



MULTIPLE REGRESSION WITH INTERACTION TERMS

Luis Castro-de-Araujo^a

11/07/2022

Virginia Institute for Psychiatric and Behavioral Genetics

^aPost-doc T32. luis.araujo@vcuhealth.org

Multiple regression recap

Interaction terms

Visualizing interactions

Marginal effects

MULTIPLE REGRESSION RECAP



- Extension of the simple linear regression model to two or more independent variables

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

- Expression = Baseline + Age + Tissue + Sex + Error
- Partial Regression Coefficients: effect on the dependent variable when increasing the i^{th} independent variable by 1 unit, **holding all other predictors constant**



- Qualitative variables are easily incorporated in regression framework through ***dummy variables***
- Simple example: sex can be coded as 0/1
- What if my categorical variable contains three levels:

$$x_1 = \begin{cases} 0 & \text{if AA} \\ 1 & \text{if AG} \\ 2 & \text{if GG} \end{cases}$$



- Previous coding would result in **colinearity**
- Solution is to set up a series of dummy variable.
- for k levels you need k-1 dummy variables

	x1	x2
AA	1	0
AG	0	1
GG	0	0

$$x_1 = \begin{cases} 1 & \text{if AA} \\ 0 & \text{otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1 & \text{if AG} \\ 0 & \text{otherwise} \end{cases}$$



Validity Does the data we're modeling matches the problem we're actually trying to solve?

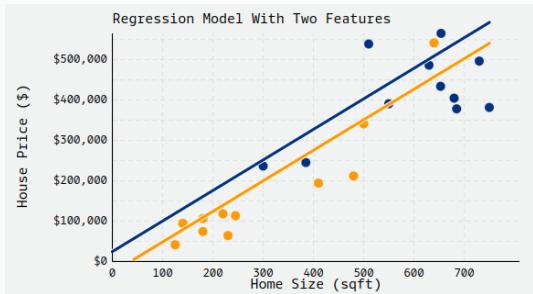
Representativeness Is the sample data used in the regression model representative of the population to which it will be applied?

Additivity and Linearity The deterministic component of a regression model is a linear function of the separate predictors: $y = B_0 + B_1x_1 + \dots + B_px_p$

Independence of Errors The errors from our model are independent.

Homoscedasticity The errors from our model have equal variance.

Normality of Errors The errors from our model are normally distributed.

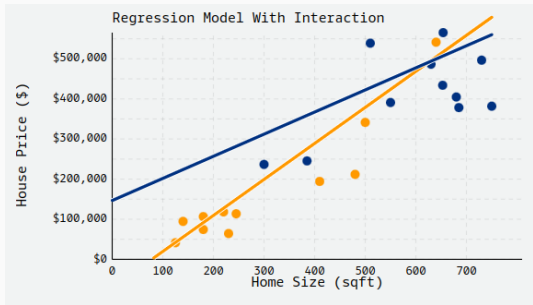


$$\text{houseprice} = -27154 + 757 * \text{sqft} + 51867 * \text{pool}$$

- In our example, we model home prices as a function of both the size of the house (sqft) and whether or not it has a pool

- intercept: -\$27,154, the predicted average housing price for houses with all $x_i = 0$. Or the cost of houses with no pools and a square-footage of zero.
- coefficient of pool: \$51,867, average expected price difference in houses of the same size (in sqft) if they do or do not have a pool. In other words, we expect, on average, houses of the same size to cost \$51,867 more if they have a pool than if they do not.
- coefficient of sqft: \$757, average expected price difference in housing price for houses that have the same value of pool but differ in size by one square-foot.
- We assume the same slope for sqft. Hence, two lines. This isn't always a valid assumption to make.

BACK TO OUR HOUSING EXAMPLE, NOW WITH INTERACTIONS



- interaction term: $-\$347$, represents the difference in the slope for sqft, comparing houses that do and do not have pools. Visually, this represents the difference between the slopes of the two lines.
- intercept: $-\$70,296$, represents the predicted housing price for houses with no pools and a square-footage of zero.
- coefficient of pool: $\$217,111$, represents the average expected difference in houses of the same size (0 sqft) that differed in whether or not they had a pool. (It's not super useful since we don't have houses with 0 square-feet).
- coefficient of sqft: $\$899$, represents the average expected difference in housing price for houses that do not have a pool (pool=0) but differ in size by one square-foot.

$$\text{houseprice} = -70296 + 899 * \text{sqft} + 217111 * \text{pool} - 347 * (\text{sqft} : \text{pool})$$

- If we believe that the slope for sqft should differ between houses that do have pools and houses that do not, we can add an interaction term to our model, (sqft:pool).

INTERACTION TERMS



- An interaction is a predictor that is some combination of the other predictors.



- Interactions are often the product of two or more predictors.
- Can be written as,

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$



- Conditional effects: the effect of a predictor on the response, holding all other predictors constant.
- Marginal effects: the effect of a predictor on the response, averaged over all values of the other predictors.



- If the conditional effects of X_1 on Y at different levels of X_2 are all the same then there is no interaction.



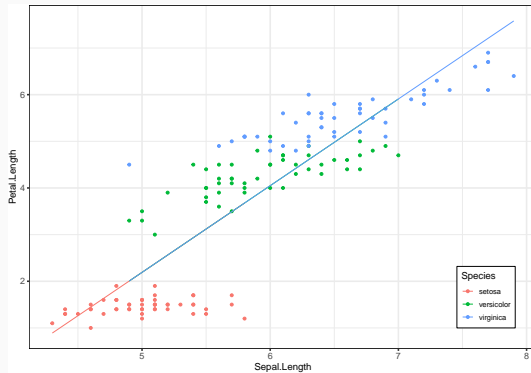
Parameter Meaning		Where people (used to) go awry
β_0	Expected value of the DV when X_1 and $X_2 == 0$	People get this
β_1	Effect of X_1 when $X_2 == 0$	Not marginal effects!
β_2	Effect of X_2 when $X_1 == 0$	Not marginal effects!
β_3	The addition to the conditional effect when both X_1 and X_2 are 1	People just look at the significance of the interaction parameter and do not calculate the underlying marginal or conditional effects or standard errors



- A common mistake that people make when interpreting interaction models is using the wrong standard errors.
- The standard errors that are printed in every regression table are the positive square roots of the diagonal elements of the variance- covariance matrix of β
- This does not matter anymore because of `margins()`



$$\widehat{petal.length}_i = \beta_0 + \beta_1 sepal.length_i$$





Creating the dummy

$$\text{setosa}_i = \begin{cases} 1 & \text{if species of flower } i = \text{setosa}, \forall i \in [1, 150] \\ 0 & \text{otherwise} \end{cases}$$

$$\text{versicolor}_i = \begin{cases} 1 & \text{if species of flower } i = \text{versicolor}, \forall i \in [1, 150] \\ 0 & \text{otherwise} \end{cases}$$

Our formula is then

$$\widehat{\text{petal.length}_i} = \beta_0 + \beta_1 \text{sepal.length}_i + \beta_2 \text{setosa}_i + \beta_3 \text{versicolor}_i$$



If it is setosa

$$\begin{aligned}\widehat{petal.length}_i &= \beta_0 + \beta_1 sepal.length_i + \beta_2 setosa_i + \beta_3 versicolor_i \\ &= \beta_0 + \beta_1 sepal.length_i + \beta_2 1 + \beta_3 0 \\ &= \beta_0 + \beta_1 sepal.length_i + \beta_2 \\ &= (\beta_0 + \beta_2) + \beta_1 sepal.length_i\end{aligned}$$

If it is versicolor

$$\begin{aligned}\widehat{petal.length}_i &= \beta_0 + \beta_1 sepal.length_i + \beta_2 setosa_i + \beta_3 versicolor_i \\ &= \beta_0 + \beta_1 sepal.length_i + \beta_2 0 + \beta_3 1 \\ &= \beta_0 + \beta_1 sepal.length_i + \beta_3 \\ &= (\beta_0 + \beta_3) + \beta_1 sepal.length_i\end{aligned}$$

If it is virginica

$$\begin{aligned}\widehat{petal.length}_i &= \beta_0 + \beta_1 sepal.length_i + \beta_2 setosa_i + \beta_3 versicolor_i \\ &= \beta_0 + \beta_1 sepal.length_i + \beta_2 0 + \beta_3 0 \\ &= \beta_0 + \beta_1 sepal.length_i \\ &= \beta_0 + \beta_1 sepal.length_i\end{aligned}$$

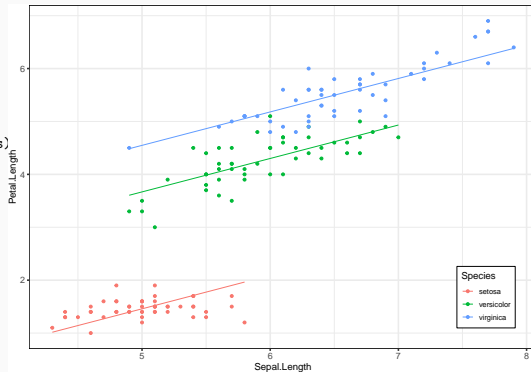
SAME SLOPE, DIFFERENT INTERCEPTS



```
iris$pred <- predict(lm(Petal.Length ~ Species+Sepal.Length,  
  data = iris))
```

plot in ggplot

```
ggplot(iris, aes(x = Sepal.Length, y = Petal.Length, color = Species))  
  geom_point() +  
  geom_line(aes(Sepal.Length, pred )) +  
  theme_luis() +  
  theme( legend.position = c(0.9, 0.15))
```





$$\widehat{petal.length}_i = \beta_0 + \beta_1 sepal.length_i + \beta_2 setosa_i + \beta_3 versicolor_i \\ + \beta_4 sepal.length_i setosa_i + \beta_5 sepal.length_i versicolor_i$$

- this will result in three unique lines depending on the species of the flower.
- both the intercepts and the slopes will be allowed to be different.

Does it make sense to retain the interaction?

```
# A tibble: 6 x 5
  term                estimate std.error statistic  p.value
  <chr>              <dbl>    <dbl>    <dbl>    <dbl>
1 (Intercept)        0.803     0.531     1.51  0.133
2 Sepal.Length        0.132     0.106     1.24  0.216
3 Speciesversicolor  -0.618     0.684    -0.904 0.368
4 Speciesvirginica   -0.193     0.658    -0.293 0.770
5 Sepal.Length:Speciesversicolor  0.555     0.128     4.33 0.0000278
6 Sepal.Length:Speciesvirginica   0.618     0.121     5.11 0.00000100
```

DOES IT MAKE SENSE TO RETAIN THE INTERACTION?



```
broom::tidy(inter) |> kable()
```

term	estimate	std.error	statistic	p.value
(Intercept)	0.803	0.531	1.512	0.133
Sepal.Length	0.132	0.106	1.244	0.216
Speciesversicolor	-0.618	0.684	-0.904	0.368
Speciesvirginica	-0.193	0.658	-0.293	0.770
Sepal.Length:Speciesversicolor	0.555	0.128	4.330	0.000
Sepal.Length:Speciesvirginica	0.618	0.121	5.111	0.000

```
anova(nointer, dum, inter)
```

Analysis of Variance Table

Model 1: Petal.Length ~ Sepal.Length

Model 2: Petal.Length ~ Sepal.Length + Species

Model 3: Petal.Length ~ Sepal.Length + Species + Sepal.Length:Species

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	148	111.5				
2	146	11.7	2	99.8	731.9	< 0.0000000000000002 ***
3	144	9.8	2	1.8	13.5	0.0000043 ***

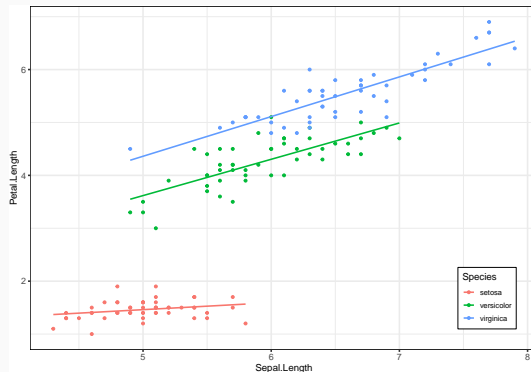
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

NOW WE CAN ADD AN INTERACTION



$$\widehat{petal.length}_i = \beta_0 + \beta_1 sepal.length_i + \beta_2 setosa_i + \beta_3 versicolor_i + \beta_4 sepal.length_i setosa_i + \beta_5 sepal.length_i versicolor_i$$

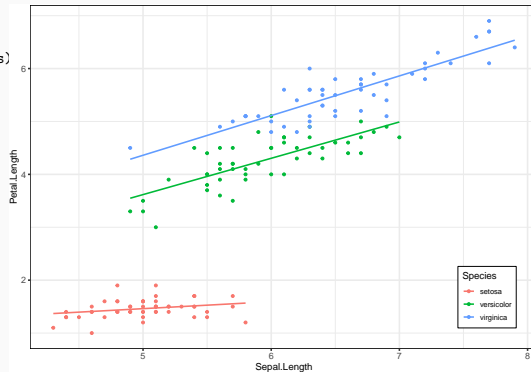
- this will result in three unique lines depending on the species of the flower.
- both the intercepts and the slopes will be allowed to be different.
- ggplot geom_smooth does this by default if color is used



NOW WE CAN ADD AN INTERACTION



```
ggplot(iris, aes(x = Sepal.Length, y = Petal.Length, color = Species))  
  geom_point() +  
  geom_smooth(method = "lm", se = FALSE) +  
  theme_luis() +  
  theme( legend.position = c(0.9, 0.15))
```



VISUALIZING INTERACTIONS



```
fit <- glm(qsec ~ wt*as.factor(cyl), data = mtcars)
summary(fit)
```

- Note: not significant, later try

```
summary(margins(fit))
```

Call:

```
glm(formula = qsec ~ wt * cyl, data = mtcars)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.1966	-0.8373	0.0499	0.8158	2.1398

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	16.726	3.118	5.36	0.00001 ***
wt	2.858	1.180	2.42	0.022 *
cyl	-0.542	0.511	-1.06	0.298
wt:cyl	-0.222	0.167	-1.33	0.193

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

(Dispersion parameter for gaussian family taken to be 1.45)

Null deviance: 98.988 on 31 degrees of freedom
Residual deviance: 40.636 on 28 degrees of freedom
AIC: 108.5

Number of Fisher Scoring iterations: 2



```
fit <- glm(mpg ~ wt*as.factor(cyl), data = mtcars)
summary(fit)
```

Call:

```
glm(formula = mpg ~ wt * as.factor(cyl), data = mtcars)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-4.151	-1.380	-0.639	1.494	5.252

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	39.57	3.19	12.39	0.00000000000021 ***
wt	-5.65	1.36	-4.15	0.00031 ***
as.factor(cyl)6	-11.16	9.36	-1.19	0.24358
as.factor(cyl)8	-15.70	4.84	-3.24	0.00322 **
wt:as.factor(cyl)6	2.87	3.12	0.92	0.36620
wt:as.factor(cyl)8	3.45	1.63	2.12	0.04344 *

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

(Dispersion parameter for gaussian family taken to be 6)

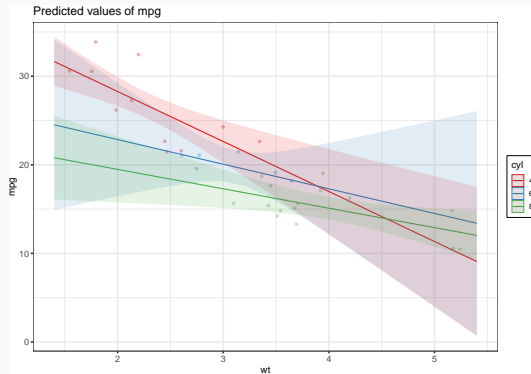
Null deviance: 1126.05 on 31 degrees of freedom
Residual deviance: 155.89 on 26 degrees of freedom
AIC: 155.5

Number of Fisher Scoring iterations: 2

REGRESSION OF MPG ON WT*CYL



```
pred <- ggpredict(fit, terms = c("wt", "cyl"))  
plot(pred, add.data = TRUE)+  
  theme_luis()
```



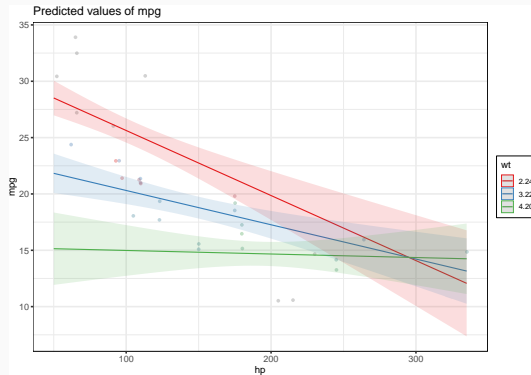
NOW, FROM WHAT POINT THE SLOPE BECOMES NON SIGNIFICANT?



Changing to $\text{mpg} \sim \text{hp} + \text{wt}$

```
fit <- glm(mpg ~ hp*wt, data = mtcars)
pred <- ggpredict(fit, terms = c("hp", "wt"))

plot(pred, add.data = TRUE) +
  theme_luis()
```

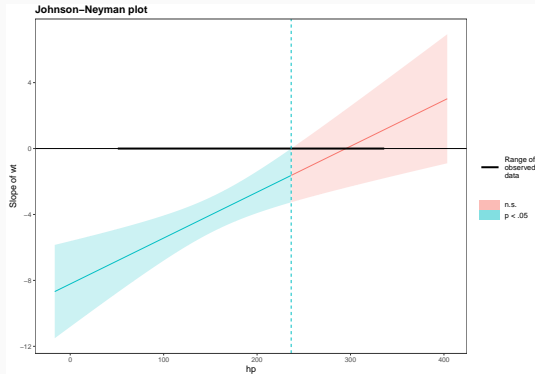


NOW, FROM WHAT POINT THE SLOPE BECOMES NON SIGNIFICANT?



JOHNSON-NEYMAN INTERVAL

```
jn <- johnson_neyman(fit, wt, hp, plot = TRUE)
jn
```



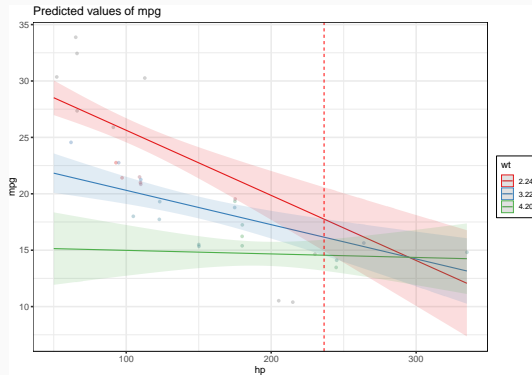
NOW, FROM WHAT POINT THE SLOPE BECOMES NON SIGNIFICANT?



JOHNSON-NEYMAN INTERVAL - Overlaid over data

```
fit <- glm(mpg ~ hp*wt, data = mtcars)
pred <- ggpredict(fit, terms = c("hp", "wt"))
jn<-johnson_neyman(fit, wt, hp , plot = TRUE)
jn_bound<-as.numeric(jn$bounds[1])

plot(pred, add.data=T) +
  geom_vline(xintercept = jn_bound, linetype = "dashed",
            color = "red")+
  theme_luis()
```



THREE-WAY INTERACTIONS



```
fit <- glm(mpg ~ hp*wt*cyl, data = mtcars)
```

```
dat <- ggpredict(fit, terms = c("hp", "wt", "cyl"))  
plot(dat, ci = FALSE)
```

Call:

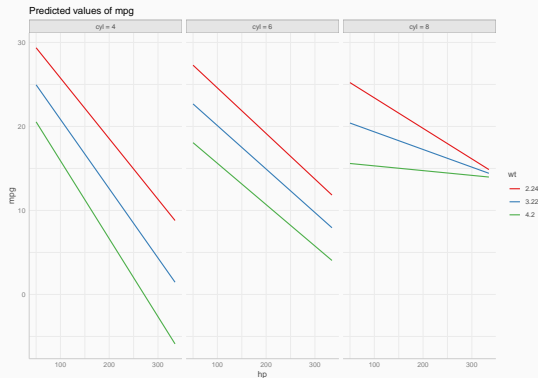
```
glm(formula = mpg ~ hp * wt * cyl, data = mtcars)
```

Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-3.352	-1.464	-0.169	1.345	4.001

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	43.96543	30.32070	1.45	0.16
hp	-0.02587	0.24000	-0.11	0.92
wt	-2.25515	10.68401	-0.21	0.83
cyl	-0.52189	6.33725	-0.08	0.94
hp:wt	-0.03666	0.09360	-0.39	0.70
hp:cyl	-0.00569	0.03850	-0.15	0.88
wt:cyl	-0.42991	1.99058	-0.22	0.83
hp:wt:cyl	0.00654	0.01375	0.48	0.64





- Marginal effects: the effect of a predictor on the response, averaged over all values of the other predictors.
- It is achieved by..

WHAT ARE THE MARGINAL EFFECTS OF THE LATEST MODEL?



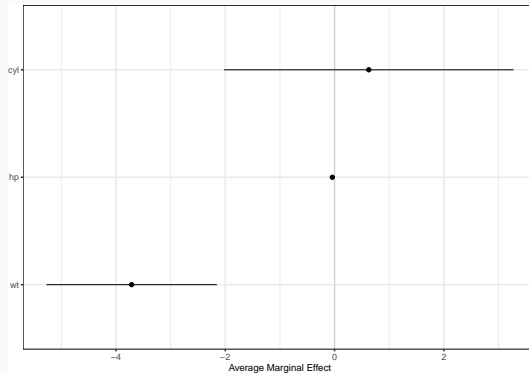
```
fit_m <- margins(fit)
summary(fit_m)
```

factor	AME	SE	z	p	lower	upper
cyl	0.6261	1.3513	0.4633	0.6431	-2.0224	3.2745
hp	-0.0402	0.0152	-2.6390	0.0083	-0.0700	-0.0103
wt	-3.7134	0.7952	-4.6696	0.0000	-5.2720	-2.1547

PLOTTING THE MARGINAL EFFECTS



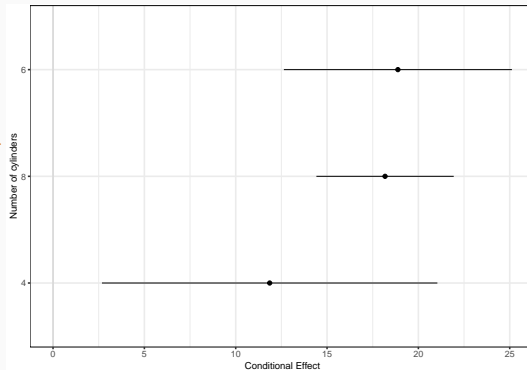
```
fit_mo <- as_tibble(summary(fit_m))  
p <- ggplot(data = fit_mo, aes(x = reorder(factor, AME),  
                               y = AME, ymin = lower, ymax = upper))  
  
p + geom_hline(yintercept = 0, color = "gray80") +  
  geom_pointrange() + coord_flip() +  
  labs(x = NULL, y = "Average Marginal Effect") +  
  theme_luis()
```



CAN WE PLOT THE CONDITIONAL EFFECTS TOO?



```
glm(mpg ~ hp*wt*as.factor(cyl), data = mtcars) %>%  
  cplot(x = "cyl", draw = F) %>%  
  ggplot( aes(x = reorder(xvals, yvals),  
              y = yvals, ymin = lower, ymax = upper)) +  
    geom_hline(yintercept = 0, color = "gray80") +  
    geom_pointrange() + coord_flip() +  
  labs(x = "Number of cylinders", y = "Conditional Effect") +  
  theme_luis()
```





- At the beginning I said, keep the interaction that is significant, but¹
 - In the conversion to probabilities (AME) the interaction may not be significant anymore, or worse
 - The interaction may be significant in the AME, but not in the original model

¹Bruin, "Deciphering Interactions in Logistic Regression," *Introduction to SAS. UCLA: Statistical Consulting Group.*; Vanhove, 2019, "Interactions in Logistic Regression Models," (2019).



```
library(umx)

model <- "
  mpg ~ a*hp + b*wt
  moderation := a*b
"

umxRAM(model, data = mtcars, tryHard="ordinal", silent = F)
```

	name	Estimate	SE	type
5	hp_with_wt	42.812	13.758	Manifest
1	a	-0.032	0.009	Manifest
2	b	-3.878	0.602	Manifest
7	one_to_mpg	37.227	1.522	Mean
8	one_to_hp	146.687	11.929	Mean
9	one_to_wt	3.217	0.170	Mean
3	mpg_with_mpg	6.095	1.524	Residual
4	hp_with_hp	4553.963	1138.504	Residual
6	wt_with_wt	0.927	0.232	Residual

Model Fit: Chi2(0) = 0, p = 1.000; CFI = 1; TLI = 1; 38

RMSEA = 0.000 (no moderation)



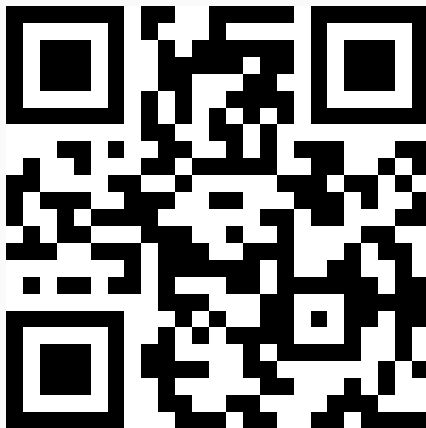
- We started reviewing multiple regression
- Then discussed the syntax and interpretation of parameters when an interaction term is included
- Finally, we discussed how to extract the marginal effects of the interaction term
- Luckily the package `margins()` makes this extremely simple.



Team

- Charles Gardner (2015)
 - Brad Verhulst (2013)
 - Joshua Pritkin.
 - Rob Kirkpatrick.
-
- Michael C Neale.
 - NIH grant no R01 DA049867 and 5T32MH-020030

Contact



- **THANK YOU**