

# CSSS508, Week 11

## Working with Model Results

Chuck Lanfear

Jun 5, 2019

Updated: Jun 7, 2019



# Topics for Today

## Displaying Model Results

- `broom`
  - Turning model output lists into dataframes
  - Summarizing models
- `ggeffects`
  - Creating counterfactual estimates
  - Plotting marginal effects
- Manual counterfactual plots
- Making regression tables
  - Using `pander` for models
  - Using `sjTable()` in `sjPlot`
- Wrapping up the course

broom

# 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
## -9.3746 -2.0569 -0.0685  2.1940  8.3739
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.4196     0.4215   3.368 0.000912 ***
## num1          0.4192     0.1084   3.867 0.000150 ***
## fac1B         0.6294     0.5423   1.161 0.247195
## fac1C         1.7502     0.5472   3.198 0.001612 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.067 on 196 degrees of freedom
## Multiple R-squared:  0.1128,    Adjusted R-squared:  0.09922
## F-statistic: 8.307 on 3 and 196 DF,  p-value: 3.15e-05
....
```

# 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
## -1.6209  -1.0734  -0.6453   1.0818   2.1621
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.80227     0.30262  -2.651 0.008023 **
## num1         0.28284     0.07915   3.573 0.000353 ***
## fac1B        0.37719     0.37295   1.011 0.311844
## fac1C        0.84695     0.38039   2.227 0.025979 *
## ---
## 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)    1.42      0.422     3.37 0.000912
## 2 num1           0.419     0.108     3.87 0.000150
## 3 fac1B          0.629     0.542     1.16 0.247
## 4 fac1C          1.75     0.547     3.20 0.00161
```

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)  -0.802     0.303     -2.65  0.00802
## 2 num1         0.283     0.0792     3.57  0.000353
## 3 fac1B        0.377     0.373      1.01  0.312
## 4 fac1C        0.847     0.380      2.23  0.0260
```



# 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
##   <dbl>      <dbl> <dbl>      <dbl>  <dbl> <int>  <dbl> <dbl> <dbl>
## 1    0.113        0.0992  3.07        8.31 3.15e-5     4 -506. 1022. 1038.
## # ... with 2 more variables: deviance <dbl>, 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
##   <dbl>   <dbl> <fct>   <dbl>   <dbl> <dbl>   <dbl> <dbl>   <dbl>
## 1  0.139   0.834  A       1.77    0.407  -1.63  0.0176   3.07  1.29e-3
## 2  6.85    4.53   A       3.32    0.555   3.53  0.0328   3.06  1.16e-2
## 3 -0.735  -1.15   B       1.57    0.440  -2.30  0.0206   3.07  3.02e-3
## 4  7.44    0.316  C       3.30    0.372   4.14  0.0147   3.06  6.91e-3
## 5  3.34    0.0845 A       1.46    0.419   1.89  0.0187   3.07  1.83e-3
## 6  1.83    0.521  C       3.39    0.369  -1.56  0.0145   3.07  9.62e-4
## # ... with 1 more variable: .std.resid <dbl>
```

# 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)    1.15     0.451      2.55 1.38e- 2
## 2 A     num1           0.436     0.126      3.46 1.07e- 3
## 3 A     fac2No         0.289     0.585      0.493 6.24e- 1
## 4 A     num2           0.743     0.0878     8.46 2.09e-11
## 5 B     (Intercept)    1.03     0.368      2.79 6.73e- 3
## 6 B     num1           0.482     0.135      3.57 6.59e- 4
## 7 B     fac2No         1.76     0.508      3.47 8.99e- 4
## 8 B     num2           0.694     0.0773     8.99 3.13e-13
## 9 C     (Intercept)    2.79     0.324      8.61 2.19e-12
## 10 C    num1           0.670     0.124      5.41 9.42e- 7
## 11 C    fac2No         0.761     0.477      1.59 1.16e- 1
## 12 C    num2           0.689     0.0880     7.83 5.26e-11
```

`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.

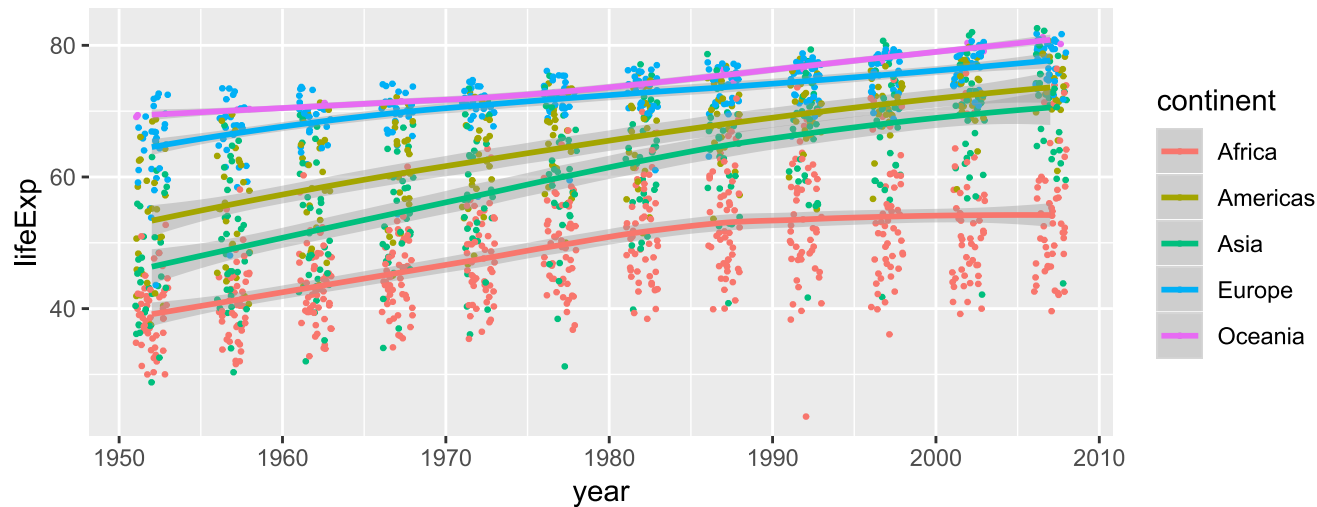
We can use `geom_smooth()` to add a layer depicting common bivariate models.

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

```
library(gapminder)
```

# Default `geom_smooth()`

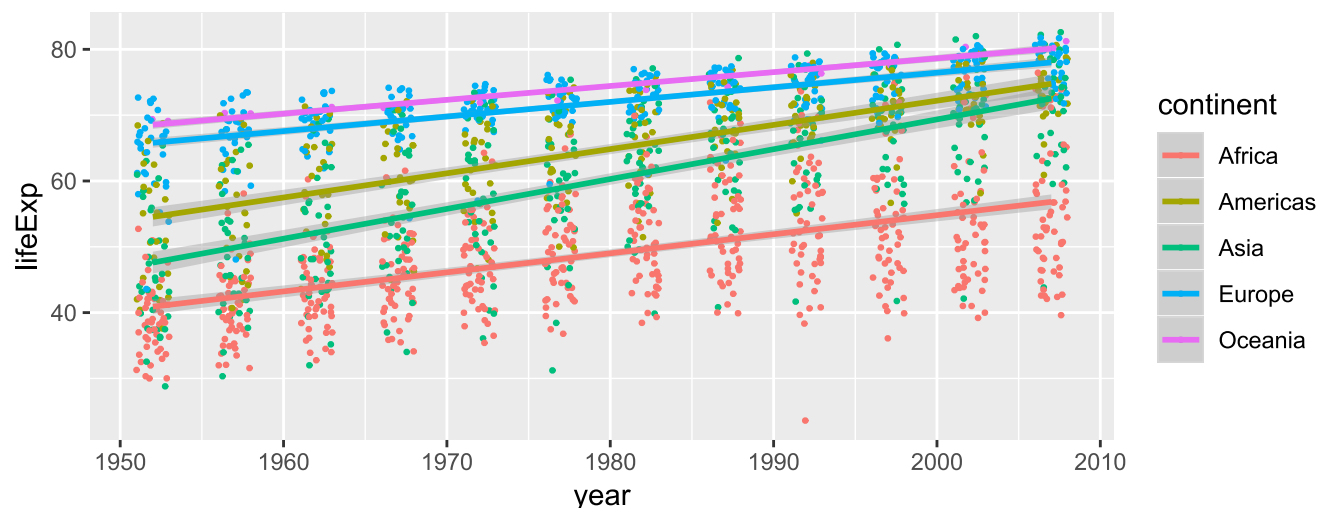
```
ggplot(data = gapminder,  
       aes(x = year, y = lifeExp, color = continent)) +  
  geom_point(position = position_jitter(1,0), size = 0.5) +  
  geom_smooth()
```



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

# Linear glm

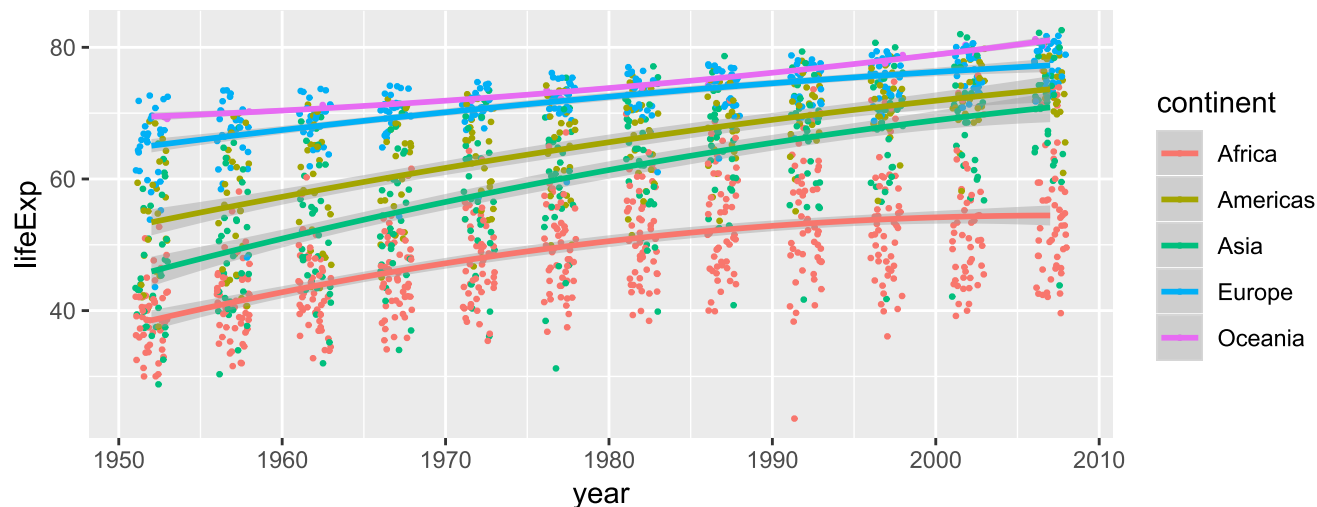
```
ggplot(data = gapminder,  
       aes(x = year, y = lifeExp, color = continent)) +  
  geom_point(position = position_jitter(1,0), size = 0.5) +  
  geom_smooth(method = "glm", formula = y ~ x)
```



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

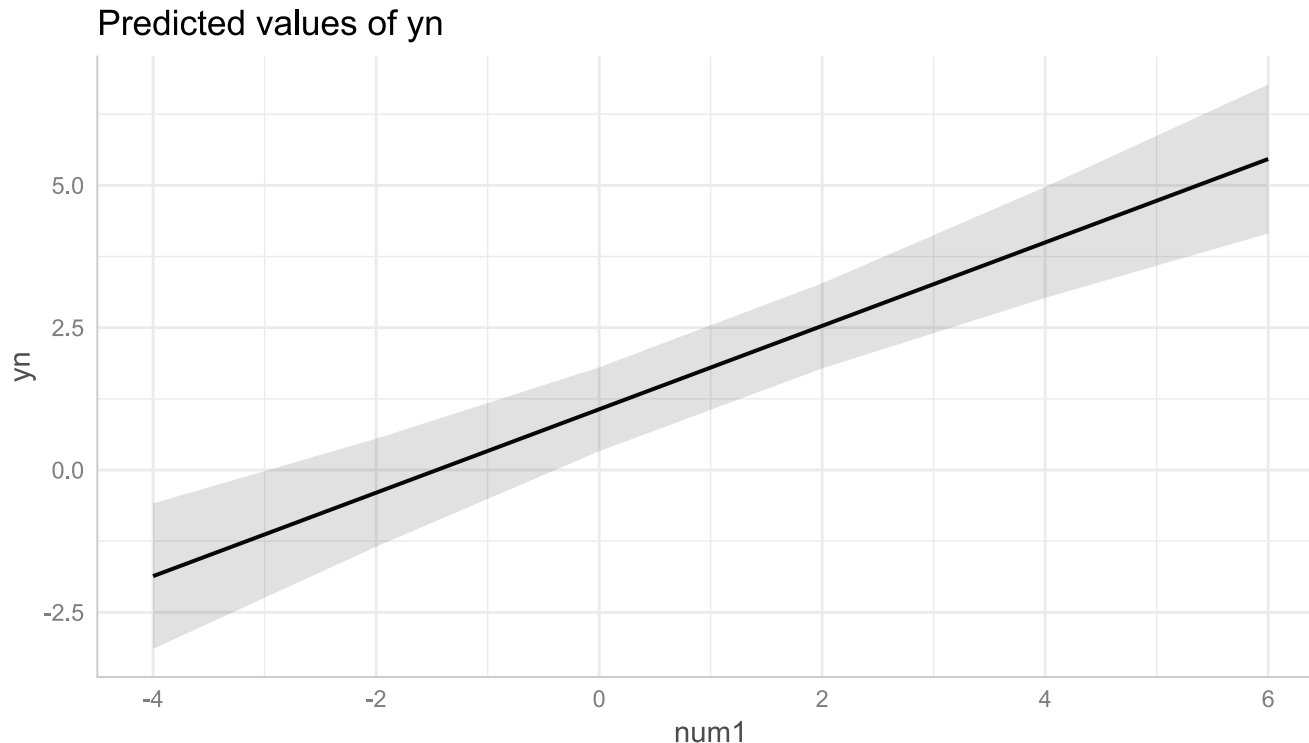
```
lm_1_est
```

```
##  
## # Predicted values of yn  
## # x = num1  
##  
##   x predicted std.error conf.low conf.high  
## -4    -1.864    0.651   -3.141   -0.587  
## -2    -0.399    0.484   -1.347    0.550  
##  0     1.066    0.376    0.329    1.803  
##  2     2.531    0.382    1.784    3.279  
##  4     3.997    0.497    3.023    4.971  
##  6     5.462    0.667    4.154    6.770  
##  
## 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)
```

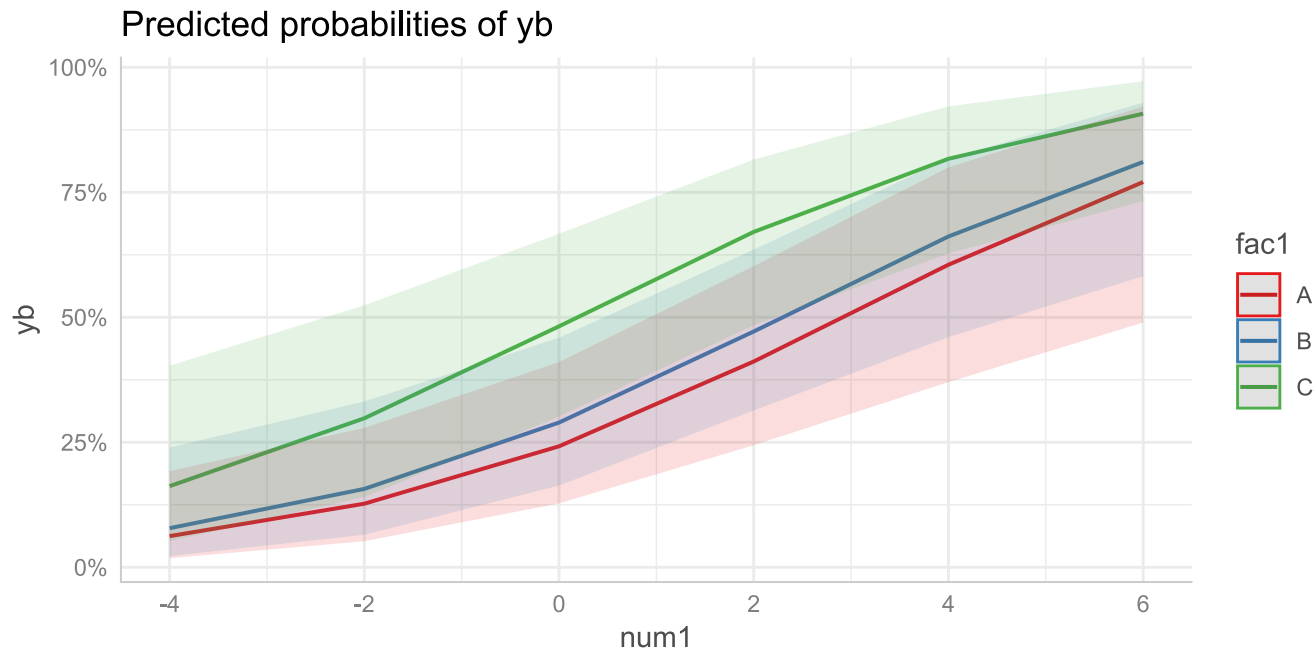




# Grouping with `ggpredict()`

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()`!

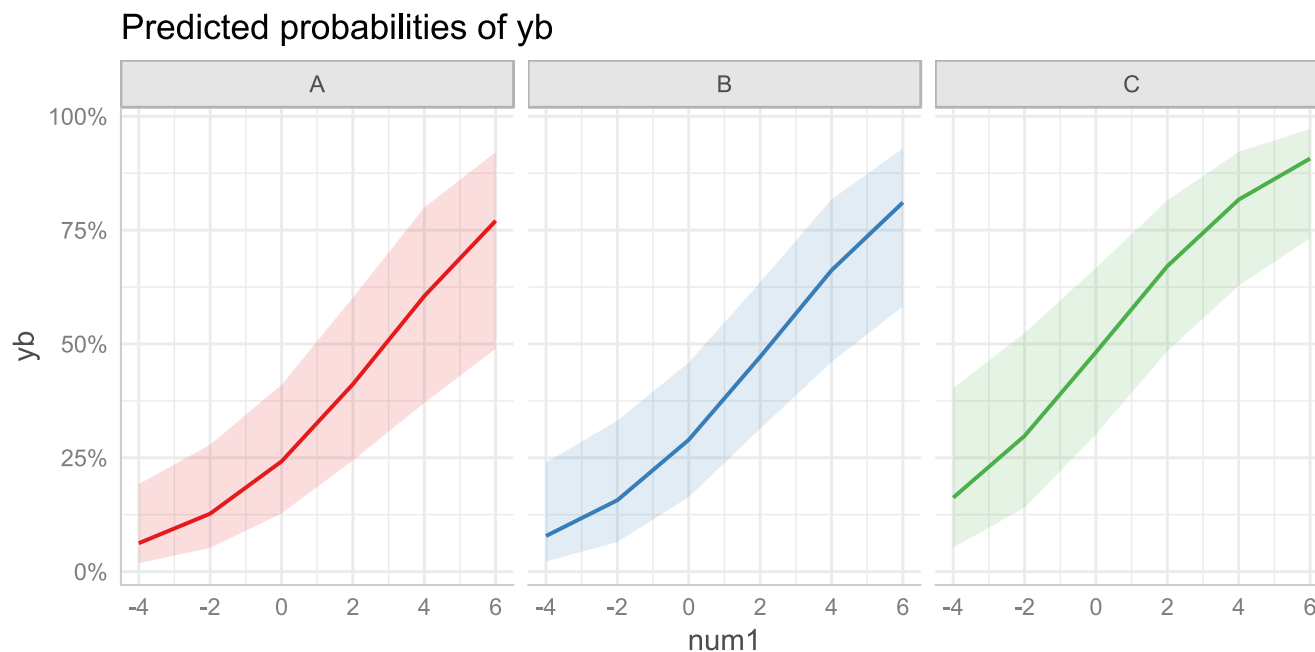
```
glm(yb ~ num1 + fac1 + num2 + fac2, data = ex_dat, family=binomial(link = "logit")) %>%  
  ggpredict(terms = c("num1", "fac1")) %>% plot()
```



# 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)
```

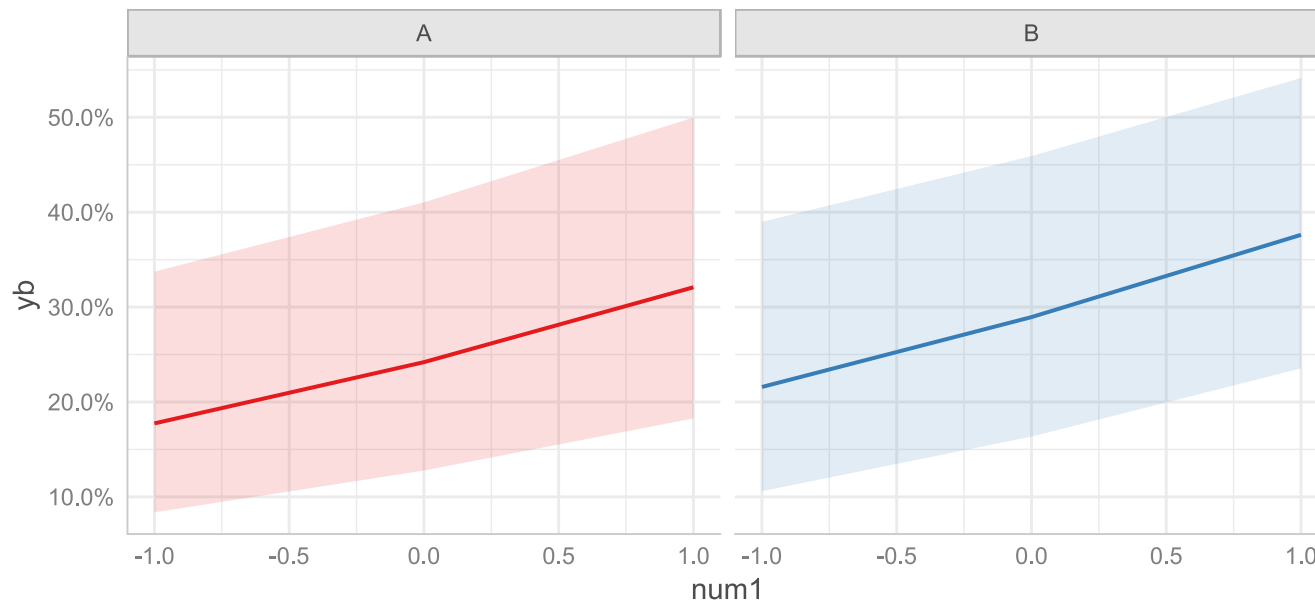


# Counterfactual Values

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)
```

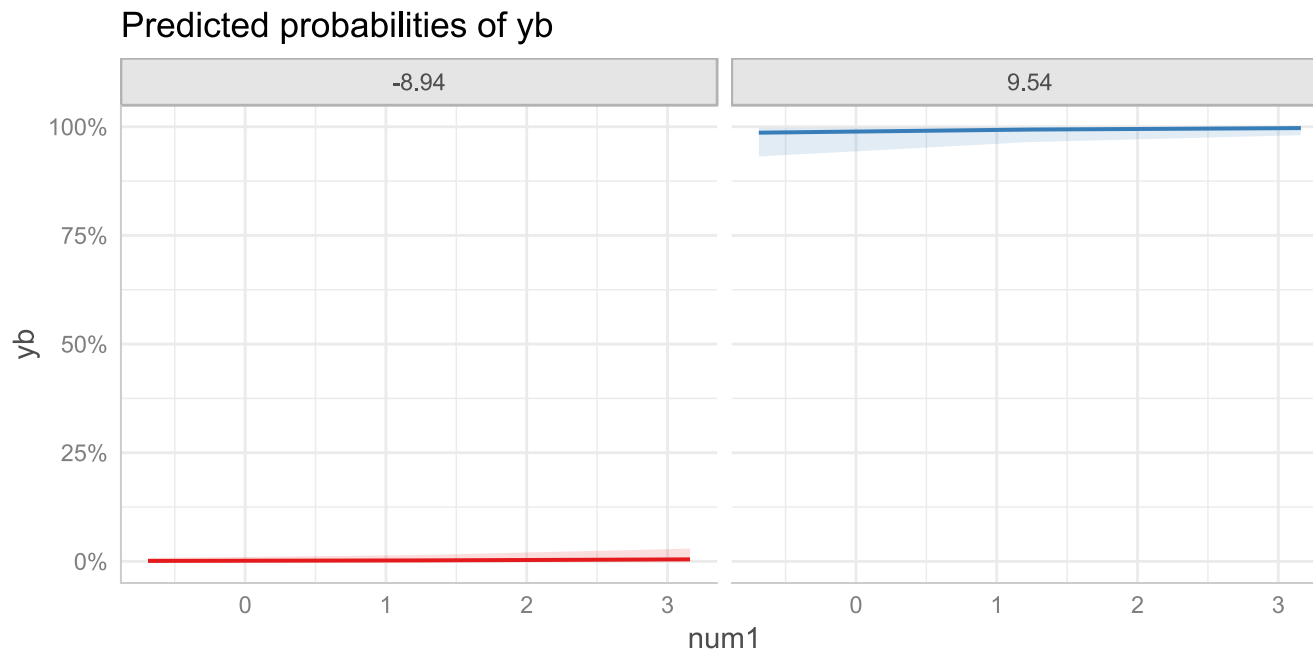
Predicted probabilities of yb



# Representative Values

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

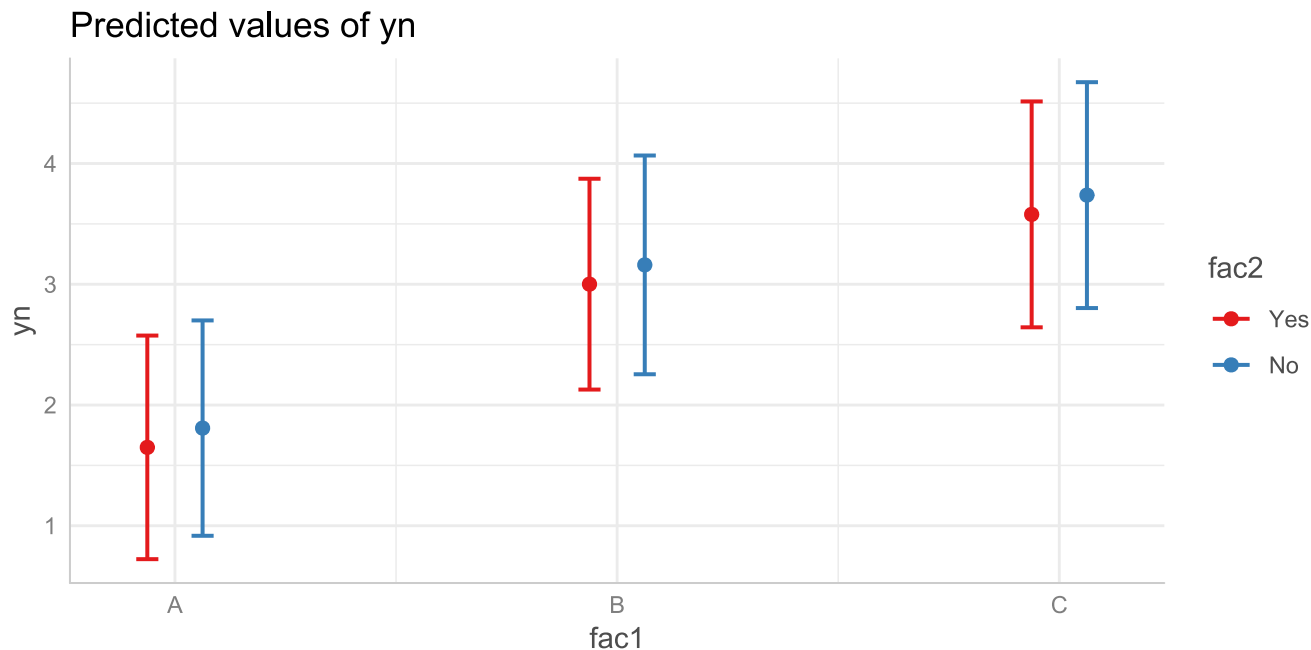
```
glm(yb ~ num1 + fac1 + num2 + fac2, data = ex_dat, family = binomial(link = "logit")) %>%  
  ggpredict(terms = c("num1 [meansd]", "num2 [minmax]")) %>% plot(facet=TRUE)
```



# 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

# An Advanced Example

Here is an example using a model from a [recent article I worked on](#).

This models the likelihood of arrest of a target in a police contact conditional on neighborhood, race of target, and race of who called the police.

```
load("data/any_arrest_data.RData")
mod_arrest <- glm(arrest ~ white_comp_wit_vict*black_arr_susp +
                  crime_type*white_comp_wit_vict + caller_type +
                  arr_susp_subj_count + comp_wit_vict_count +
                  black_arr_susp*neighb_type + crime_type*neighb_type +
                  serious_rate + pbl + pot + dis + year,
                  family = binomial(link = "logit"),
                  data = any_arrest_data)
```

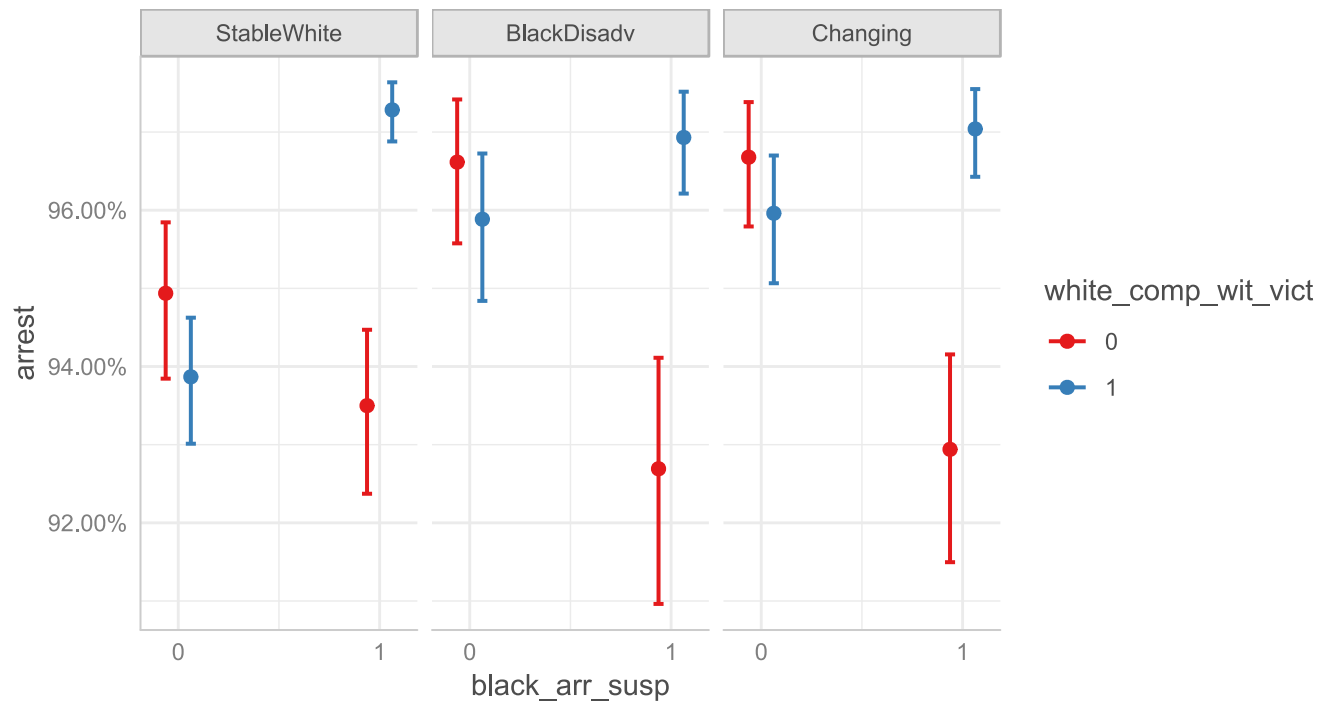
There are a lot of interactions here:

- Target Race x Caller Race
- Crime Type x Caller Race
- Target Race x Neighborhood Type
- Crime Type x Neighborhood Type

# ggeffects Output

```
mod_arrest %>% ggpredict(terms = c("black_arr_susp",  
  "white_comp_wit_vict", "neighb_type")) %>% plot()
```

Predicted probabilities of arrest





# A Complex Example

`ggpredict()` can only handle three variables in its `terms=` argument.

For my article, I wanted to plot estimates across counterfactual values of all four variables in my interaction terms:

- Caller Race
- Target Race
- Crime Type
- Neighborhood Type

How could I do this?

Stats + Math + Code = ♥

# Some Background

Given we've estimate a model, consider the following:

1.  $\hat{Y} = X\hat{\beta}$ , where  $X$  is the model matrix and  $\hat{\beta}$  is the coefficients.
2.  $\hat{\beta}$  is a vector of *random variables* whose estimated distribution is described by parameter variance-covariance matrix  $\Sigma$ .

Using this, we can do the following:

1. Extract the model matrix  $X$ , estimated coefficients ( $\hat{\beta}$ ), and  $\Sigma$  from our fitted model.
2. Make lots of random parameter draws centered on  $\hat{\beta}$  and distributed according to  $\Sigma$ .
3. Multiply *each* of these draws by *counterfactual*  $X$  values to get  $\hat{Y}$  values.
4. Take the 2.5% and 97.% quantiles of these  $\hat{Y}$  values.

This produces a *simulated* mean and confidence interval. This is called the **percentile method**, a type of *bootstrapping*.

# Simulating Coefficients

We can make random draws from our estimated distribution of parameters using `MASS::mvrnorm()` which takes three main arguments:

1. `n`: The number of draws
2. `mu`: mean—our coefficient estimates—obtained via `coef()`.
3. `Sigma`: a covariance matrix, obtained via `vcov()`.

```
sim_params <- MASS::mvrnorm(n = 10000,  
                             mu = coef(mod_arrest),  
                             Sigma = vcov(mod_arrest))  
sim_params[1:6, 1:4]
```

```
##      (Intercept) white_comp_wit_vict1 black_arr_susp1 crime_typeNuisance  
## [1,]      2.649874          -0.24837403          -0.4269412          -0.4063155  
## [2,]      2.788417          -0.15061167          -0.1473641          -0.4704289  
## [3,]      2.704444          -0.34618448          -0.3178364          -0.3906221  
## [4,]      2.609430          -0.09684517          -0.2334463          -0.4445482  
## [5,]      2.793919          -0.21388214          -0.2317792          -0.5622488  
## [6,]      2.733762          -0.26654232          -0.2280260          -0.5613452
```

# Counterfactual Values

Next we need a data frame with our counterfactual values.

We want one row (or *scenario*) per estimate to plot, and all variables at their means *except* the ones we are varying. We also don't want impossible values; `neighb_type` values are mutually exclusive.

```
x_values <- colMeans(model.matrix(mod_arrest)) # vars at mean
n_scen  <- (2*2*2*3) # Number of scenarios
x_frame <- setNames(data.frame(matrix(x_values, nrow=n_scen,
                                     ncol=length(x_values),
                                     byrow=T)), names(x_values))
cf_vals <- arrangements::permutations(c(0,1), k=5, replace=T)
cf_vals <- cf_vals[cf_vals[,4]+cf_vals[,5]!=2,] # Remove impossible vals
colnames(cf_vals) <- c("white_comp_wit_vict1", "black_arr_susp1",
                     "crime_typeNuisance", "neighb_typeBlackDisadv",
                     "neighb_typeChanging")
x_frame[colnames(cf_vals)] <- cf_vals # assign to countefactual df
```

`permutations()` is a quick way to get all combinations of some values.

# What Do We Have?

```
glimpse(x_frame)
```

```
## Observations: 24
## Variables: 24
## $ `(Intercept)`           <dbl> 1, 1, 1, 1, 1, 1, ...
## $ white_comp_wit_vict1     <dbl> 0, 0, 0, 0, 0, 0, ...
## $ black_arr_susp1          <dbl> 0, 0, 0, 0, 0, 0, ...
## $ crime_typeNuisance       <dbl> 0, 0, 0, 1, 1, 1, ...
## $ caller_typeVictim        <dbl> 0.8019442, 0.80194...
## $ caller_typeWitness       <dbl> 0.08516245, 0.0851...
## $ arr_susp_subj_count      <dbl> 1.552955, 1.552955...
## $ comp_wit_vict_count      <dbl> 1.571604, 1.571604...
## $ neighb_typeBlackDisadv   <dbl> 0, 0, 1, 0, 0, 1, ...
## $ neighb_typeChanging     <dbl> 0, 1, 0, 0, 1, 0, ...
## $ serious_rate             <dbl> 2.157661e-17, 2.15...
## $ pbl                      <dbl> 1.185962e-16, 1.18...
## $ pot                      <dbl> -2.808814e-17, -2....
## $ dis                      <dbl> 9.987197e-18, 9.98...
## $ year2009                 <dbl> 0.2829368, 0.28293...
## $ year2010                 <dbl> 0.1670504, 0.16705...
## $ year2011                 <dbl> 0.09201842, 0.0920...
## $ year2012                 <dbl> 0.1232284, 0.12322...
## $ `white_comp_wit_vict1:black_arr_susp1` <dbl> 0.4452034, 0.44520...
## $ `white_comp_wit_vict1:crime_typeNuisance` <dbl> 0.1614479, 0.16144...
## $ `black_arr_susp1:neighb_typeBlackDisadv` <dbl> 0.04893835, 0.0489...
## $ `black_arr_susp1:neighb_typeChanging` <dbl> 0.111691, 0.111691...
## $ `crime_typeNuisance:neighb_typeBlackDisadv` <dbl> 0.0154771, 0.01547...
## $ `crime_typeNuisance:neighb_typeChanging` <dbl> 0.03300077, 0.0330...
```

# Fixing Interactions

Our main variables are correct... but we need to make our interaction terms.

The interaction terms in the model matrix have specific form `var1:var2`.

Their counterfactual values are just equal to the products of their components.

```
x_frame <- x_frame %>%  
  mutate(  
    `white_comp_wit_vict1:black_arr_susp1` = white_comp_wit_vict1*black_arr_susp1,  
    `white_comp_wit_vict1:crime_typeNuisance` = white_comp_wit_vict1*crime_typeNuisance,  
    `black_arr_susp1:neighb_typeBlackDisadv` = black_arr_susp1*neighb_typeBlackDisadv,  
    `black_arr_susp1:neighb_typeChanging` = black_arr_susp1*neighb_typeChanging,  
    `crime_typeNuisance:neighb_typeBlackDisadv` = crime_typeNuisance*neighb_typeBlackDisadv,  
    `crime_typeNuisance:neighb_typeChanging` = crime_typeNuisance*neighb_typeChanging,  
    `black_arr_susp1:neighb_typeBlackDisadv` = black_arr_susp1*neighb_typeBlackDisadv,  
    `black_arr_susp1:neighb_typeChanging` = black_arr_susp1*neighb_typeChanging)
```

# Fixed

```
glimpse(x_frame)
```

```
## Observations: 24
## Variables: 24
## $ `(Intercept)`           <dbl> 1, 1, 1, 1, 1, 1, ...
## $ white_comp_wit_vict1     <dbl> 0, 0, 0, 0, 0, 0, ...
## $ black_arr_susp1          <dbl> 0, 0, 0, 0, 0, 0, ...
## $ crime_typeNuisance       <dbl> 0, 0, 0, 1, 1, 1, ...
## $ caller_typeVictim        <dbl> 0.8019442, 0.80194...
## $ caller_typeWitness       <dbl> 0.08516245, 0.0851...
## $ arr_susp_subj_count      <dbl> 1.552955, 1.552955...
## $ comp_wit_vict_count      <dbl> 1.571604, 1.571604...
## $ neighb_typeBlackDisadv   <dbl> 0, 0, 1, 0, 0, 1, ...
## $ neighb_typeChanging      <dbl> 0, 1, 0, 0, 1, 0, ...
## $ serious_rate             <dbl> 2.157661e-17, 2.15...
## $ pbl                       <dbl> 1.185962e-16, 1.18...
## $ pot                       <dbl> -2.808814e-17, -2....
## $ dis                       <dbl> 9.987197e-18, 9.98...
## $ year2009                  <dbl> 0.2829368, 0.28293...
## $ year2010                  <dbl> 0.1670504, 0.16705...
## $ year2011                  <dbl> 0.09201842, 0.0920...
## $ year2012                  <dbl> 0.1232284, 0.12322...
## $ `white_comp_wit_vict1:black_arr_susp1` <dbl> 0, 0, 0, 0, 0, 0, ...
## $ `white_comp_wit_vict1:crime_typeNuisance` <dbl> 0, 0, 0, 0, 0, 0, ...
## $ `black_arr_susp1:neighb_typeBlackDisadv` <dbl> 0, 0, 0, 0, 0, 0, ...
## $ `black_arr_susp1:neighb_typeChanging` <dbl> 0, 0, 0, 0, 0, 0, ...
## $ `crime_typeNuisance:neighb_typeBlackDisadv` <dbl> 0, 0, 0, 0, 0, 1, ...
## $ `crime_typeNuisance:neighb_typeChanging` <dbl> 0, 0, 0, 0, 1, 0, ...
```

# Estimates!

Then we just multiply our parameters by our counterfactual data:

```
sims_logodds <- sim_params %*% t(as.matrix(x_frame))  
sims_logodds[1:6, 1:6]
```

```
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]  
## [1,] 2.177016 2.582403 2.589639 1.770700 2.297997 2.157879  
## [2,] 2.059036 2.527301 2.685111 1.588607 2.183050 2.127477  
## [3,] 2.279709 2.798278 2.722473 1.889087 2.531622 2.265328  
## [4,] 1.992662 2.504421 2.434744 1.548114 2.112916 1.995474  
## [5,] 2.145578 2.531906 2.512737 1.583329 2.066743 1.716077  
## [6,] 2.158352 2.622928 2.660889 1.597007 2.236660 2.181941
```

```
dim(sims_logodds)
```

```
## [1] 10000      24
```

Now we log-odds 10,000 estimates each (rows) of 24 counterfactual scenarios (columns).



# Getting Probabilities

The model for this example is a *logistic regression*, which produces estimates in *log-odds* ( $\ln(\text{Odds}(x))$ ).

We can convert these to probabilities based on two identities:

$$1. \text{Odds}(x) = e^{\ln(\text{Odds}(x))}$$

$$2. \text{Pr}(x) = \frac{\text{Odds}(x)}{(1 + \text{Odds}(x))}$$

```
sims_prob <- exp(sims_logodds) / (1 + exp(sims_logodds))  
sims_prob[1:6, 1:6]
```

```
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]  
## [1,] 0.8981664 0.9297205 0.9301918 0.8545447 0.9087110 0.8964028  
## [2,] 0.8868574 0.9260337 0.9361423 0.8304200 0.8987170 0.8935453  
## [3,] 0.9071826 0.9425827 0.9383398 0.8686514 0.9263291 0.9059645  
## [4,] 0.8800245 0.9244512 0.9194386 0.8246412 0.8921523 0.8803210  
## [5,] 0.8952548 0.9263485 0.9250299 0.8296755 0.8876285 0.8476229  
## [6,] 0.8964467 0.9323227 0.9346789 0.8315996 0.9034937 0.8986160
```

# A Quick Function

We are going to want to grab the mean and 95% confidence interval from our simulation estimates.

Here's a quick function to do it and make it pretty.

```
extract_pe_ci <- function(x){  
  vals <- c(mean(x), quantile(x, probs=c(.025, .975)))  
  names(vals) <- c("PE", "LB", "UB")  
  return(vals)  
}
```

This returns a length 3 vector with the following names:

- **PE** for *point estimate*
- **LB** for *lower bound* of the confidence interval
- **UB** for *upper bound*

# Prep for Plotting

First we extract our point estimates and confidence intervals by *applying* `extract_pe_ci()` to each column of estimated probabilities.

```
estimated_pes <- as.data.frame( t(apply(sims_prob, 2, extract_pe_ci)))
```

Then I add columns describing the scenarios to color, group, and facet over based on the counterfactual values.

```
estimated_pes$`Reporter`      <- ifelse(cf_vals[,1]==1, "Any White", "All Black")
estimated_pes$`Target`        <- ifelse(cf_vals[,2]==1, "Any Black", "All White")
estimated_pes$`Crime Type`    <- ifelse(cf_vals[,3]==1, "Nuisance Crime", "Serious Crime")
estimated_pes$`Neighborhood` <- case_when(
  cf_vals[,4]==1 ~ "Disadvantaged",
  cf_vals[,5]==1 ~ "Changing",
  TRUE ~ "Stable White")
```

# Final Tidy Data

```
estimated_pes %>% mutate_if(is.numeric, round, digits=3) # round for display
```

##	PE	LB	UB	Reporter	Target	Crime Type	Neighborhood
## 1	0.894	0.877	0.909	All Black	All White	Serious Crime	Stable White
## 2	0.929	0.914	0.942	All Black	All White	Serious Crime	Changing
## 3	0.927	0.909	0.943	All Black	All White	Serious Crime	Disadvantaged
## 4	0.831	0.799	0.859	All Black	All White	Nuisance Crime	Stable White
## 5	0.897	0.872	0.919	All Black	All White	Nuisance Crime	Changing
## 6	0.879	0.845	0.909	All Black	All White	Nuisance Crime	Disadvantaged
## 7	0.866	0.852	0.879	All Black	Any Black	Serious Crime	Stable White
## 8	0.855	0.834	0.874	All Black	Any Black	Serious Crime	Changing
## 9	0.850	0.825	0.874	All Black	Any Black	Serious Crime	Disadvantaged
## 10	0.790	0.758	0.820	All Black	Any Black	Nuisance Crime	Stable White
## 11	0.799	0.760	0.833	All Black	Any Black	Nuisance Crime	Changing
## 12	0.765	0.712	0.813	All Black	Any Black	Nuisance Crime	Disadvantaged
## 13	0.873	0.866	0.880	Any White	All White	Serious Crime	Stable White
## 14	0.914	0.900	0.927	Any White	All White	Serious Crime	Changing
## 15	0.912	0.896	0.928	Any White	All White	Serious Crime	Disadvantaged
## 16	0.749	0.732	0.766	Any White	All White	Nuisance Crime	Stable White
## 17	0.842	0.811	0.869	Any White	All White	Nuisance Crime	Changing
## 18	0.816	0.774	0.853	Any White	All White	Nuisance Crime	Disadvantaged
## 19	0.941	0.937	0.946	Any White	Any Black	Serious Crime	Stable White
## 20	0.936	0.927	0.945	Any White	Any Black	Serious Crime	Changing
## 21	0.934	0.922	0.944	Any White	Any Black	Serious Crime	Disadvantaged
## 22	0.875	0.863	0.885	Any White	Any Black	Nuisance Crime	Stable White
## 23	0.880	0.856	0.901	Any White	Any Black	Nuisance Crime	Changing
## 24	0.857	0.825	0.886	Any White	Any Black	Nuisance Crime	Disadvantaged

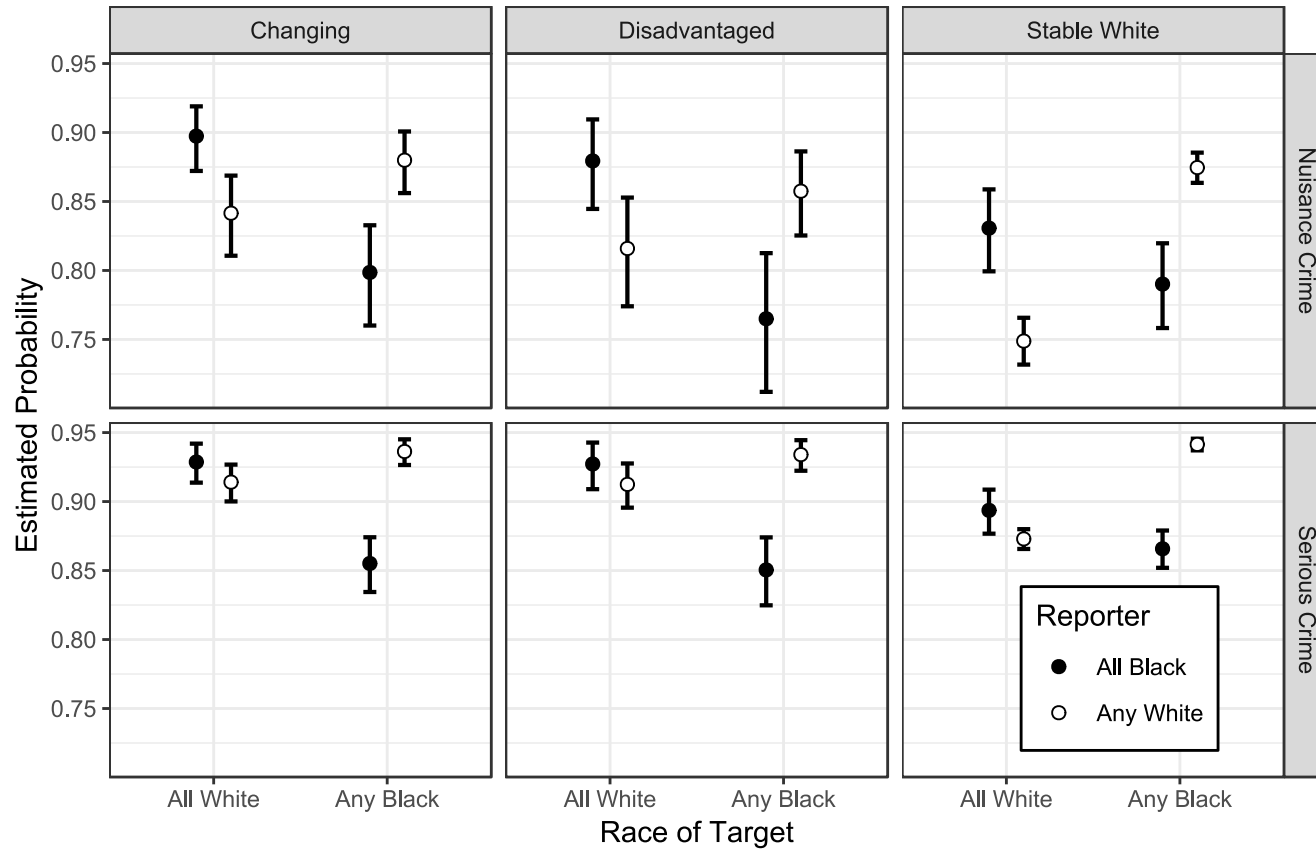
# Plot Code

Finally we plot estimates (PE) as points with error bars (UB, LB) stratified on Target and Reporter and faceted by Crime Type and Neighborhood.

```
ggplot(estimated_pes, aes(x = Target, y = PE, group = Reporter)) +  
  facet_grid(`Crime Type` ~ Neighborhood) +  
  geom_errorbar(aes(ymin = LB, ymax = UB),  
               position = position_dodge(width = .4),  
               size = 0.75, width = 0.15) +  
  geom_point(shape = 21, aes(fill = Reporter),  
            position = position_dodge(width = .4),  
            size = 2) +  
  scale_fill_manual("Reporter", values=c("Any White" = "white",  
                                         "All Black" = "black")) +  
  ggtitle("Figure 3. Probability of Arrest",  
         subtitle = "by Reporter and Target Race, Neighborhood and Crime Type") +  
  xlab("Race of Target") + ylab("Estimated Probability") +  
  theme_bw() + theme(legend.position = c(0.86, 0.15),  
                    legend.background = element_rect(color = 1))
```

# Plot

Figure 3. Probability of Arrest  
by Reporter and Target Race, Neighborhood and Crime Type



# 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)
<b>(Intercept)</b>	37.23	1.599	23.28	2.565e-20
<b>wt</b>	-3.878	0.6327	-6.129	1.12e-06
<b>hp</b>	-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)
<b>(Intercept)</b>	37.23	1.599	23.28	2.565e-20
<b>wt</b>	-3.878	0.6327	-6.129	1.12e-06
<b>hp</b>	-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)
```

# 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 <sup>2</sup> / adj. R <sup>2</sup>	.753 / .745			.827 / .815			.840 / .823		

# sjPlot does a lot more

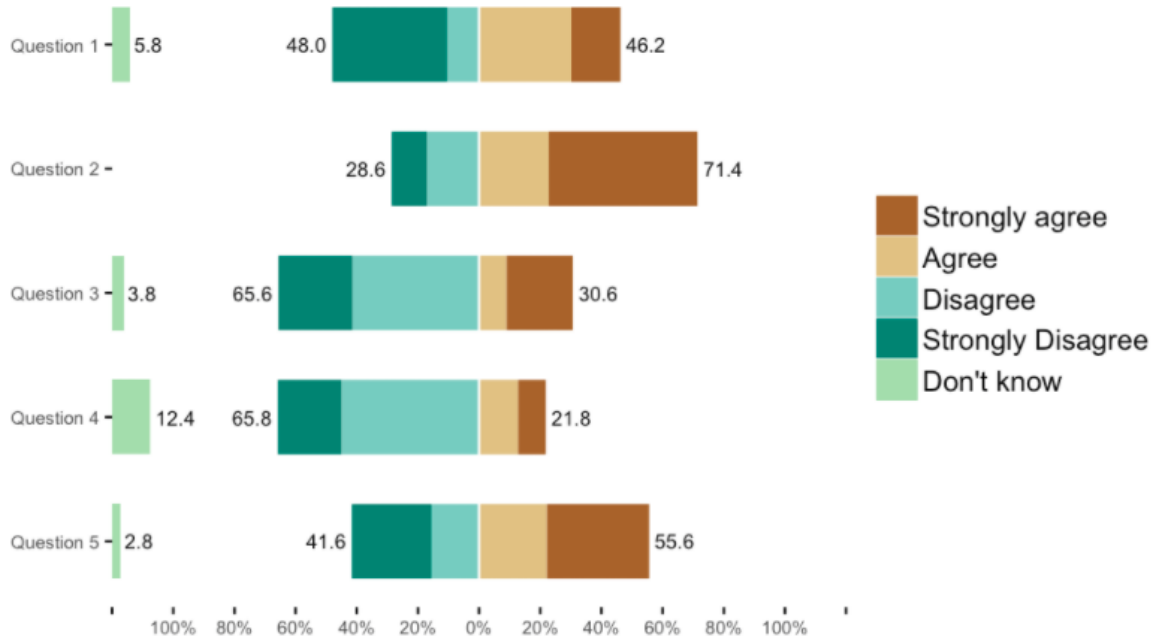
The `sjPlot` package does *a lot* more than just make pretty tables. It is a rabbit hole of *incredibly* powerful and useful functions for displaying descriptive and inferential results.

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<sup>A</sup>T<sub>E</sub>X*—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<sup>A</sup>T<sub>E</sub>X* using additive syntax similar to `ggplot2` and `dplyr`.

`stargazer` produces nicely formatted *L<sup>A</sup>T<sub>E</sub>X* tables but is idiosyncratic.

If you want to edit *L<sup>A</sup>T<sub>E</sub>X* documents, you can do it in R using Sweave documents (`.Rnw`). Alternatively, you may want to work in a dedicated *L<sup>A</sup>T<sub>E</sub>X* editor. I recommend [Overleaf](#) for this purpose.

RMarkdown has support for a fair amount of basic *L<sup>A</sup>T<sub>E</sub>X* syntax if you aren't trying to get too fancy!

Another approach I have used is to manually format *L<sup>A</sup>T<sub>E</sub>X* 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.

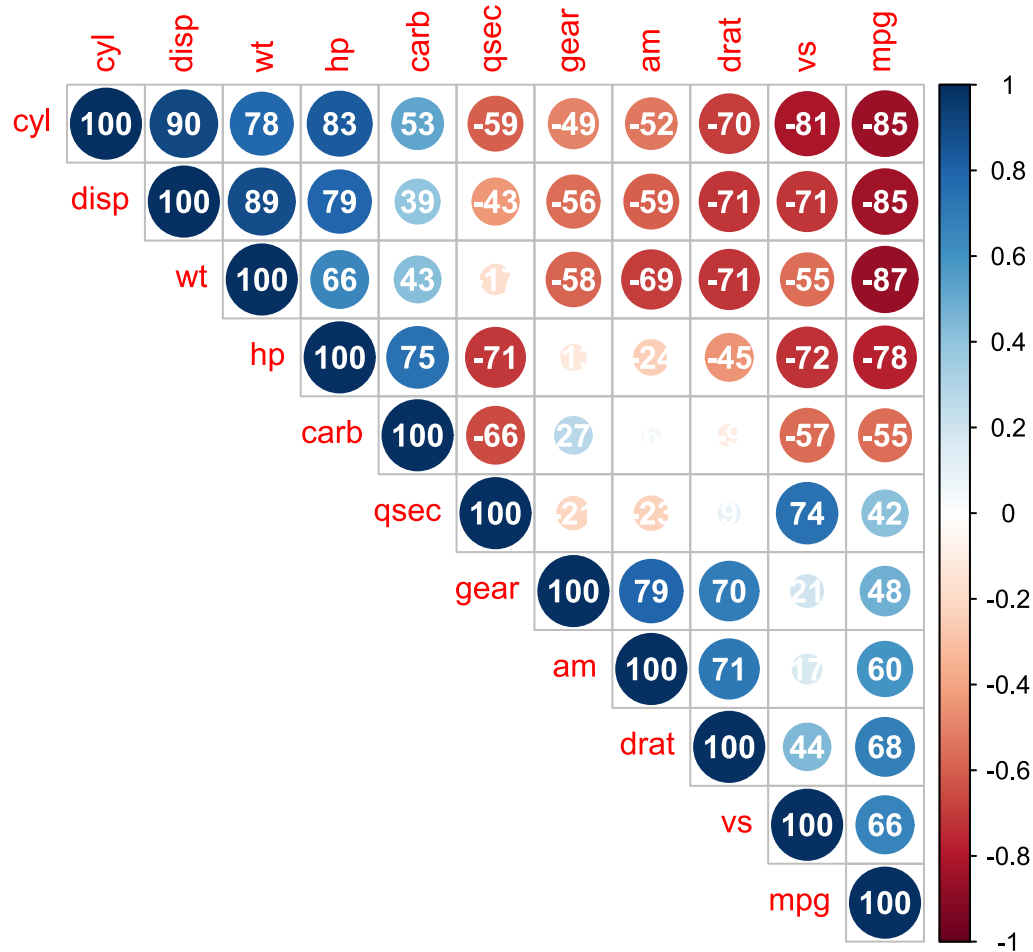
These make visualizing the correlations between variables in a data set easier.

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



# 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** <sup>1</sup>
- have (only) finished SOC506 - **CSSS510: Maximum Likelihood**
- want to master visualization - **CSSS569: Visualizing Data**
- study events or durations - **CSSS544: Event History Analysis** <sup>2</sup>
- 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!