

R-basics for 431 (Small Version)

Thomas E. Love

August 30, 2016

Contents

1	The chickwts study	1
1.1	Summarizing the distribution of a categorical variable, feed	2
1.2	Describing the distribution of a continuous variable, weight , numerically	2
1.3	Drawing an attractive histogram of the weight data	3
1.4	Drawing a Boxplot of the Weights	6
1.5	Drawing a Normal Q-Q plot of the Weights	8
2	The Orange Study	10
2.1	Numerical Summary	10
2.2	Correlation between Age and Circumference	11
2.3	Scatterplot predicting Circumference using Age across all Trees	12
2.4	Scatterplot with Linear Fit predicting Circumference using Age across all Trees	13
2.5	The Linear Model	14
2.6	Some Other Scatterplots: Assessing each Tree separately	15
3	Getting Data into R from Excel or another Software Package: The Fundamentals	17

Working through this document will definitely help you get rolling on Assignment 1.

To begin, we'll load the packages (libraries) that we will use in our analyses.

```
knitr::opts_chunk$set(comment=NA)
library(dplyr); library(mosaic); library(ggplot2)
```

1 The chickwts study

The **chickwts** data, available as part of the base installation of R (in the **datasets** package) describe an experiment conducted to measure and compare the effectiveness of various feed supplements on the growth rate of chickens. For more on the **chickwts** data, type `?(chickwts)` into the R console.

We'll begin by placing the data in a tibble called **chick**.

```
chick <- tbl_df(chickwts)
chick
```

```
# A tibble: 71 × 2
  weight    feed
  <dbl>    <fctr>
1    179 horsebean
2    160 horsebean
3    136 horsebean
4    227 horsebean
5    217 horsebean
6    168 horsebean
```

```

7      108 horsebean
8      124 horsebean
9      143 horsebean
10     140 horsebean
# ... with 61 more rows

```

- The `weight` variable is numeric (double-precision) and gives the chick's weight.
- The `feed` variable is categorical (a factor in R) and gives the feed type.

1.1 Summarizing the distribution of a categorical variable, `feed`

The regular `summary` function can provide some useful results.

```
summary(chick$feed)
```

```

      casein horsebean   linseed  meatmeal   soybean sunflower
      12         10         12         11         14         12

```

There are lots of ways to generate a table for a factor (categorical variable) like this.

```

chick %>%
  select(feed) %>%
  table() %>%
  addmargins()

```

```

      casein horsebean   linseed  meatmeal   soybean sunflower      Sum
      12         10         12         11         14         12       71

```

1.2 Describing the distribution of a continuous variable, `weight`, numerically

The regular `summary` function provides a five-number summary, plus the mean.

```
summary(chick$weight)
```

```

      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
108.0    204.5    258.0    261.3   323.5    423.0

```

The `favstats` function from the `mosaic` library produces a more extensive set of numerical summaries.

```
mosaic::favstats(chick$weight)
```

```

min    Q1 median    Q3 max    mean    sd n missing
108 204.5    258 323.5 423 261.3099 78.0737 71      0

```

Here is a smaller numerical summary of the weights broken down by feed category.

```

table1 <- chick %>%
  group_by(feed) %>%
  summarize(mean(weight), sd(weight), median(weight))

table1

# A tibble: 6 × 4
      feed `mean(weight)` `sd(weight)` `median(weight)`
  <fctr>      <dbl>      <dbl>      <dbl>
1 casein    323.5833    64.43384    342.0
2 horsebean 160.2000    38.62584    151.5

```

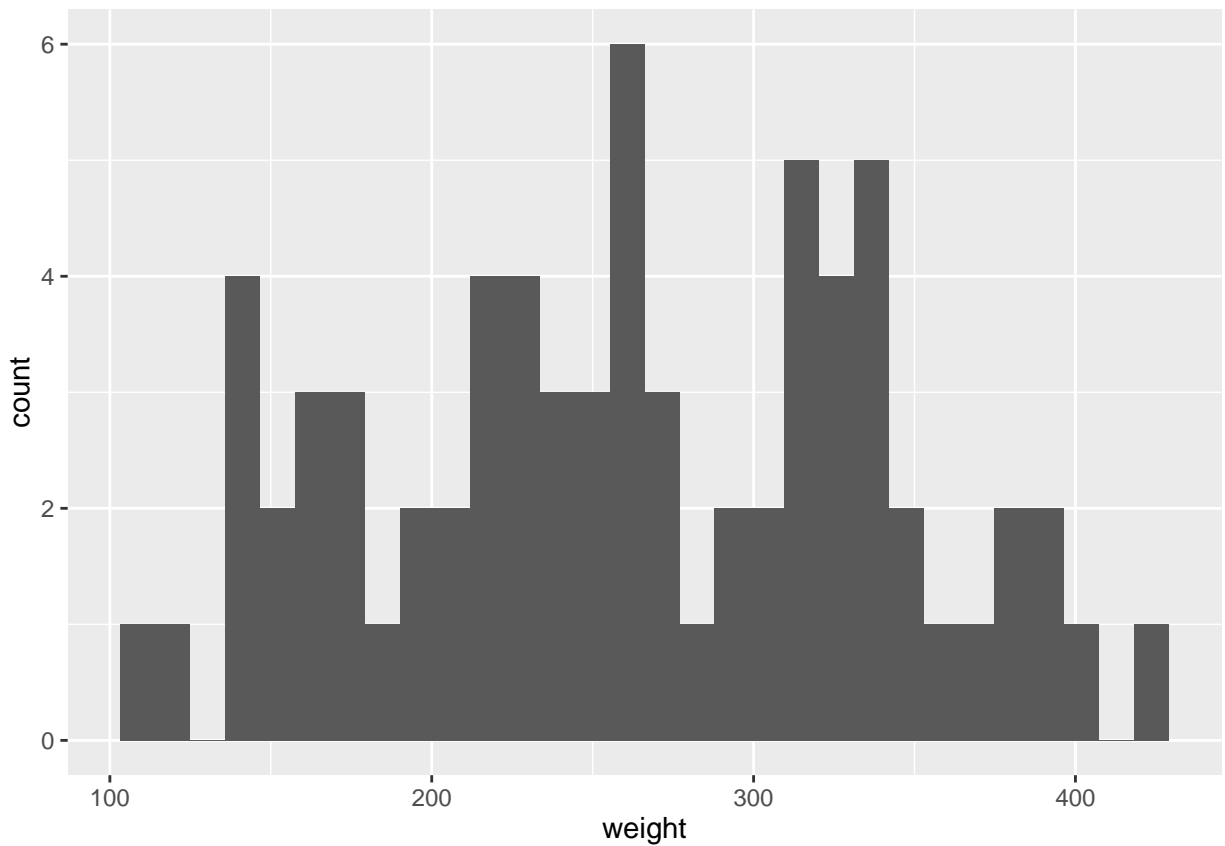
3	linseed	218.7500	52.23570	221.0
4	meatmeal	276.9091	64.90062	263.0
5	soybean	246.4286	54.12907	248.0
6	sunflower	328.9167	48.83638	328.0

1.3 Drawing an attractive histogram of the weight data

Here is the default approach.

```
ggplot(chick, aes(x = weight)) +  
  geom_histogram()
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



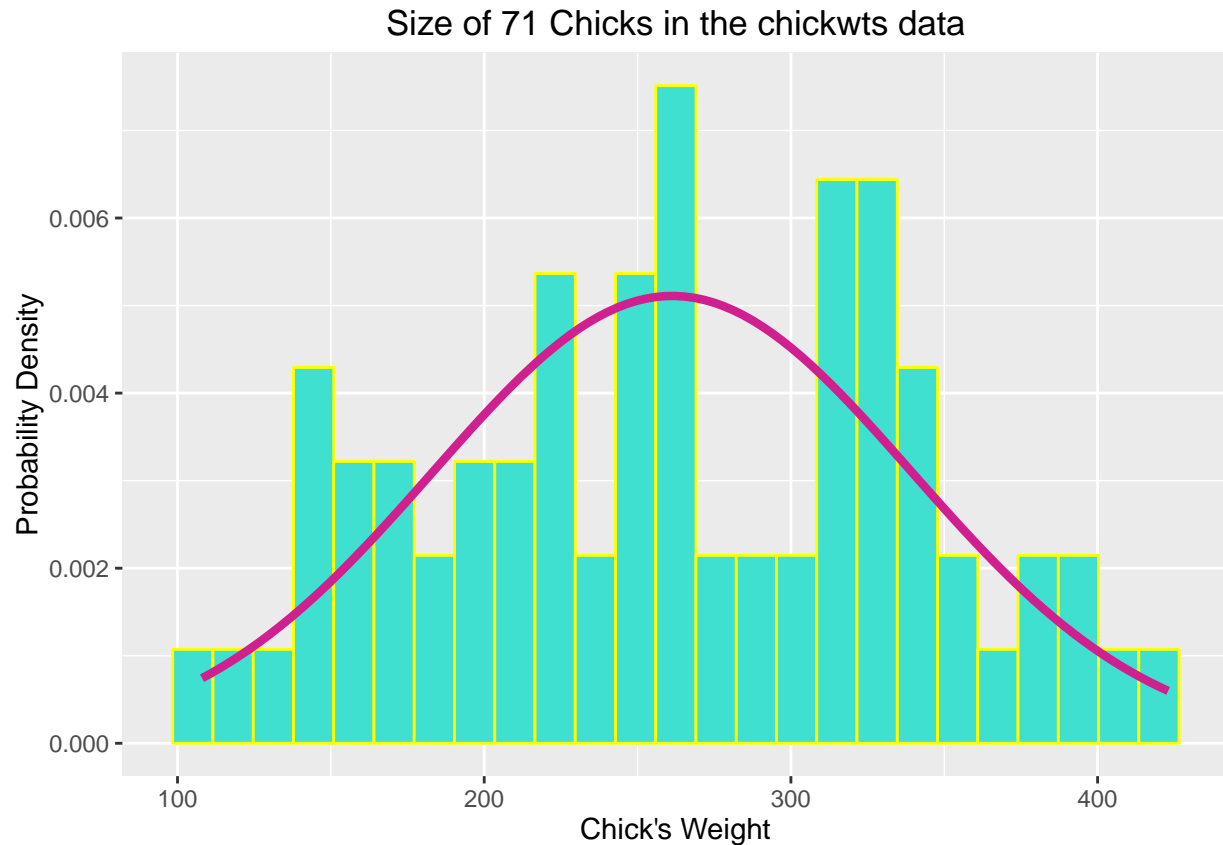
Let's make that slightly more attractive, revise the labels, and place a title.

```
ggplot(chick, aes(x = weight)) +  
  geom_histogram(bins = 20, color = "yellow", fill = "turquoise") +  
  labs(x = "Chick's Weight", y = "Number of Chicks",  
       title = "Size of 71 Chicks in the chickwts data")
```



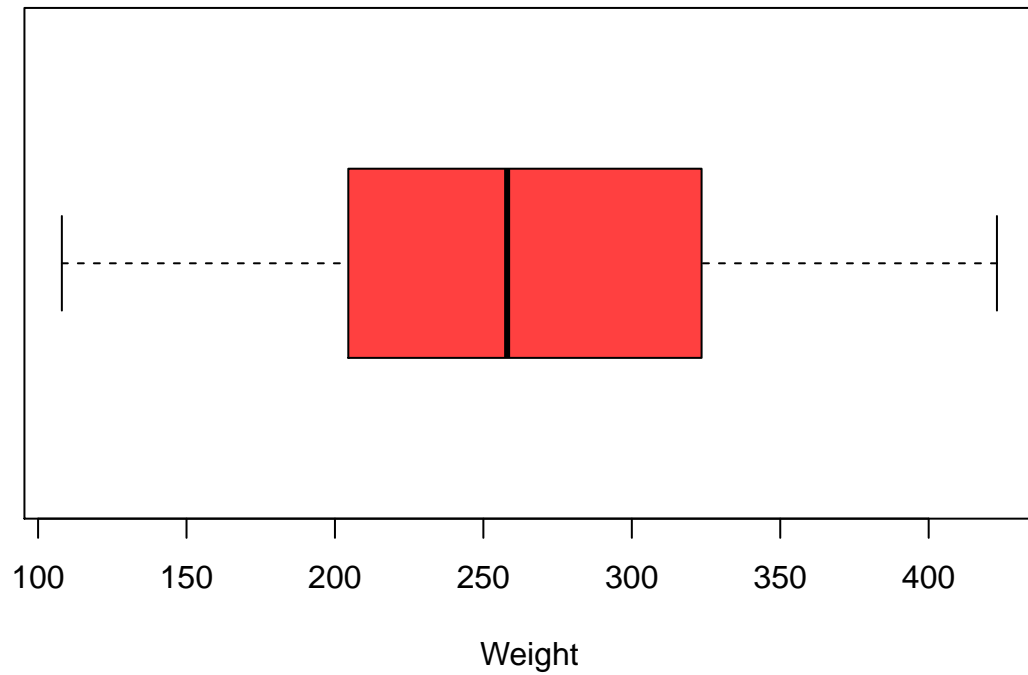
Another option would be to plot the density function, rather than the raw counts, and compare it directly to what we would expect from a Normal model with the same mean and standard deviation as the weights in the chick data.

```
ggplot(chick, aes(x = weight)) +  
  geom_histogram(aes(y = ..density..), bins = 25, color = "yellow", fill = "turquoise") +  
  stat_function(fun = dnorm,  
               args = list(mean = mean(chick$weight), sd = sd(chick$weight)),  
               lwd = 1.5, col = "violetred") +  
  labs(x = "Chick's Weight", y = "Probability Density",  
       title = "Size of 71 Chicks in the chickwts data")
```



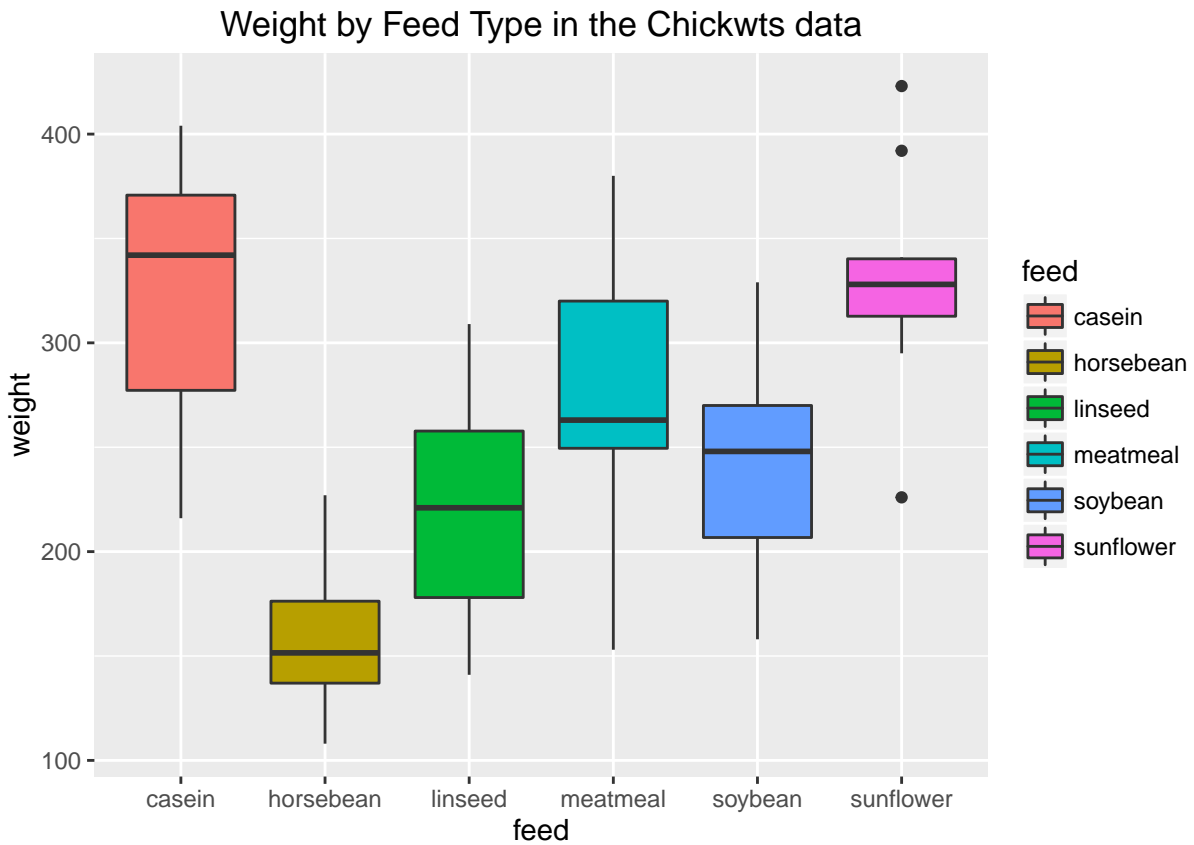
1.4 Drawing a Boxplot of the Weights

```
boxplot(chick$weight, col = "brown1", horizontal = TRUE, xlab = "Weight")
```



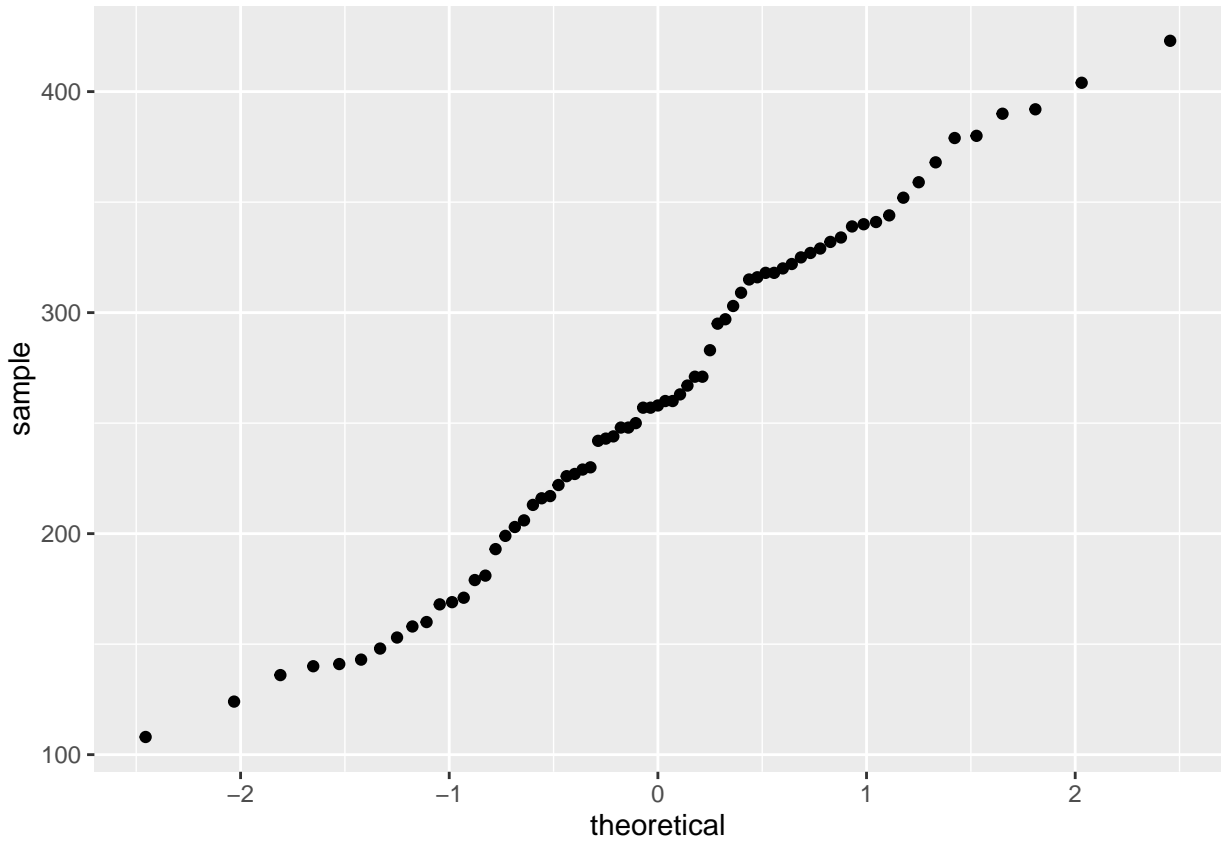
A more interesting boxplot would compare the weights across the various types of feed.

```
ggplot(chick, aes(x = feed, y = weight, fill = feed)) +  
  geom_boxplot() +  
  labs(title = "Weight by Feed Type in the Chickwts data")
```



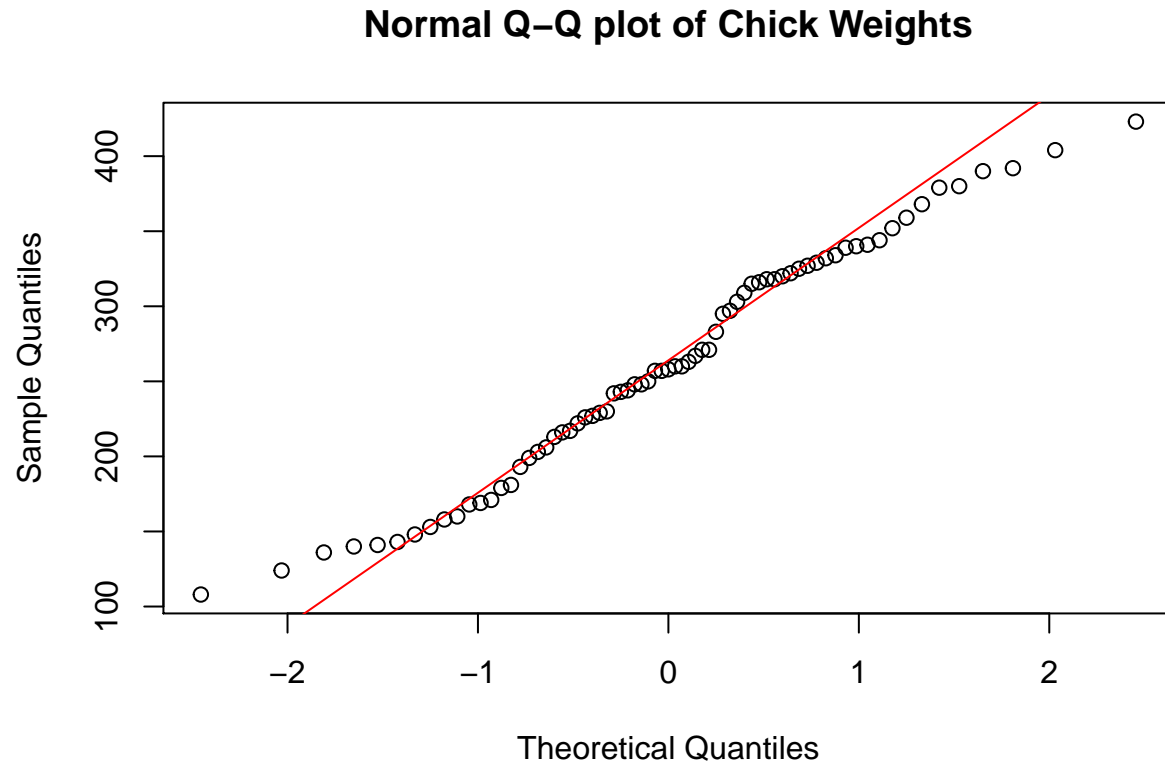
1.5 Drawing a Normal Q-Q plot of the Weights

```
ggplot(chick, aes(sample = weight)) +  
  geom_point(stat="qq")
```



An alternative method for obtaining a Normal Q-Q plot would be the following.

```
qqnorm(chick$weight, main = "Normal Q-Q plot of Chick Weights")  
qqline(chick$weight, col = "red")
```



2 The Orange Study

The Orange data frame has 35 rows and 3 columns of records of the growth of orange trees. Let's get the data into a tibble.

```
orange <- tbl_df(Orange)
orange
```

```
# A tibble: 35 × 3
  Tree   age circumference
*   <ord> <dbl>         <dbl>
1     1  118           30
2     1  484           58
3     1  664           87
4     1 1004          115
5     1 1231          120
6     1 1372          142
7     1 1582          145
8     2  118           33
9     2  484           69
10    2  664          111
# ... with 25 more rows
```

- **tree** is an ordinal factor, which indicates the tree on which the measurement was made. The ordering is by increasing maximum diameter of the five trees.
- **age** is a numerical variable, containing the age of the tree as measured in days since 1968-12-31.
- **circumference** is a numerical variable, containing the trunk circumference (probably at “breast height”) in mm.

2.1 Numerical Summary

And here's the standard numerical summary for the full data set.

```
summary(orange)
```

```
Tree      age      circumference
3:7  Min.   : 118.0   Min.       : 30.0
1:7  1st Qu.: 484.0   1st Qu.: 65.5
5:7  Median :1004.0   Median :115.0
2:7  Mean    : 922.1   Mean      :115.9
4:7  3rd Qu.:1372.0   3rd Qu.:161.5
      Max.    :1582.0   Max.       :214.0
```

Next, we'll look at the mean age and circumference, within each of the seven measurements per tree.

```
orange %>%
  group_by(Tree) %>%
  summarize(mean(age), mean(circumference))
```

```
# A tibble: 5 × 3
  Tree `mean(age)` `mean(circumference)`
  <ord>         <dbl>         <dbl>
1     3    922.1429         94.00000
2     1    922.1429         99.57143
3     5    922.1429        111.14286
4     2    922.1429        135.28571
```

```
5      4      922.1429      139.28571
```

It looks like each of the trees was measured at exactly the same time (age).

```
table(orange$age, orange$Tree)
```

```
      3 1 5 2 4
118   1 1 1 1 1
484   1 1 1 1 1
664   1 1 1 1 1
1004  1 1 1 1 1
1231  1 1 1 1 1
1372  1 1 1 1 1
1582  1 1 1 1 1
```

Yes, each tree was measured at precisely the same five times.

2.2 Correlation between Age and Circumference

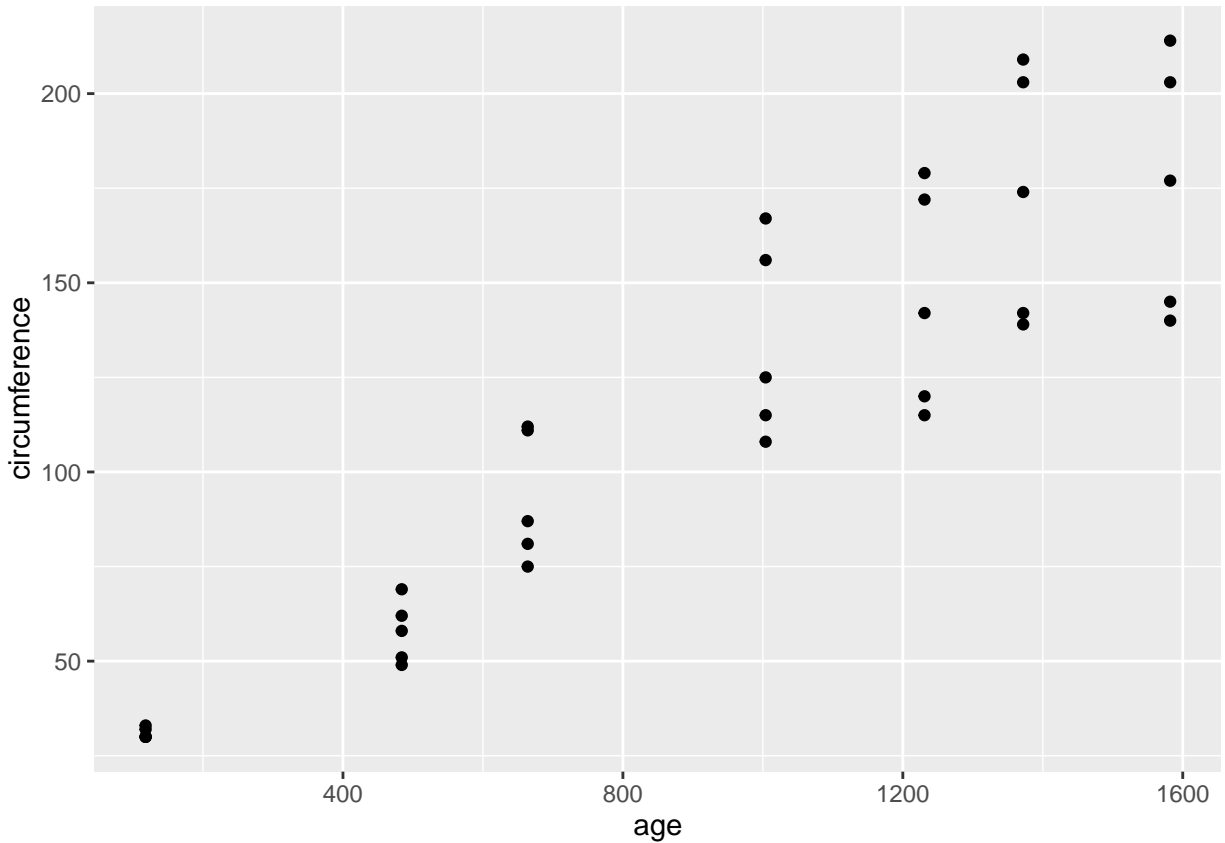
```
cor(orange$age, orange$circumference)
```

```
[1] 0.9135189
```

The Pearson correlation of age and circumference is 0.91 which is pretty strong, indicating that we'd expect to see a fairly positive and mostly linear association in a scatterplot. So, let's see if that's what we get.

2.3 Scatterplot predicting Circumference using Age across all Trees

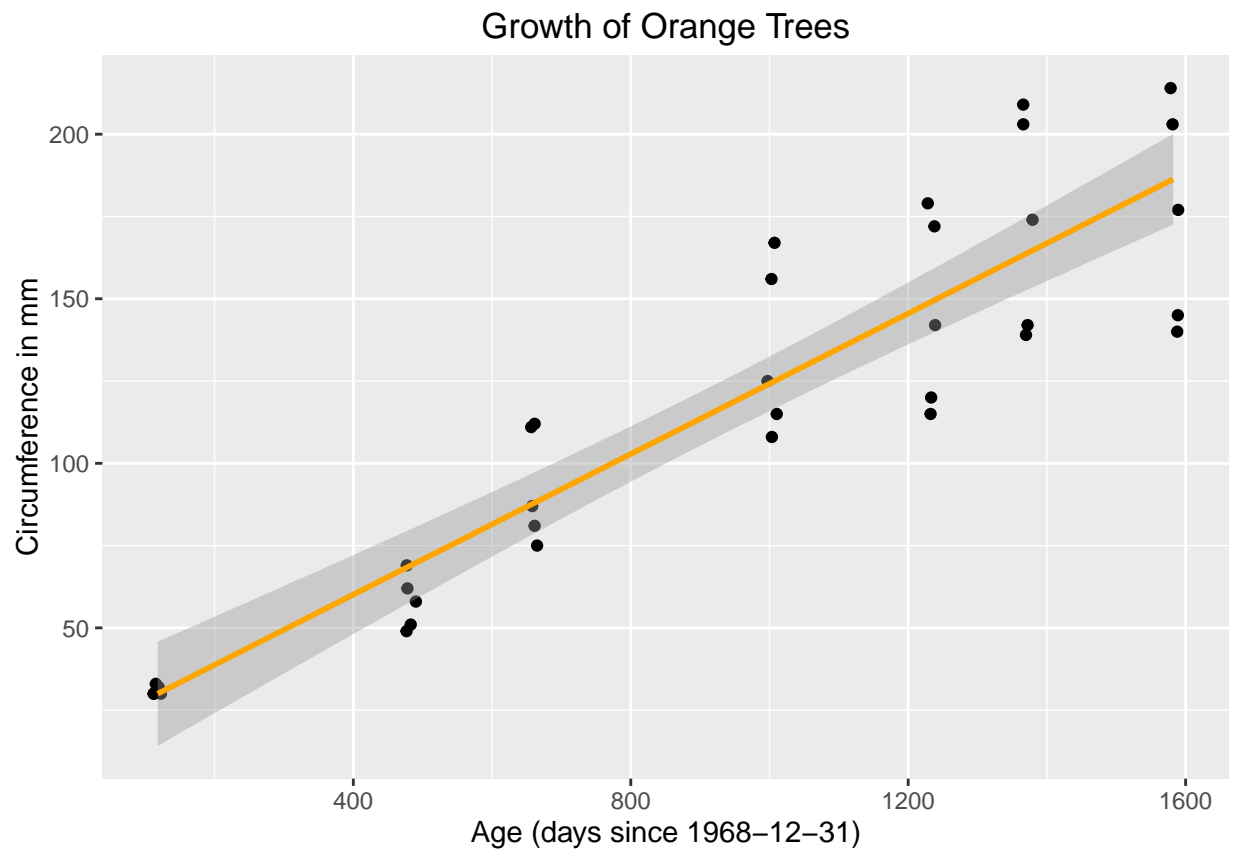
```
ggplot(orange, aes(x = age, y = circumference)) +  
  geom_point()
```



OK. Let's add a linear model to this plot, as well as some better labels, and we'll change from mapping the points as observed to using `geom_jitter` to add a little horizontal (x-axis) jitter to the points, so that we don't have so much overlap.

2.4 Scatterplot with Linear Fit predicting Circumference using Age across all Trees

```
ggplot(orange, aes(x = age, y = circumference)) +  
  geom_jitter(width = 20, height = 0) +  
  geom_smooth(method = "lm", col = "orange") +  
  labs(title = "Growth of Orange Trees",  
        x = "Age (days since 1968-12-31)", y = "Circumference in mm")
```



2.5 The Linear Model

The linear model fitted here is summarized below:

```
model1 <- lm(circumference ~ age, data = orange)
summary(model1)
```

Call:

```
lm(formula = circumference ~ age, data = orange)
```

Residuals:

Min	1Q	Median	3Q	Max
-46.310	-14.946	-0.076	19.697	45.111

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	17.399650	8.622660	2.018	0.0518 .
age	0.106770	0.008277	12.900	1.93e-14 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 23.74 on 33 degrees of freedom

Multiple R-squared: 0.8345, Adjusted R-squared: 0.8295

F-statistic: 166.4 on 1 and 33 DF, p-value: 1.931e-14

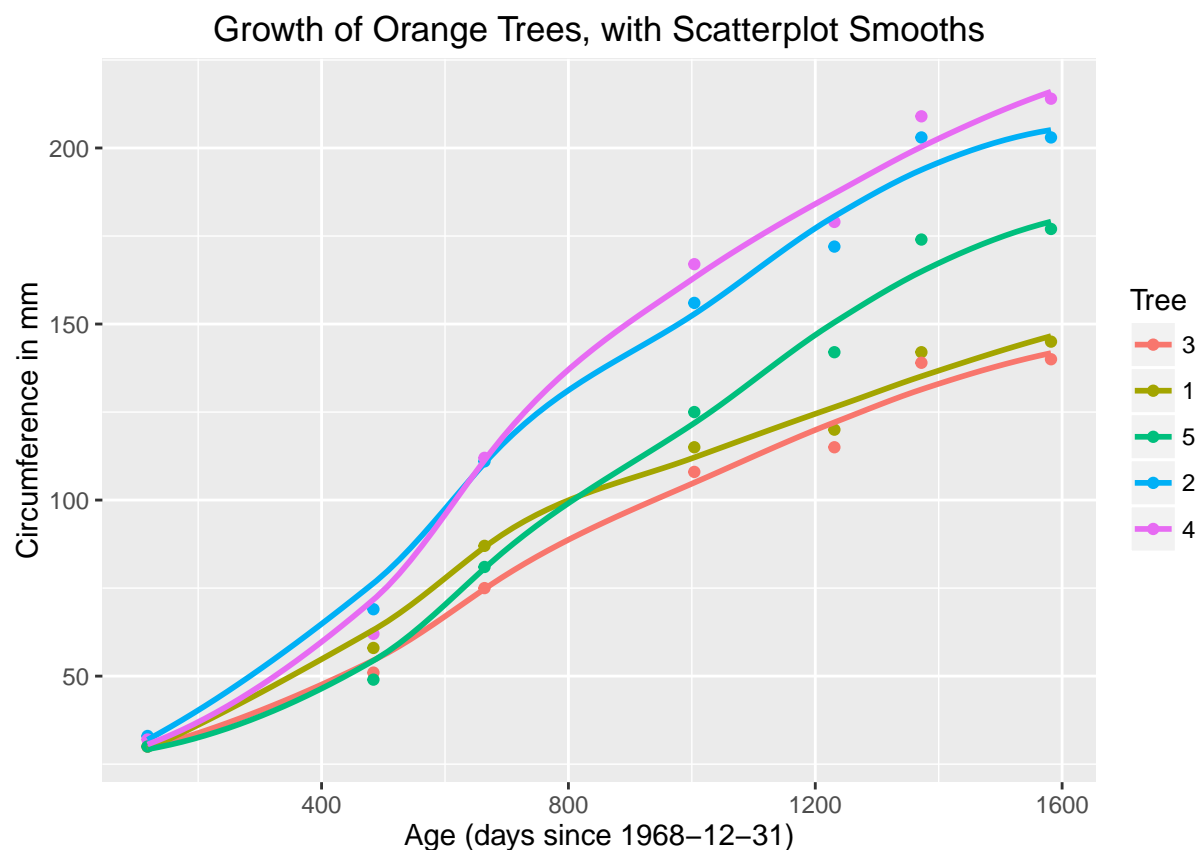
So the linear regression model is: $\text{circumference} = 17.4 + 0.107 \text{ age}$.

So our predicted circumference for a tree of age 1000 days is 124 mm.

2.6 Some Other Scatterplots: Assessing each Tree separately

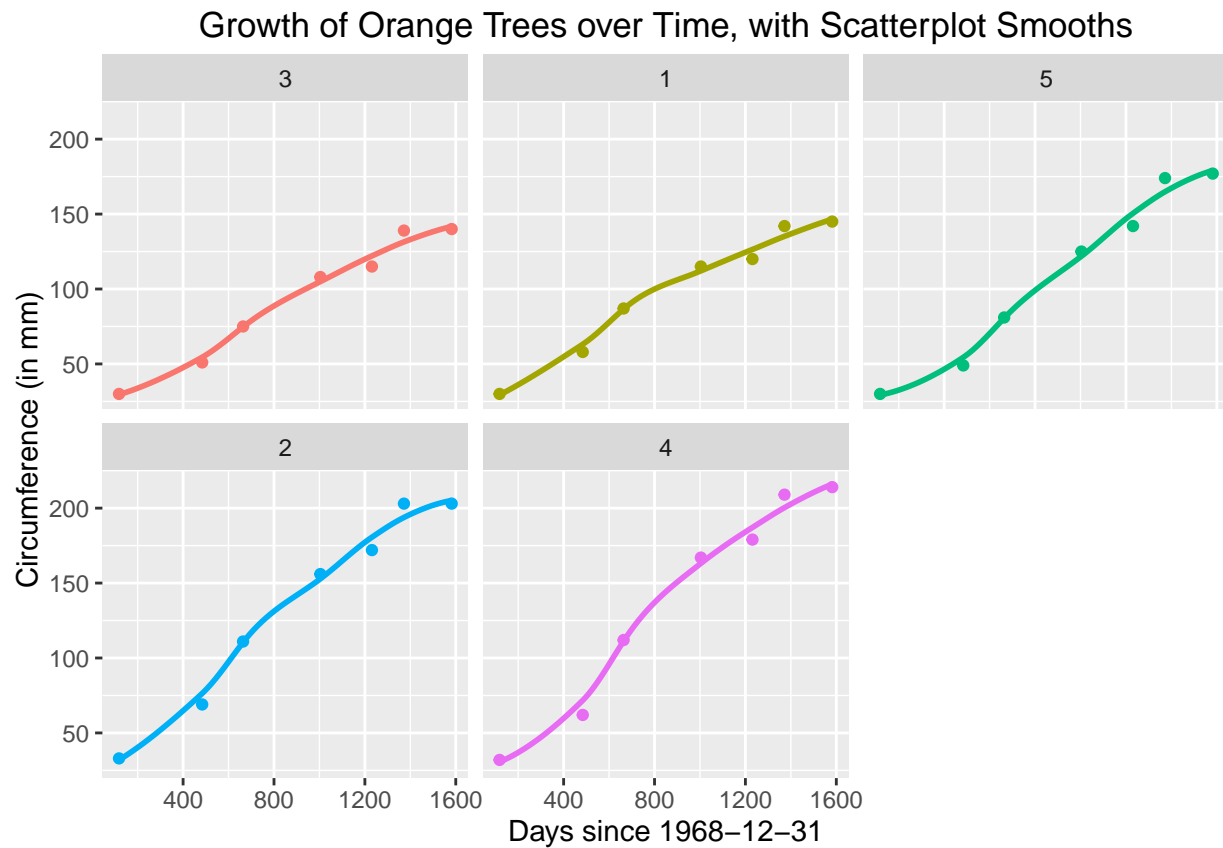
As a third option, let's fit separate smooth (loess) curves to each of the five individual trees, and plot each of them in different colors.

```
ggplot(orange, aes(x = age, y = circumference, col = Tree)) +  
  geom_point() +  
  geom_smooth(se = FALSE) +  
  labs(title = "Growth of Orange Trees, with Scatterplot Smooths",  
        x = "Age (days since 1968-12-31)", y = "Circumference in mm")
```



Or we could facet the plots, showing multiple scatterplots, one for each Tree.

```
ggplot(orange, aes(x = age, y = circumference, col = Tree)) +  
  geom_point() +  
  geom_smooth(se = FALSE) +  
  facet_wrap(~ Tree) +  
  guides(col = FALSE) +  
  labs(title = "Growth of Orange Trees over Time, with Scatterplot Smooths",  
       x = "Days since 1968-12-31", y = "Circumference (in mm)")
```



3 Getting Data into R from Excel or another Software Package: The Fundamentals

The easiest way to get data from another software package into R is to save the file (from within the other software package) in a form that R can read. What you want is to end up with an Excel file that looks like this...

	A	B	C	D
1	patient	drug	gender	response
2	MW	A	M	23
3	TT	B	F	15
4	KH	B	M	18
5	GC	A	M	29
6	DS	B	F	34
7	HJ	B	F	15
8	KM	A	M	7
9	RS	A	M	19
10	DG	A	F	22
11				

Figure 1: An Excel sheet with a tidy data set

This *tidy* data set contains:

- one row for each subject
- variables that indicate characteristics of each of the subjects

The variable names are in the first row, and the data are in the remaining rows (2-10 in this small example). Categorical variables are most easily indicated by letters (drug A or B, for instance) while continuous variables, like response, are indicated by numbers. Leave missing cells blank or use the symbol **NA**, rather than indicating them with, say, -99 or some other code.

Within Excel, this file can be saved as a **.csv** (comma-separated text file) or just as an Excel **.XLS** file, and then imported directly into R, via RStudio by clicking Import Dataset under the Workspace tab, then selecting From Text File. If you've saved the file in Excel as a **.csv** file, RStudio will generally make correct guesses about how to import the file. Once imported, you just need to save the workspace when you quit RStudio and you'll avoid the need to re-import.