

CSSS508, Week 10

Model Results and Reproducibility

Chuck Lanfear

Jun 3, 2019

Updated: Mar 29, 2020



Topics for Today

Working with Model Results

- Tidy model output with `broom`
- Visualizing models with `ggeffects`
- Making regression tables

Reproducible Research

Best Practices

- Organization
- Portability
- Version Control

Wrapping up the course

Working with Model Results

broom

`broom` is a package that "tidies up" the output from models such as `lm()` and `glm()`.

It has a small number of key functions:

- `tidy()` - Creates a dataframe summary of a model.
- `augment()` - Adds columns—such as fitted values—to the data used in the model.
- `glance()` - Provides one row of fit statistics for models.

```
library(broom)
```

Model Output is a List

`lm()` and `summary()` produce lists as output, which cannot go directly into tidyverse functions, particularly those in `ggplot2`.

```
lm_1 <- lm(yn ~ num1 + fac1, data = ex_dat)
summary(lm_1)
```

```
##
## Call:
## lm(formula = yn ~ num1 + fac1, data = ex_dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.198 -2.231 -0.235  1.911  7.386
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.503      0.351    4.28 2.9e-05 ***
## num1           0.584      0.103    5.66 5.2e-08 ***
## fac1B          1.205      0.503    2.40  0.018 *
## fac1C          1.142      0.501    2.28  0.024 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.94 on 196 degrees of freedom
## Multiple R-squared:  0.162,    Adjusted R-squared:  0.149
## F-statistic: 12.6 on 3 and 196 DF,  p-value: 1.43e-07
....
```

Model Output Varies!

Each type of model also produces somewhat different output, so you can't just reuse the same code to handle output from every model.

```
glm_1 <- glm(yb ~ num1 + fac1, data = ex_dat, family=binomial(link="logit"))
summary(glm_1)
```

```
##
## Call:
## glm(formula = yb ~ num1 + fac1, family = binomial(link = "logit"),
##      data = ex_dat)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.801  -1.077  -0.589   1.065   1.836
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.6306     0.2595  -2.43  0.01511 *
## num1          0.3069     0.0797   3.85  0.00012 ***
## fac1B         0.3626     0.3582   1.01  0.31137
## fac1C         0.4828     0.3605   1.34  0.18048
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
.....
```

broom::tidy()

`tidy()` produces the similar output, but as a dataframe.

```
lm_1 %>% tidy()
```

```
## # A tibble: 4 x 5
##   term          estimate std.error statistic    p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)    1.50      0.351     4.28 0.0000286
## 2 num1           0.584     0.103     5.66 0.0000000521
## 3 fac1B          1.21      0.503     2.40 0.0175
## 4 fac1C          1.14      0.501     2.28 0.0237
```

Each type of model (e.g. `glm`, `lmer`) has a different *method* with its own additional arguments. See `?tidy.lm` for an example.

broom::tidy()

This output is also completely identical between different models.

This can be very useful and important if running models with different test statistics... or just running a lot of models!

```
glm_1 %>% tidy()
```

```
## # A tibble: 4 x 5
##   term          estimate std.error statistic  p.value
##   <chr>         <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)  -0.631     0.260     -2.43  0.0151
## 2 num1         0.307     0.0797     3.85  0.000119
## 3 fac1B        0.363     0.358     1.01  0.311
## 4 fac1C        0.483     0.360     1.34  0.180
```


broom::glance()

`glance()` produces dataframes of fit statistics for models.

If you run many models, you can compare each model row-by-row in each column... or even plot their different fit statistics to allow holistic comparison.

```
glance(lm_1)
```

```
## # A tibble: 1 x 11
##   r.squared adj.r.squared sigma statistic p.value    df logLik   AIC    BIC deviance
##   <dbl>      <dbl> <dbl>      <dbl>  <dbl> <int>  <dbl> <dbl> <dbl>    <dbl>
## 1    0.162        0.149   2.94        12.6 1.43e-7     4  -497. 1004. 1021.    1688.
## # ... with 1 more variable: df.residual <int>
```

broom augment()

`augment()` takes values generated by a model and adds them back to the original data. This includes fitted values, residuals, and leverage statistics.

```
augment(lm_1) %>% head()
```

```
## # A tibble: 6 x 10
##       yn      num1 fac1  .fitted .se.fit .resid  .hat .sigma  .cooksd .std.resid
##   <dbl>   <dbl> <fct>   <dbl>   <dbl> <dbl>  <dbl> <dbl>   <dbl>      <dbl>
## 1 -1.52    3.44    A      3.51    0.396 -5.03  0.0182  2.92  0.0139    -1.73
## 2 -0.980    1.69    A      2.49    0.329 -3.47  0.0126  2.93  0.00452   -1.19
## 3  4.90   -0.647    C      2.27    0.413  2.64  0.0198  2.94  0.00415    0.907
## 4  6.05    5.53    B      5.94    0.603  0.114  0.0422  2.94  0.0000173  0.0396
## 5  2.77    0.00818 C      2.65    0.391  0.119  0.0177  2.94  0.00000754 0.0409
## 6  3.01   -1.92    A      0.381   0.462  2.63  0.0248  2.94  0.00524    0.908
```

The Power of broom

The real advantage of `broom` becomes apparent when running many models at once. Here we run separate models for each level of `fac1`:

```
ex_dat %>% group_by(fac1) %>% do(tidy(lm(yn ~ num1 + fac2 + num2, data = .)))
```

```
## # A tibble: 12 x 6
## # Groups:   fac1 [3]
##   fac1 term          estimate std.error statistic  p.value
##   <fct> <chr>          <dbl>      <dbl>      <dbl>    <dbl>
## 1 A     (Intercept)    1.09      0.362      3.02 3.42e- 3
## 2 A     num1           0.423     0.111      3.81 2.73e- 4
## 3 A     fac2No         0.903     0.478      1.89 6.24e- 2
## 4 A     num2           0.705     0.0686     10.3 4.39e-16
## 5 B     (Intercept)    1.63      0.461      3.54 8.12e- 4
## 6 B     num1           0.748     0.161      4.64 2.23e- 5
## 7 B     fac2No         0.638     0.549      1.16 2.50e- 1
## 8 B     num2           0.662     0.0889      7.45 6.80e-10
## 9 C     (Intercept)    3.10      0.288     10.8 2.94e-15
## 10 C    num1           0.577     0.0968      5.95 1.80e- 7
## 11 C    fac2No         0.0877     0.415      0.211 8.33e- 1
## 12 C    num2           0.750     0.0833      9.00 1.79e-12
```

`do()` repeats whatever is inside it once for each level of the variable(s) in `group_by()` then puts them together as a data frame.

Plotting Model Results

geom_smooth()

I have used `geom_smooth()` in many past examples.

`geom_smooth()` generates "smoothed conditional means" including loess curves and generalized additive models (GAMs).

Note, however, that most regression models are conditional mean models, such as ordinary least squares, generalized linear models.

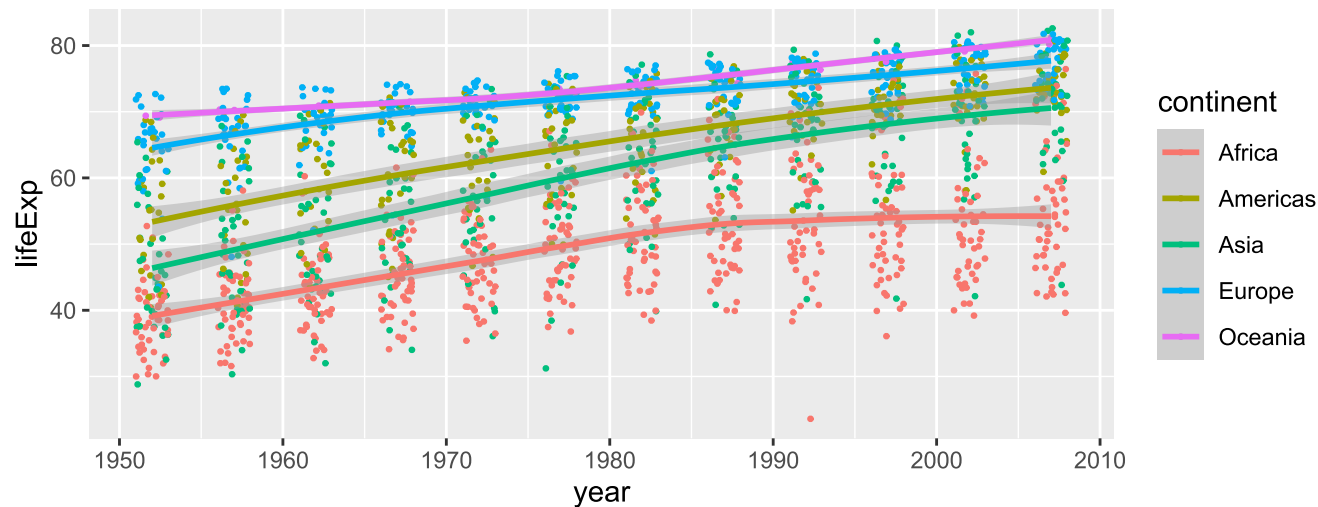
We can use `geom_smooth()` to add a layer depicting common bivariate models.

We'll look at this with the `gapminder` data from Week 2.

```
library(gapminder)
```

Default `geom_smooth()`

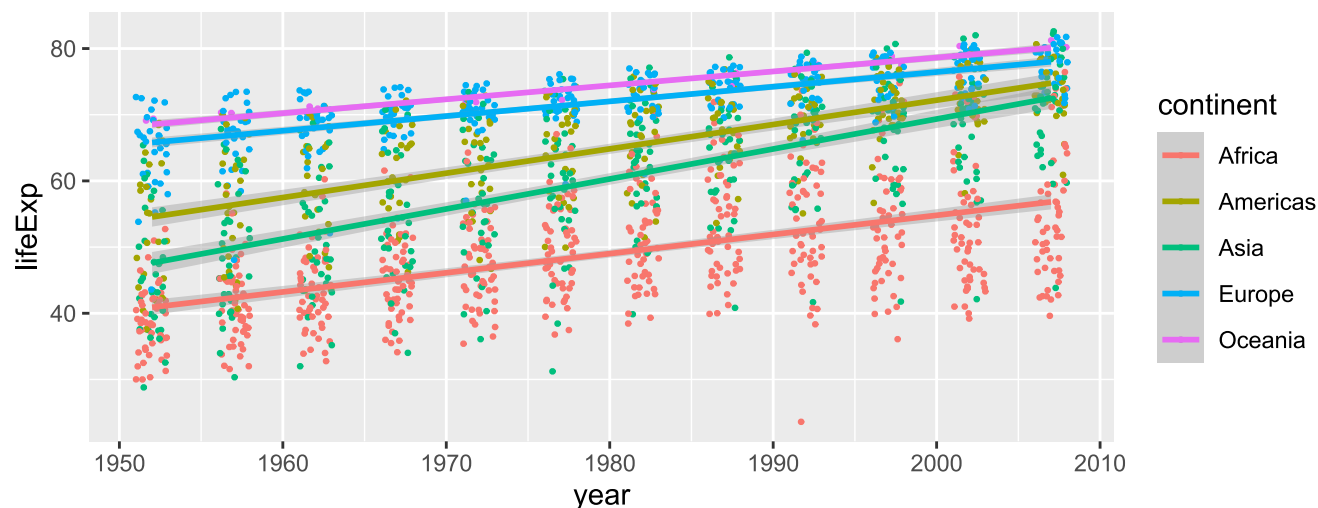
```
ggplot(data = gapminder,  
       aes(x = year, y = lifeExp, color = continent)) +  
  geom_point(position = position_jitter(1,0), size = 0.5) +  
  geom_smooth()
```



By default, `geom_smooth()` chooses either a loess smoother ($N < 1000$) or a GAM depending on the number of observations.

Linear glm

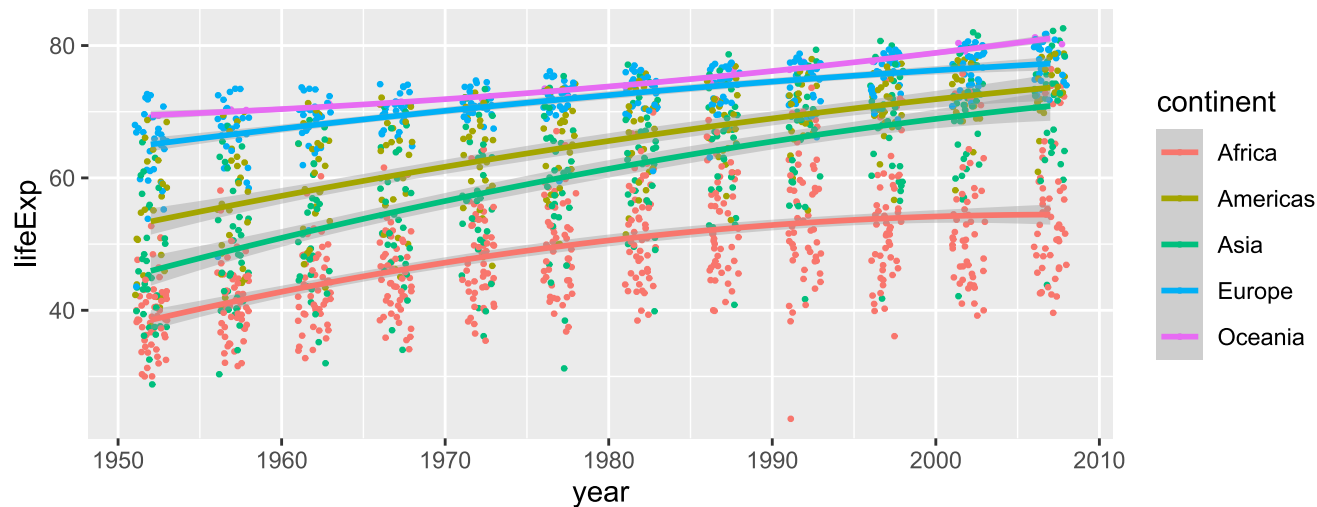
```
ggplot(data = gapminder,  
       aes(x = year, y = lifeExp, color = continent)) +  
  geom_point(position = position_jitter(1,0), size = 0.5) +  
  geom_smooth(method = "glm", formula = y ~ x)
```



We could also fit a standard linear model using either `method = "glm"` or `method = "lm"` and a formula like `y ~ x`.

Polynomial glm

```
ggplot(data = gapminder,  
       aes(x = year, y = lifeExp, color = continent)) +  
  geom_point(position = position_jitter(1,0), size = 0.5) +  
  geom_smooth(method = "glm", formula = y ~ poly(x, 2))
```



`poly(x, 2)` produces a quadratic model which contains a linear term (`x`) and a quadratic term (`x2`).

More Complex Models

What if we want something more complex than a bivariate model?

What if we have a statistically complex model, like nonlinear probability model or multilevel model?

We need to go beyond `geom_smooth()`!

But first, vocab!

We are often interested in what might happen if some variables take particular values, often ones not seen in the actual data.

When we set variables to certain values, we refer to them as **counterfactual values** or just **counterfactuals**.

For example, if we know nothing about a new observation, our prediction for that estimate is often based on assuming every variable is at its mean.

Sometimes, however, we might have very specific questions which require setting (possibly many) combinations of variables to particular values and making an estimate or prediction.

Providing specific estimates, conditional on values of covariates, is a nice way to summarize results, particularly for models with unintuitive parameters (e.g. logit models).

ggeffects

ggeffects

If we want to look at more complex models, we can use `ggeffects` to create and plot tidy *marginal effects*.

That is, tidy dataframes of *ranges* of predicted values that can be fed straight into `ggplot2` for plotting model results.

We will focus on two `ggeffects` functions:

- `ggpredict()` - Computes predicted values for the outcome variable at margins of specific variables.
- `plot.ggeffects()` - A plot method for `ggeffects` objects (like `ggredict()` output)

```
library(ggeffects)
```

Quick Simulated Data

To best show off `ggeffects`, I need a data frame with numeric and categorical variables with strong relationships. It is easiest to just simulate it:

```
ex_dat <- data.frame(num1 = rnorm(200, 1, 2),  
                     fac1 = sample(c(1, 2, 3), 200, TRUE),  
                     num2 = rnorm(200, 0, 3),  
                     fac2 = sample(c(1, 2))) %>%  
  mutate(yn = num1 * 0.5 + fac1 * 1.1 + num2 * 0.7 +  
          fac2 - 1.5 + rnorm(200, 0, 2)) %>%  
  mutate(yb = as.numeric(yn > mean(yn))) %>%  
  mutate(fac1 = factor(fac1, labels = c("A", "B", "C")),  
         fac2 = factor(fac2, labels = c("Yes", "No")))
```

Now we can get `ggpredicting`!

ggpredict()

When you run `ggpredict()`, it produces a dataframe with a row for every unique value of a supplied predictor ("independent") variable (`term`).

Each row contains an expected (estimated) value for the outcome ("dependent") variable, plus confidence intervals.

```
lm_1 <- lm(yn ~ num1 + fac1, data = ex_dat)
lm_1_est <- ggpredict(lm_1, terms = "num1")
```

If desired, the argument `interval="prediction"` will give predicted intervals instead.

ggpredict() output

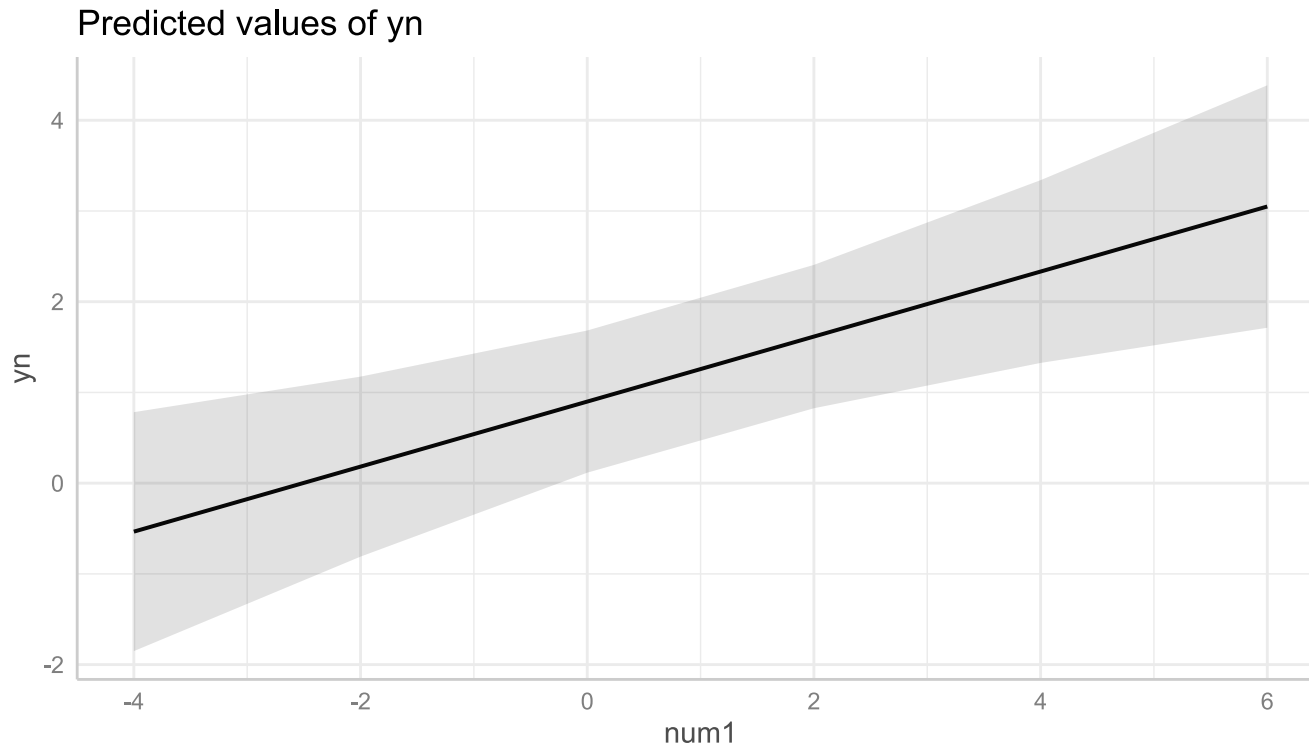
```
lm_1_est
```

```
##  
## # Predicted values of yn  
## # x = num1  
##  
## x | Predicted | SE | 95% CI  
## -----  
## -4 | -0.53 | 0.67 | [-1.85, 0.78]  
## -2 | 0.18 | 0.51 | [-0.81, 1.17]  
## 0 | 0.90 | 0.40 | [ 0.12, 1.68]  
## 2 | 1.62 | 0.40 | [ 0.83, 2.41]  
## 4 | 2.33 | 0.51 | [ 1.33, 3.34]  
## 6 | 3.05 | 0.68 | [ 1.71, 4.39]  
##  
## Adjusted for:  
## * fac1 = A
```

plot() for ggpredict()

ggeffects features a `plot()` method, `plot.ggeffects()`, which produces a ggplot when you give `plot()` output from `ggpredict()`.

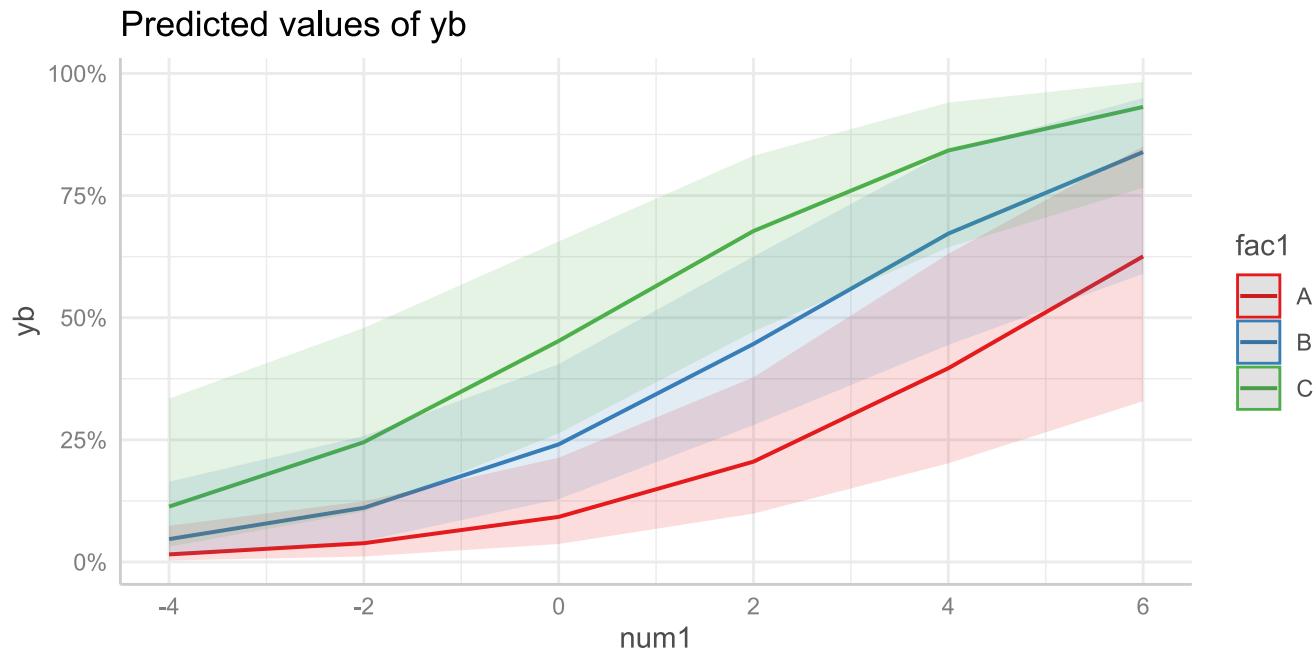
```
plot(lm_1_est)
```



Grouping with `ggpredict()`

When using a vector of `terms`, `ggeffects` will plot the first along the x-axis and use others for *grouping*. Note we can pipe a model into `ggpredict()`!

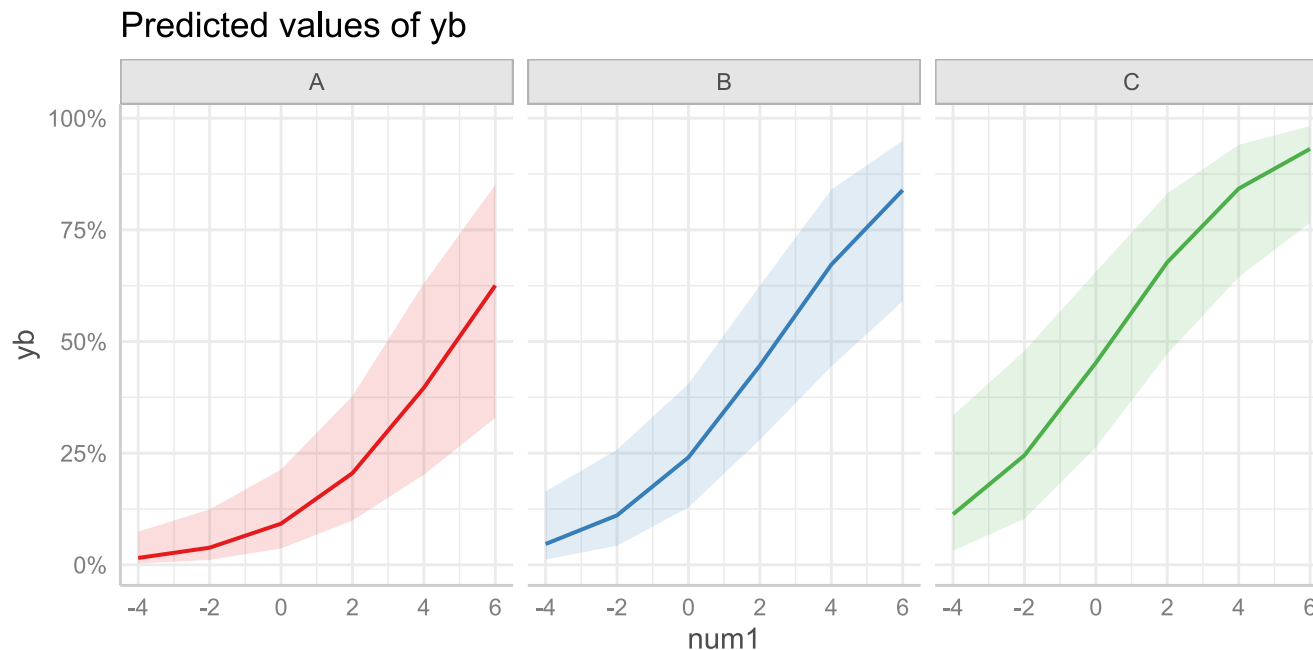
```
glm(yb ~ num1 + fac1 + num2 + fac2, data = ex_dat, family=binomial(link = "logit")) %>%  
  ggpredict(terms = c("num1", "fac1")) %>% plot()
```



Faceting with `ggpredict()`

You can add `facet=TRUE` to the `plot()` call to facet over *grouping terms*.

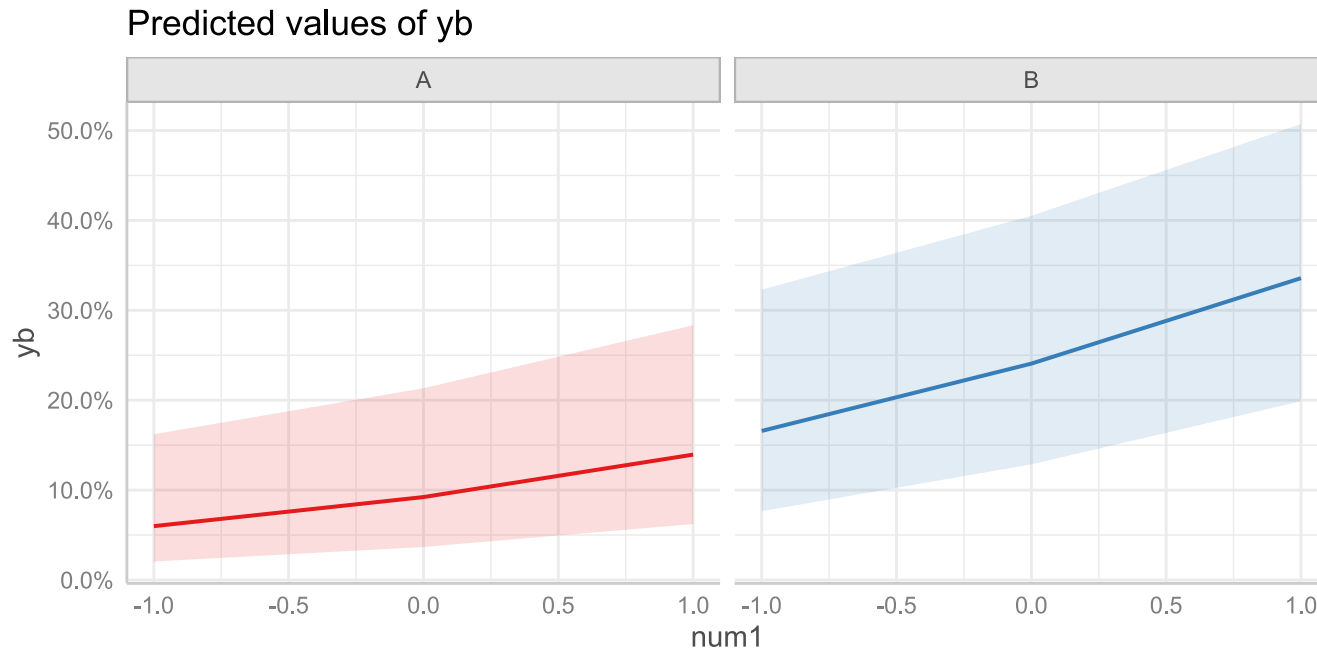
```
glm(yb ~ num1 + fac1 + num2 + fac2, data = ex_dat, family = binomial(link = "logit")) %>%  
  ggpredict(terms = c("num1", "fac1")) %>% plot(facet=TRUE)
```



Counterfactual Values

You can add values in square brackets in the `terms=` argument to specify counterfactual values.

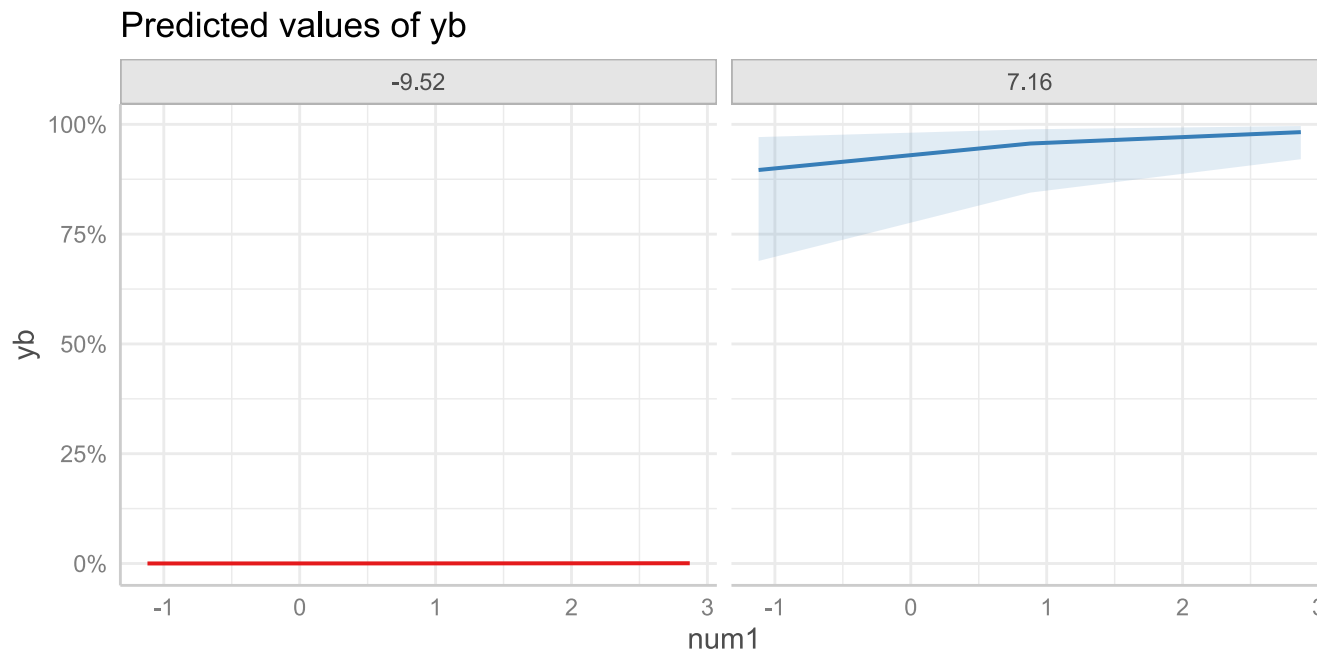
```
glm(yb ~ num1 + fac1 + num2 + fac2, data=ex_dat, family=binomial(link="logit")) %>%  
  ggpredict(terms = c("num1 [-1,0,1]", "fac1 [A,B]")) %>% plot(facet=TRUE)
```



Representative Values

You can also use `[meansd]` or `[minmax]` to set representative values.

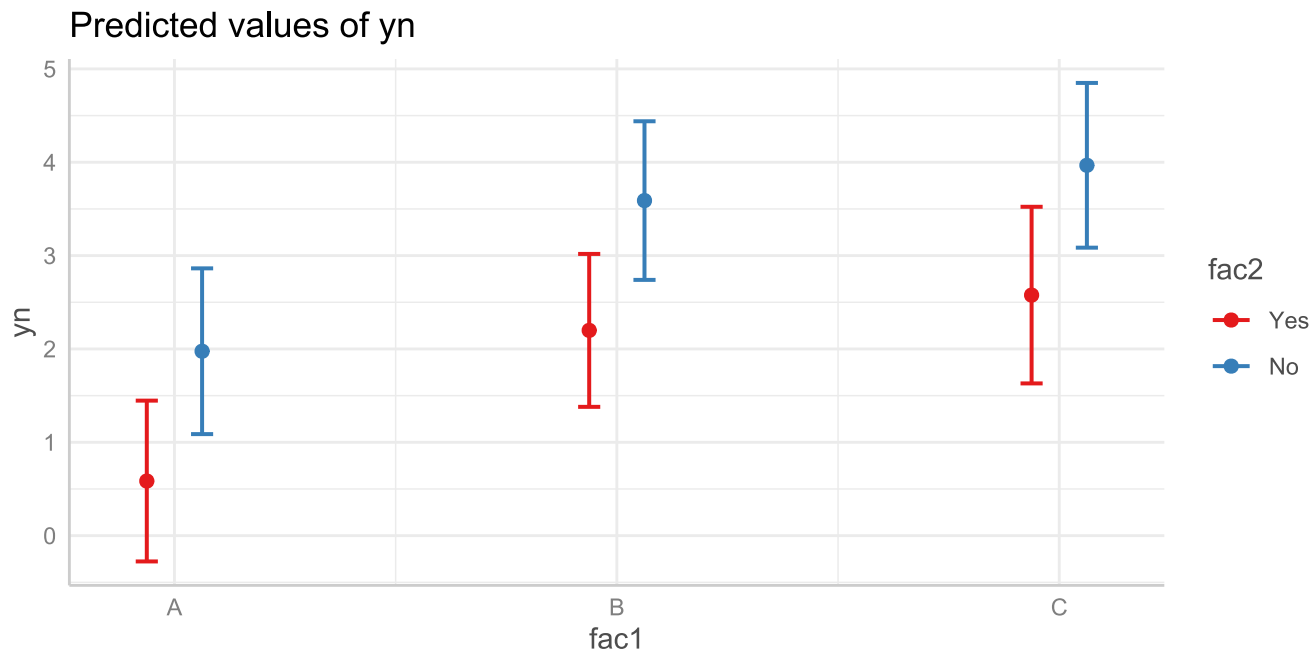
```
glm(yb ~ num1 + fac1 + num2 + fac2, data = ex_dat, family = binomial(link = "logit")) %>%  
  ggpredict(terms = c("num1 [meansd]", "num2 [minmax]")) %>% plot(facet=TRUE)
```



Dot plots with `ggpredict()`

`ggpredict` will produce dot plots with error bars for categorical predictors.

```
lm(yn ~ fac1 + fac2, data = ex_dat) %>%  
  ggpredict(terms=c("fac1", "fac2")) %>% plot()
```



Notes on ggeffects

There is a lot more to the `ggeffects` package that you can see in [the package vignette](#) and the [github repository](#). This includes, but is not limited to:

- Predicted values for polynomial and interaction terms
- Getting predictions from models from dozens of other packages
- Sending `ggeffects` objects to `ggplot2` to freely modify plots

If you need to do something more complex then `ggeffects` allows, see the [Advanced Counterfactuals](#) slides here.

Making Tables

pander Regression Tables

We've used `pander` to create nice tables for dataframes. But `pander` has *methods* to handle all sort of objects that you might want displayed nicely.

This includes model output, such as from `lm()`, `glm()`, and `summary()`.

```
library(pander)
```


pander() and lm()

You can send an `lm()` object straight to `pander`:

```
pander(lm_1)
```

	Estimate	Std. Error	t value	Pr(>t)
(Intercept)	37.23	1.599	23.28	2.565e-20
wt	-3.878	0.6327	-6.129	1.12e-06
hp	-0.03177	0.00903	-3.519	0.001451

Table: Fitting linear model: `mpg ~ wt + hp`

pander() and summary()

You can do this with `summary()` as well, for added information:

```
pander(summary(lm_1))
```

	Estimate	Std. Error	t value	Pr(>t)
(Intercept)	37.23	1.599	23.28	2.565e-20
wt	-3.878	0.6327	-6.129	1.12e-06
hp	-0.03177	0.00903	-3.519	0.001451
Observations	Residual Std. Error	R^2	Adjusted R^2	
32	2.593	0.8268	0.8148	

Table: Fitting linear model: $\text{mpg} \sim \text{wt} + \text{hp}$

sjPlot

`pander` tables are great for basic `rmarkdown` documents, but they're not generally publication ready.

The `sjPlot` package produces `html` tables that look more like those you may find in journal articles.

```
library(sjPlot)
```

```
## Install package "strengjacke" from GitHub (`devtools::install_github("str"
```

sjPlot Tables

`tab_model()` will produce tables for most models.

```
model_1 <- lm(mpg ~ wt, data = mtcars)
tab_model(model_1)
```

	mpg		
	<i>B</i>	<i>CI</i>	<i>p</i>
(Intercept)	37.29	33.45 – 41.12	<.001
wt	-5.34	-6.49 – -4.20	<.001
Observations	32		
R ² / adj. R ²	.753 / .745		

Multi-Model Tables with `sjTable`

Often in journal articles you will see a single table that compares multiple models.

Typically, authors will start with a simple model on the left, then add variables, until they have their most complex model on the right.

The `sjPlot` package makes this easy to do: just give `tab_model()` more models!

Multiple `tab_model()`

```
model_2 <- lm(mpg ~ hp + wt, data = mtcars)
model_3 <- lm(mpg ~ hp + wt + factor(am), data = mtcars)
tab_model(model_1, model_2, model_3)
```

	mpg			mpg			mpg		
	<i>B</i>	<i>CI</i>	<i>p</i>	<i>B</i>	<i>CI</i>	<i>p</i>	<i>B</i>	<i>CI</i>	<i>p</i>
(Intercept)	37.29	33.45 – 41.12	<.001	37.23	33.96 – 40.50	<.001	34.00	28.59 – 39.42	<.001
wt	-5.34	-6.49 – -4.20	<.001	-3.88	-5.17 – -2.58	<.001	-2.88	-4.73 – -1.02	.004
hp				-0.03	-0.05 – -0.01	.001	-0.04	-0.06 – -0.02	<.001
factor(am) (1)							2.08	-0.74 – 4.90	.141
Observations	32			32			32		
R ² / adj. R ²	.753 / .745			.827 / .815			.840 / .823		

sjPlot does a lot more

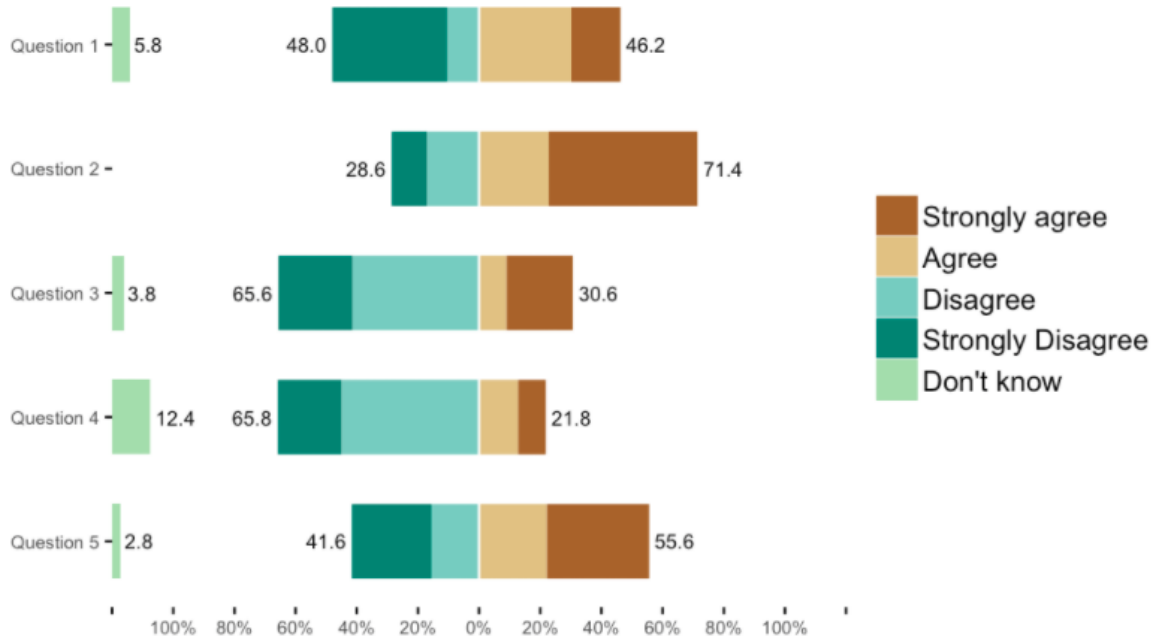
The `sjPlot` package does *a lot* more than just make pretty tables. It is a rabbit hole of *incredibly* powerful and useful functions for displaying descriptive and inferential results.

View the [package website](#) for extensive documentation.

`sjPlot` is a bit more complicated than `ggeffects` but can do just about everything it can do as well; they were written by the same author!

`sjPlot` is fairly new but offers a fairly comprehensive solution for `ggplot` based publication-ready social science data visualization. All graphical functions in `sjPlot` are based on `ggplot2`, so it should not take terribly long to figure out.

sjPlot Example: Likert plots



sjPlot Example: Crosstabs

<i>elder's dependency</i>	<i>carer's level of education</i>			<i>Total</i>
	low level of education	intermediate level of education	high level of education	
independent	21	76	10	107
	19.6 %	71 %	9.3 %	100 %
	1.4 %	5.1 %	0.7 %	7.2 %
slightly dependent	72	238	68	378
	19 %	63 %	18 %	100 %
	4.9 %	16.1 %	4.6 %	25.6 %
moderately dependent	106	289	103	498
	21.3 %	58 %	20.7 %	100 %
	7.2 %	19.5 %	7 %	33.7 %
severely dependent	118	296	84	498
	23.7 %	59.4 %	16.9 %	100 %
	8 %	20 %	5.7 %	33.7 %
<i>Total</i>	317	899	265	1481
	21.5 %	60.7 %	18 %	100 %
	21.5 %	60.7 %	18 %	100 %

$\chi^2=8.658 \cdot df=6 \cdot \Phi_C=.072 \cdot p=.194$

LaTeX Tables

For tables in *L^AT_EX*—as is needed for `.pdf` files—I recommend looking into the `gt`, `stargazer`, or `kableExtra` packages.

`gt` and `kableExtra` allow the construction of complex tables in either HTML or *L^AT_EX* using additive syntax similar to `ggplot2` and `dplyr`.

`stargazer` produces nicely formatted *L^AT_EX* tables but is idiosyncratic.

If you want to edit *L^AT_EX* documents, you can do it in R using Sweave documents (`.Rnw`). Alternatively, you may want to work in a dedicated *L^AT_EX* editor. I recommend [Overleaf](#) for this purpose.

RMarkdown has support for a fair amount of basic *L^AT_EX* syntax if you aren't trying to get too fancy!

Another approach I have used is to manually format *L^AT_EX* tables but use in-line R calls to fill in the values dynamically. This gets you the *exact* format you want but without forcing you to update values any time something changes.

Bonus: `corrplot`

The `corrplot` package has functions for displaying correlograms.

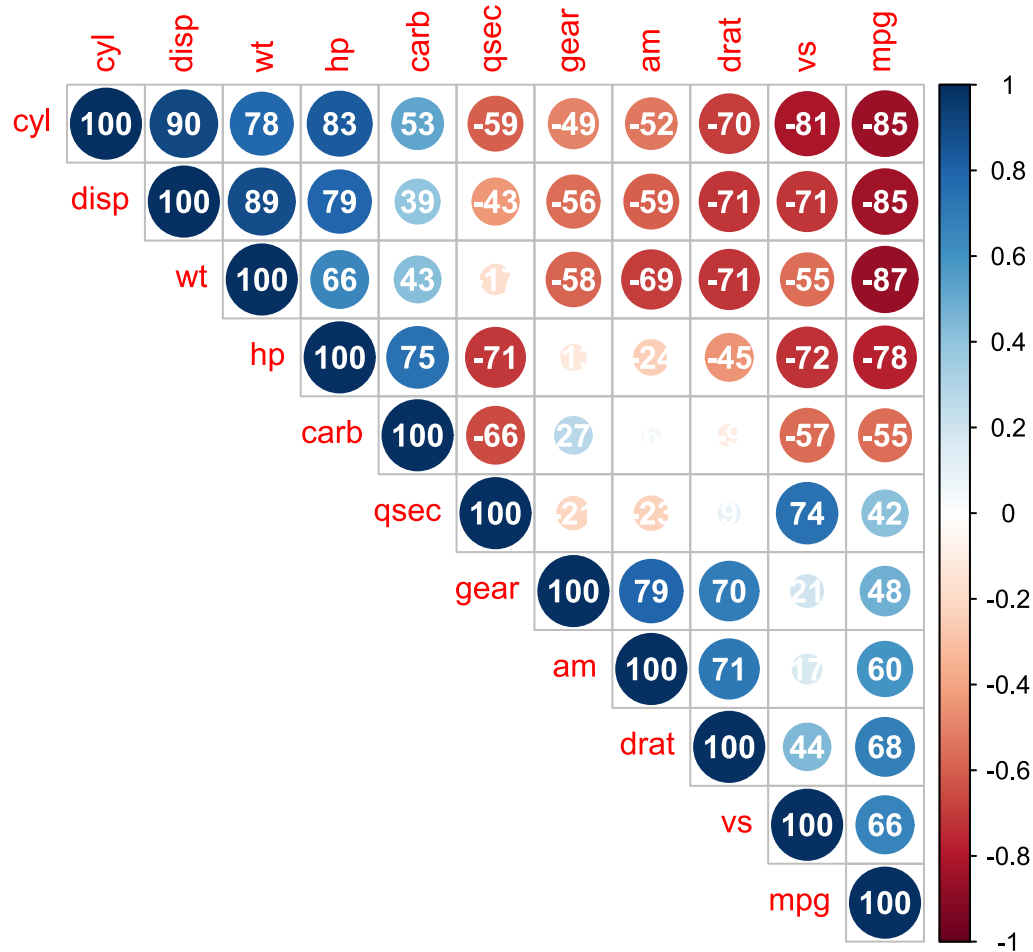
These make visualizing the correlations between variables in a data set easier.

The first argument is a call to `cor()`, the base R function for generating a correlation matrix.

[See the vignette for customization options.](#)

```
library(corrplot)
corrplot(
  cor(mtcars),
  addCoef.col = "white",
  addCoefasPercent=T,
  type="upper",
  order="AOE")
```

Correlogram



Reproducible Research

Why Reproducibility?

Reproducibility is not *replication*.

- **Replication** is running a new study to show if and how results of a prior study hold.
- **Reproducibility** is about rerunning *the same study* and getting the *same results*.

Reproducible studies can still be *wrong*... and in fact reproducibility makes proving a study wrong *much easier*.

Reproducibility means:

- Transparent research practices.
- Minimal barriers to verifying your results.

Any study that isn't reproducible can be trusted only on faith.

Reproducibility Definitions

Reproducibility comes in three forms (Stodden 2014):

1. **Empirical:** Repeatability in data collection.
2. **Statistical:** Verification with alternate methods of inference.
3. **Computational:** Reproducibility in cleaning, organizing, and presenting data and results.

R is particularly well suited to enabling **computational reproducibility**.¹

They will not fix flawed research design, nor offer a remedy for improper application of statistical methods.

Those are the difficult, non-automatable things you want skills in.

[1] Python is equally well suited.

Computational Reproducibility

Elements of computational reproducibility:

- Shared data
 - Researchers need your original data to verify and replicate your work.
- Shared code
 - Your code must be shared to make decisions transparent.
- Documentation
 - The operation of code should be either self-documenting or have written descriptions to make its use clear.
- **Version Control**
 - Documents the research process.
 - Prevents losing work and facilitates sharing.

Levels of Reproducibility

For academic papers, degrees of reproducibility vary:

1. "Read the article"
2. Shared data with documentation
3. Shared data and all code
4. **Interactive document**
5. **Research compendium**
6. Docker compendium: Self-contained ecosystem

Interactive Documents

Interactive documents—like R Markdown docs—combine code and text together into a self-contained document.

- Load and process data
- Run models
- Generate tables and plots in-line with text
- In-text values automatically filled in

Interactive documents allow a reader to examine your computational methods within the document itself; in effect, they are self-documenting.

By re-running the code, they reproduce your results on demand.

Common Platforms:

- **R:** R Markdown ([an example of mine](#))
- **Python:** Jupyter Notebooks

Research Compendia

A **research compendium** is a portable, reproducible distribution of an article or other project.

Research compendia feature:

- An interactive document as the foundation
- Files organized in a recognizable structure (e.g. an R package)
- Clear separation of data, method, and output. *Data are read only.*
- A well-documented or even *preserved* computational environment (e.g. Docker)

`rrtools` by UW's [Ben Markwick](#) provides a simplified workflow to accomplish this in R.

[Here is an example compendium of mine.](#)

Bookdown

`bookdown`—which is integrated into `rrtools`—can generate documents in the proper format for articles, theses, books, or dissertations.

`bookdown` provides an accessible alternative to writing *L^AT_EX* for typesetting and reference management.

You can integrate citations and automate reference page generation using bibtex files (such as produced by Zotero).

`bookdown` supports `.html` output for ease and speed and also renders `.pdf` files through *L^AT_EX* for publication-ready documents.

For University of Washington theses and dissertations, consider Ben Marwick's [`huskydown` package](#) which uses Markdown but renders via a UW approved *L^AT_EX* template.

Best Practices

Organization and Portability

Organization Systems

Organizing research projects is something you either do accidentally—and badly—or purposefully with some upfront labor.

Uniform organization makes switching between or revisiting projects easier.

I suggest something like the following:

```
project/  
  readme.md  
  data/  
    derived/  
      processed_data.RData  
    raw/  
      core_data.csv  
  docs/  
    paper.Rmd  
  syntax/  
    functions.R  
    models.R
```

1. There is a clear hierarchy
 - Written content is in docs
 - Code is in syntax
 - Data is in data
2. Naming is uniform
 - All lower case
 - Words separated by underscores
3. Names are self-descriptive

Workflow versus Project

To summarize Jenny Bryan, one should separate workflow from projects.

Workflow

- The software you use to write your code (e.g. RStudio)
- The location you store a project
- The specific computer you use
- The code you ran earlier or typed into your console

Project

- The raw data
- The code that operates on your raw data
- The packages you use
- The output files or documents

Projects *should not modify anything outside of the project* nor need to be modified by someone else (or future you) to run.

Projects *should be independent of your workflow.*

Portability

For research to be reproducible, it must also be *portable*. Portable software operates *independently of workflow* such as fixed file locations.

Do Not:

- Use `setwd()` in scripts or .Rmd files.
- Use *absolute paths* except for *fixed, immovable sources* (secure data).
 - `read_csv("C:/my_project/data/my_data.csv")`
- Use `install.packages()` in script or .Rmd files.
- Use `rm(list=ls())` anywhere but your console.

Do:

- Use RStudio projects (or the [here_package](#)) to set directories.
- Use *relative paths* to load and save files:
 - `read_csv("../data/my_data.csv")`
- Load all required packages using `library()`.
- Clear your workspace when closing RStudio.
 - Set *Tools > Global Options... > Save workspace...* to **Never**

Divide and Conquer

Often you do not want to include all the code for a project in a single `.Rmd` file:

- The code takes too long to knit.
- The file is so long it is difficult to read.

There are two ways to deal with this:

1. Use separate `.R` scripts or `.Rmd` files which save results from complicated parts of a project, then load these results in the main `.Rmd` file.
 - This is good for loading and cleaning large data.
 - Also for running slow models.
2. Use `source()` to run external `.R` scripts.
 - This can be used to run large files that aren't impractically slow.
 - Also good for loading project-specific functions.

The Way of Many Files

I find it beneficial to break projects into *many* files:

- Scripts with specialized functions.
- Scripts to load and clean each set of variables.
- Scripts to run each set of models and make tables and plots.
- A main .Rmd that runs some or all of these to reproduce the entire project.

Splitting up a project carries benefits:

- Once a portion of the project is done and in its own file, *it is out of your way*.
- If you need to make changes, you don't need to search through huge files.
- Entire sections of the project can be added or removed quickly (e.g. converted to an appendix of an article)
- **It is the only way to build a proper *pipeline* for a project.**

Pipelines

Professional researchers and teams design projects as a **pipeline**.

A **pipeline** is a series of consecutive processing elements (scripts and functions in R).

Each stage of a pipeline...

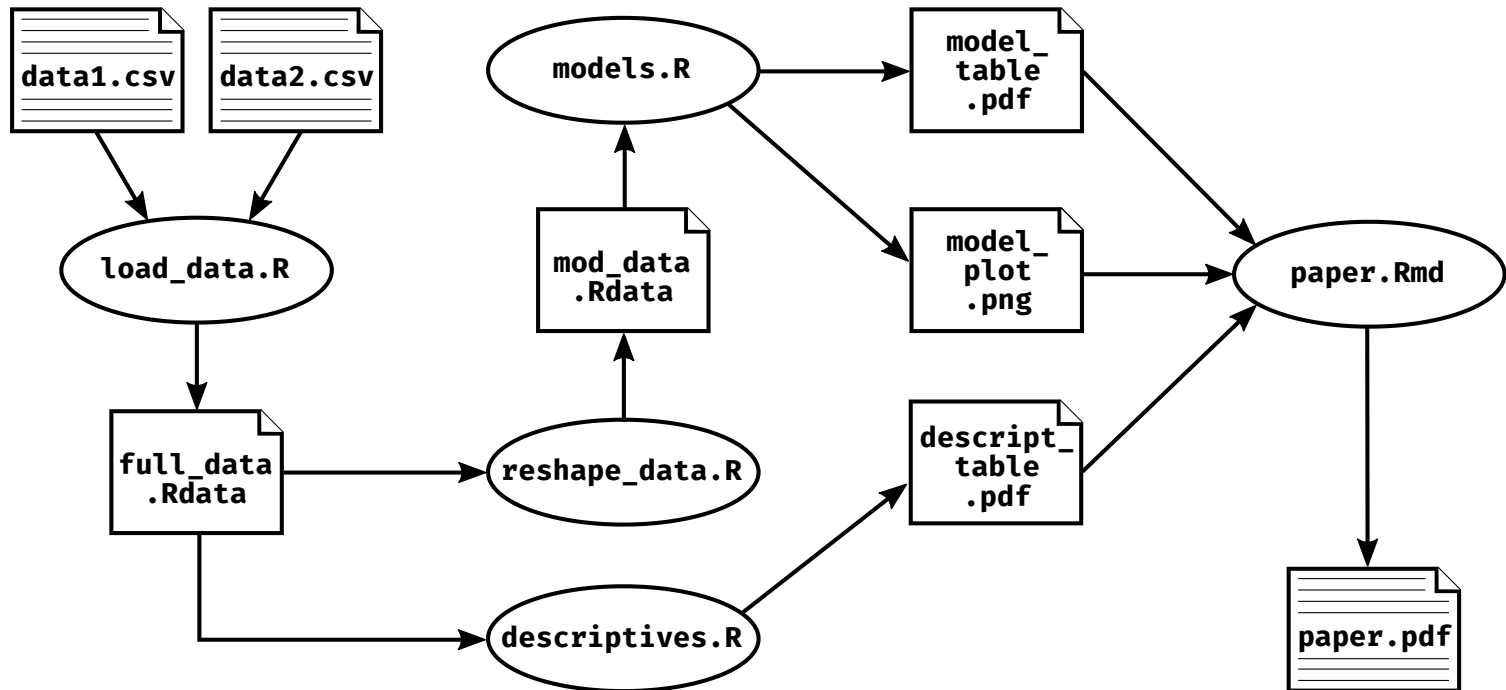
1. Has clearly defined inputs and outputs
2. Does not modify its inputs.
3. Produces the exact same output every time it is re-run.

This means...

1. When you modify one stage, you only need to rerun *subsequent stages*.
2. Different people can work on each stage.
3. Problems are isolated within stages.
4. You can depict your project as a *directed graph* of **dependencies**.

Example Pipeline

Every stage (oval) has an unambiguous input and output. Everything that precedes a given stage is a **dependency**—something required to run it.



Tools

Some opinionated advice

On Formats

Avoid "closed" or commercial software and file formats except where absolutely necessary.

Use open source software and file formats.

- It is always better for *science*:
 - People should be able to explore your research without buying commercial software.
 - You do not want your research to be inaccessible when software is updated.
- It is often just *better*.
 - It is usually updated more quickly
 - It tends to be more secure
 - It is rarely abandoned

The ideal: Use software that reads and writes *raw text*.

Text

Writing and formatting documents are two completely separate jobs.

- Write first
- Format later
- [Markdown](#) was made for this

Word processors—like Microsoft Word—try to do both at the same time, usually badly.

They waste time by leading you to format instead of writing.

Find a good modular text editor and learn to use it:

- [Atom](#)
- [Sublime](#) (Commercial)
- Emacs
- Vim

Version Control

Version Control

Version control originates in collaborative software development.

The Idea: All changes ever made to a piece of software are documented, saved automatically, and revertible.

Version control allows all decisions ever made in a research project to be documented automatically.

Version control can:

1. Protect your work from destructive changes
2. Simplify collaboration by merging changes
3. Document design decisions
4. Make your research process transparent

Git and GitHub

`git` is the dominant platform for version control, and [GitHub](#) is a free (and now Microsoft owned) platform for hosting **repositories**.

Repositories are folders on your computer where all changes are tracked by Git.

Once satisfied with changes, you "commit" them then "push" them to a remote repository that stores your project.

Others can copy your project ("pull"), and if you permit, make suggestions for changes.

Constantly committing and pulling changes automatically generates a running "history" that documents the evolution of a project.

`git` is integrated into RStudio under the *Tools* menu. [It requires some setup.](#)¹

[1] You can also use the [GitHub desktop application](#).

GitHub as a CV

Beyond archiving projects and allowing sharing, GitHub also serves as a sort of curriculum vitae for the programmer.

By allowing others to view your projects, you can display competence in programming and research.

If you are planning on working in the private sector, an active GitHub profile will give you a leg up on the competition.

If you are aiming for academia, a GitHub account signals technical competence and an interest in research transparency.

Wrapping up the Course

What You've Learned

A lot!

- How to get data into R from a variety of formats
- How to do "data custodian" work to manipulate and clean data
- How to make pretty visualizations
- How to automate with loops and functions
- How to combine text, calculations, plots, and tables into dynamic R Markdown reports
- How to acquire and work with spatial data

What Comes Next?

- Statistical inference (e.g. more CSSS courses)
 - Functions for hypothesis testing, hierarchical/mixed effect models, machine learning, survey design, etc. are straightforward to use... once data are clean
 - Access output by working with list structures (like from regression models) or using `broom` and `ggeffects`
- Practice, practice, practice!
 - Replicate analyses you've done in Excel, SPSS, or Stata
 - Think about data using `dplyr` verbs, tidy data principles
 - R Markdown for reproducibility
- More advanced projects
 - Using version control (git) in RStudio
 - Interactive Shiny web apps
 - Write your own functions and put them in a package

Course Plugs

If you...

- have no stats background yet - **SOC504: Applied Social Statistics**
- want to learn more social science computing - **SOC590: Big Data and Population Processes** ¹
- have (only) finished SOC506 - **CSSS510: Maximum Likelihood**
- want to master visualization - **CSSS569: Visualizing Data**
- study events or durations - **CSSS544: Event History Analysis** ²
- want to use network data - **CSSS567: Social Network Analysis**
- want to work with spatial data - **CSSS554: Spatial Statistics**
- want to work with time series - **CSSS512: Time Series and Panel Data**

[1] We're hoping to offer that again soon!

[2] Also a great maximum likelihood introduction.

Thank you!