

# An Sweave Demo

Charles J. Geyer

June 6, 2012

This is a demo for using the **Sweave** command in R. To get started make a regular  $\text{\LaTeX}$  file (like this one) but give it the suffix **.Rnw** instead of **.tex** and then turn it into a  $\text{\LaTeX}$  file (**foo.tex**) with the (unix) command

```
R CMD Sweave foo.Rnw
```

So you can do

```
latex foo
xdvi foo
```

and so forth.

So now we have a more complicated file chain

$$\text{foo.Rnw} \xrightarrow{\text{Sweave}} \text{foo.tex} \xrightarrow{\text{latex}} \text{foo.dvi} \xrightarrow{\text{xdvi}} \text{view of document}$$

and what have we accomplished other than making it twice as annoying to the WYSIWYG crowd (having to run both **Sweave** and **latex** to get anything that looks like the document)?

Well, we can now include R in our document. Here's a simple example

```
> 2 + 2

[1] 4
```

What I actually typed in **foo.Rnw** was

```
<<two>>=
2 + 2
@
```

This is not  $\text{\LaTeX}$ . It is a “code chunk” to be processed by **Sweave**. When **Sweave** hits such a thing, it processes it, runs R to get the results, and stuffs (by default) the output in the  $\text{\LaTeX}$  file it is creating. The  $\text{\LaTeX}$  between code chunks is copied verbatim (except for **Sexpr**, about which see below). Hence to create a Rnw document you just write plain old  $\text{\LaTeX}$  interspersed with “code chunks” which are plain old R.

Plots get a little more complicated. First we make something to plot (simulate regression data).

```
> n <- 50
> x <- seq(1, n)
> a.true <- 3
> b.true <- 1.5
> y.true <- a.true + b.true * x
> s.true <- 17.3
> y <- y.true + s.true * rnorm(n)
> out1 <- lm(y ~ x)
> summary(out1)
```

Call:

```
lm(formula = y ~ x)
```

Residuals:

Min	1Q	Median	3Q	Max
-42.355	-9.931	0.241	13.755	29.553

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	12.6084	4.7988	2.627	0.0115 *
x	1.3077	0.1638	7.985	2.32e-10 ***

---

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

Residual standard error: 16.71 on 48 degrees of freedom

Multiple R-squared: 0.5705, Adjusted R-squared: 0.5615

F-statistic: 63.76 on 1 and 48 DF, p-value: 2.322e-10

(for once we won't show the code chunk itself, look at `foo.Rnw` if you want to see what the actual code chunk was).

Figure 1 (p. 3) is produced by the following code

```
> plot(x, y)
> abline(out1)
```

Note that `x`, `y`, and `out1` are remembered from the preceding code chunk. We don't have to regenerate them. All code chunks are part of one R "session".

Now this was a little tricky. We did this with two code chunks, one visible and one invisible. First we did

```
<<label=fig1plot,include=FALSE>>=
plot(x, y)
abline(out1)
@
```

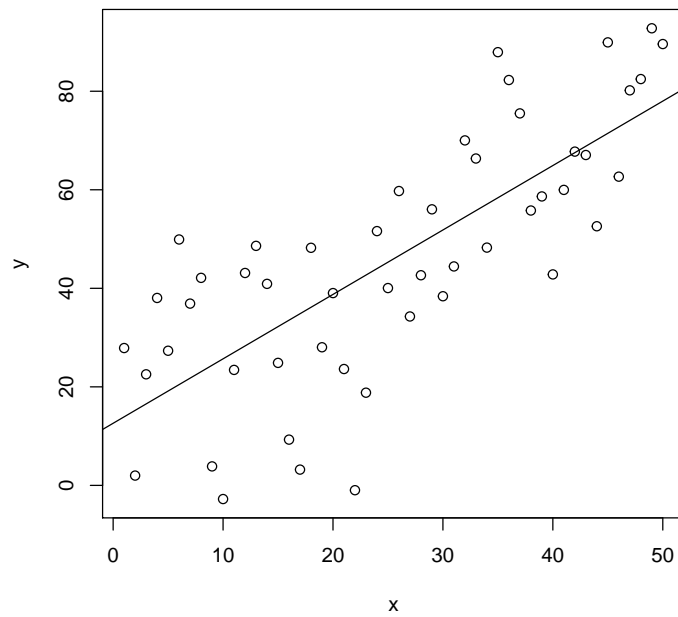


Figure 1: Scatter Plot with Regression Line

where the `include=FALSE` indicates that the output (text and graphics) should not go here (they will be some place else) and the `label=fig1plot` gives the code chunk a name (to be used later). And “later” is almost immediate. Next we did

```
\begin{figure}
\begin{center}
<<label=fig1,fig=TRUE,echo=FALSE>>=
<<fig1plot>>
@
\end{center}
\caption{Scatter Plot with Regression Line}
\label{fig:one}
\end{figure}
```

In this code chunk the `fig=TRUE` indicates that the chunk generates a figure. Sweave automatically makes both EPS and PDF files for the figure and automatically generates an appropriate L<sup>A</sup>T<sub>E</sub>X `\includegraphics` command to include the plot in the `figure` environment. The `echo=FALSE` in the code chunk means just what it says (we’ve already seen the code—it was produced by the preceding chunk—and we don’t want to see it again, especially not in our figure). The `<<fig1plot>>` is an example of “code chunk reuse”. It means that we reuse the code of the code chunk named `fig1plot`. It is important that we observe the DRY/SPOT rule (*don’t repeat yourself* or *single point of truth*) and only have one bit of code for generating the plot. What the reader sees is guaranteed to be the code that made the plot. If we had used cut-and-paste, just repeating the code, the duplicated code might get out of sync after edits. The rest of this should be recognizable to anyone who has ever done a L<sup>A</sup>T<sub>E</sub>X figure.

So making a figure is a bit more complicated in some ways but much simpler in others. Note the following virtues

- The figure is guaranteed to be the one described by the text (at least by the R in the text).
- No messing around with sizing or rotations. It just works!

Note that if you don’t care to show the R code to make the figure, it is simpler still. Figure 2 (p. 5) shows another plot. What I actually typed in `foo.Rnw` was

```
\begin{figure}
\begin{center}
<<label=fig2,fig=TRUE,echo=FALSE>>=
out3 <- lm(y ~ x + I(x^2) + I(x^3))
plot(x, y)
curve(predict(out3, newdata=data.frame(x=x)), add = TRUE)
@
```

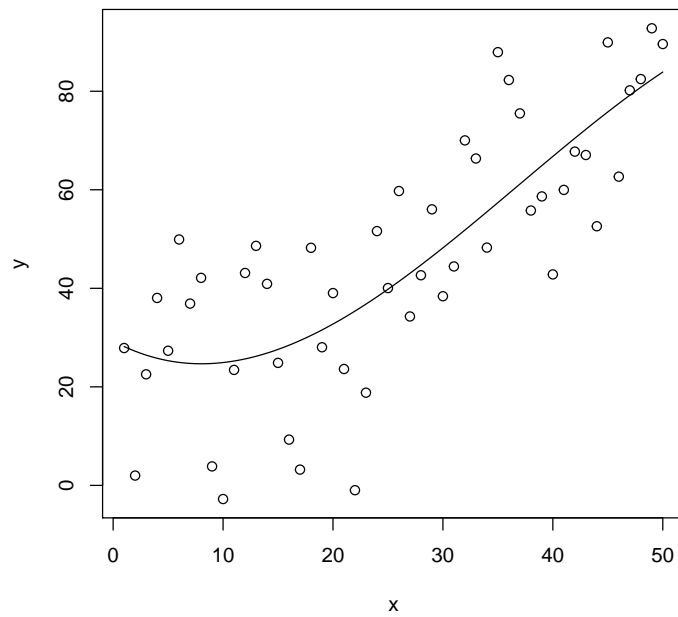


Figure 2: Scatter Plot with Cubic Regression Curve

```

\end{center}
\caption{Scatter Plot with Cubic Regression Curve}
\label{fig:two}
\end{figure}

```

Now we just included the code for the plot in the figure (with `echo=FALSE` so it doesn't show).

Also note that every time we rerun **Sweave** Figures 1 and 2 change, the latter conspicuously (because the simulated data are random). Everything just works. This should tell you the main virtue of Sweave. It's always correct. There is never a problem with stale cut-and-paste.

Simple numbers can be plugged into the text with the `\Sexpr` command, for example, the quadratic and cubic regression coefficients in the preceding regression were  $\beta_2 = 0.0837$  and  $\beta_3 = -0.0008$ . Just magic! What I actually typed in `foo.Rnw` was

```

in the preceding regression
were $\beta_2 = \Sexpr{round(out3$coef[3], 4)}$
and $\beta_3 = \Sexpr{round(out3$coef[4], 4)}$.

```

The `xtable` command is used to make tables. (The following is the **Sweave** of another code chunk that we don't explicitly show. Look at `foo.Rnw` for details.)

```

> out2 <- lm(y ~ x + I(x^2))
> foo <- anova(out1, out2, out3)
> foo

```

#### Analysis of Variance Table

```

Model 1: y ~ x
Model 2: y ~ x + I(x^2)
Model 3: y ~ x + I(x^2) + I(x^3)
  Res.Df  RSS Df Sum of Sq    F Pr(>F)
1     48 13407
2     47 12260   1   1146.62 4.3588 0.04239 *
3     46 12101   1    159.19 0.6052 0.44060
---

```

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

```

> class(foo)

[1] "anova"      "data.frame"

> dim(foo)

[1] 3 6

> foo <- as.matrix(foo)
> foo

```

Table 1: ANOVA Table

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	48	13406.63				
2	47	12260.01	1	1146.62	4.359	0.042
3	46	12100.82	1	159.19	0.605	0.441

```

Res.Df      RSS Df Sum of Sq      F      Pr(>F)
1      48 13406.63 NA      NA      NA      NA
2      47 12260.01 1 1146.6190 4.358753 0.04238526
3      46 12100.82 1  159.1918 0.605151 0.44060232

```

So now we are ready to turn the matrix `foo` into Table 1 using the R chunk

```

<<label=tab1,echo=FALSE,results=tex>>=
library(xtable)
print(xtable(foo, caption = "ANOVA Table", label = "tab:one",
  digits = c(0, 0, 2, 0, 2, 3, 3)), table.placement = "tbp",
  caption.placement = "top")
@

```

(note the difference between arguments to the `xtable` function and to the `xtable` method of the `print` function).

To summarize, **Sweave** is terrific, so important that soon we'll not be able to get along without it. It's virtues are

- The numbers and graphics you report are actually what they are claimed to be.
- Your analysis is reproducible. Even years later, when you've completely forgotten what you did, the whole write-up, every single number or pixel in a plot is reproducible.
- Your analysis actually works—at least in this particular instance. The code you show actually executes without error.
- Toward the end of your work, with the write-up almost done you discover an error. Months of rework to do? No! Just fix the error and rerun **Sweave** and **latex**. One single problem like this and you will have all the time invested in **Sweave** repaid.
- This methodology provides discipline. There's nothing that will make you clean up your code like the prospect of actually revealing it to the world.

Whether we're talking about homework, a consulting report, a textbook, or a research paper. If they involve computing and statistics, this is the way to do it.