email.

What are the `minimum` and `maximum` number of characters in emails from this dataset? Don't forget to *Knit html* after each chunk to make sure the output is being produced correctly.

```
min(email50$____);max(____$_____)
```

Use the `range()` function to find these numbers in another way.

```
____(email50$_____)
```

What is the `mean` and `variance` of this variable?

```
____(email50$____)  
var(____$_____)
```

Univariate visualizations (one variable)

Let's visualize the distribution of `num_char`.

Histograms

1. Create a histogram of `num_char` using the base `hist(x)` plotting function.

```
hist(_____)
```

2. Use the `breaks=` argument to create 3 breaks in the histogram. How many bins are created?

```
_____(_____, breaks=3)
```

3. Does a histogram with 13 bins provide a different description of the data distribution than when there are only 4 bins? Draw one and find out.

```
_____(_____, breaks=___)
```

Boxplots

1. Calculate the five number `summary`.

```
summary(_____)
```

2. Draw a boxplot using `boxplot(x)`.

```
boxplot(_____)
```

3. Redraw the boxplot horizontally.

```
boxplot(_____, horizontal=_____)
```

Bivariate visualizations (two variables)

This is when multiple variables are plotted against each other. Or one variable is plotted against levels of another variable.

Scatterplot

We would expect that the more characters in the email corresponds with more `line_breaks`. Let's create a scatterplot with `line_breaks` on the y axis and `num_char` on the x. The syntax for a scatterplot is `plot(y~x)` OR `plot(x, y)`. The first notation is called **model notation** and we will be using it frequently in this class.

```
plot(_____ ~ _____)
```

Grouped boxplot

Can you use the `num_char` or `line_breaks` to determine if an email is spam or not? Creating grouped boxplots can help us see if one group has higher measurements compared to the other. Again we use the model notation and type `boxplot(c~g)`, where in this case `c` is the continuous variable and `g` is the grouping variable `spam`.

1. Create a boxplot of `num_char` against `spam`.

```
boxplot(____$num_char ~ email50$____)
```

2. Create a horizontal grouped boxplot of `line_breaks` against `spam`.

```
boxplot(email50$_____ ~ _____$spam, horizontal=_____)
```