

MV Regression: Indicators and Interaction Terms

Consider a variable that takes value 0 or 1 that want to add to a regression model

Ex: Sex of penguins

Example: Male 0 or 1

1.7% of world popⁿ estimated to be intersex (people, not penguins)

$$y = b_0 + b_1 x + b_2 I_{\text{Male}_i} + \epsilon_i$$

Annotations:
 b_0 : weight
 $b_1 x$: bill length
 I_{Male_i} : value 1 for males, 0 otherwise
 ϵ_i : error $N(0, \sigma)$

This is an indicator variable

For a male penguin the equation becomes

$$y_i = b_0 + b_1 x_i + b_2 \times 1 + \epsilon_i \Rightarrow y_i = (b_0 + b_2) + b_1 x_i + \epsilon_i$$

For a female (non-male) the equation becomes

$$y_i = b_0 + b_1 x_i + b_2 \times 0 + \epsilon_i \Rightarrow y_i = b_0 + b_1 x_i + \epsilon_i$$

Essentially, male penguins have a different intercept in the equation.

What if we want to use more than 2 categories?

Species: Adelie, Gentoo, Chinstrap

Make 3 indicators

$$y_i = b_0 + b_1 x_i + b_2 I_{A_i} + b_3 I_{G_i} + b_4 I_{C_i} + \epsilon_i$$

Adelie

Gentoo

Chinstrap

..

..

$$y_i = (b_0 + b_2) + b_1 x_i + \epsilon_i$$

Adelie

Gentoo

Chinstrap

$$y_i = (b_0 + b_2) + b_1 x_i + \epsilon_i$$

$$y_i = (b_0 + b_3) + b_1 x_i + \epsilon_i$$

$$y_i = (b_0 + b_4) + b_1 x_i + \epsilon_i$$

But... we don't need 4 different parameters to describe 3 different intercepts.

We choose a baseline (Adelie) and remove the indicator for the baseline. [We roll the old b_2 into the intercept term]

$$y_i = b_0 + b_1 x_i + b_2 I_{G_i} + b_3 I_{C_i} + \epsilon_i$$

$$y_i = b_0 + b_1 x_i + \epsilon_i$$

What if we want different slopes for each species?

We can use interaction terms.

value of 0 or 1

Equation becomes

$$y_i = b_0 + b_1 I_{G_i} + b_2 I_{C_i} + b_3 x_i + b_4 x_i I_{G_i} + b_5 x_i I_{C_i} + \epsilon_i$$

3 equations

$$A: y_i = b_0 + b_3 x_i + \epsilon_i$$

$$G: y_i = (b_0 + b_1) + (b_3 + b_4) x_i + \epsilon_i$$

$$C: y_i = (b_0 + b_2) + (b_3 + b_5) x_i + \epsilon_i$$

represents diff in slope for Gentoo penguins compared to Adelie.

diff for chinstrap compared to Adelie.

We can also have interactions between continuous variables

Ex bike rentals we might wonder how temp and humidity affect rentals on their own,

how they interact together.

and humidity affect returns on ...
but also how they interact together:

$$Y_i = b_0 + b_1 \text{Temp} + b_2 \text{Humidity} + b_3 \text{Temp} * \text{Humidity} + \varepsilon_i$$