

①

Regression Models with categorical predictor variables

common examples:

- sex (male / female)
- SES (low / middle / high)
- intervention (treatment / control)

example: $y = \text{weight}$

(turkey data) age is a continuous (numerical) input.

(Georgia
Virginia
Wisconsin) state of origin is a categorical input (g, v, w)

data:

group		
<u>1</u>	<u>2</u> ...	<u>c</u>
(X_{11}, Y_{11})	(X_{21}, Y_{21})	(X_{c1}, Y_{c1})
\vdots	\vdots ... \vdots	\vdots
(X_{1n_1}, Y_{1n_1})	(X_{2n_2}, Y_{2n_2})	(X_{cn_c}, Y_{cn_c})

We will model the categorical input using indicator variables.

$$\text{Let } I_1 = 1\{g\} = \begin{cases} 1, & \text{if state} = \text{'georgia'} \\ 0, & \text{otherwise} \end{cases}$$

$$I_2 = 1\{v\} = \begin{cases} 1, & \text{if state} = \text{'virginia'} \\ 0, & \text{otherwise} \end{cases}$$

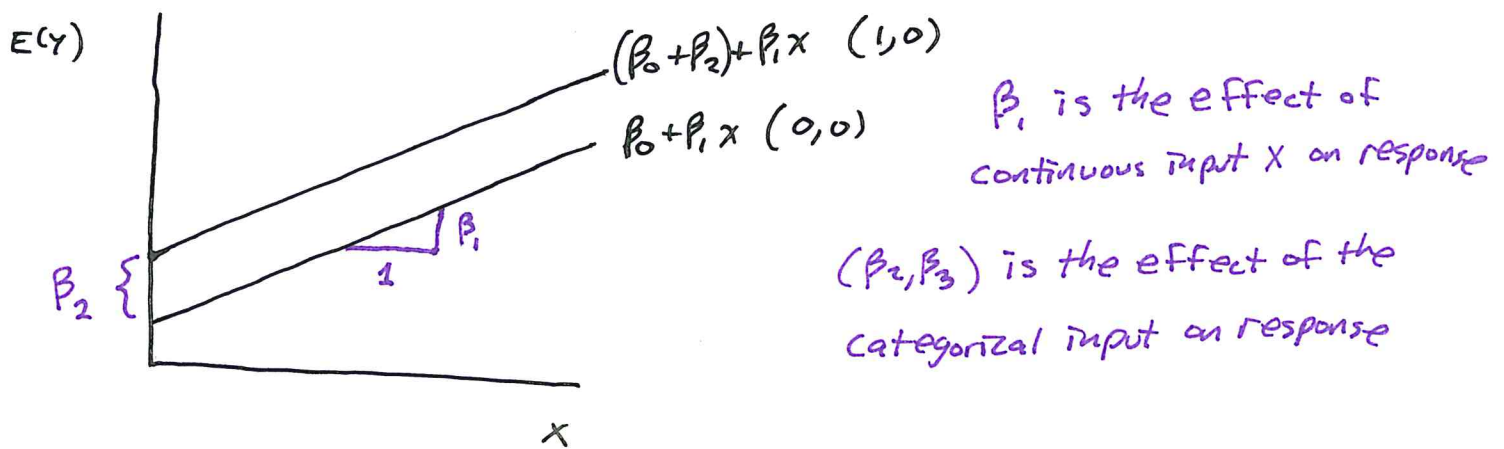
$$(I_1, I_2) = (1, 0) \text{ if } g, (0, 1) \text{ if } v, (0, 0) \text{ if } w$$

②

model : $E(Y) = \beta_0 + \beta_1 X + \beta_2 I_1 + \beta_3 I_2$

$$\Leftrightarrow E(Y) = \begin{cases} (\beta_0 + \beta_2) + \beta_1 X & , \text{ if } g \\ (\beta_0 + \beta_3) + \beta_1 X & , \text{ if } v \\ \beta_0 + \beta_1 X & , \text{ if } w \end{cases}$$

(Note that we use one less indicator variable than number of categories)



parameter interpretation :

$\beta_1 = \frac{\partial E(Y)}{\partial X_1}$. β_1 is the difference in mean weight from a 1 year increase in age, with origin held constant.

(effect of the continuous input)

$\beta_2 = E(Y | g, X) - E(Y | w, X)$ (effect of categorical input)
(is defined by multiple parameters)

β_2 is the difference in mean weight between Georgia and Wisconsin, with age held constant.

③

$$\beta_3 = E(Y | v, X) - E(Y | w, X) \quad \left(\begin{array}{l} \text{diff. in mean weight} \\ \text{between } v \text{ and } w. \end{array} \right)$$

$$\beta_2 - \beta_3 = E(Y | g, X) - E(Y | v, X) \quad \left(\begin{array}{l} \text{diff. in mean weight} \\ \text{between } g \text{ and } v \end{array} \right)$$

Test for categorical input effects: $(H_0: \beta_2 = \beta_3 = 0)$

$$(F): E(Y) = \beta_0 + \beta_1 X + \beta_2 \mathbb{I}_1 + \beta_3 \mathbb{I}_2 \quad \left[\begin{array}{l} m_A = \text{lm}(Y \sim X + \text{group}) \\ \text{OR} \\ m_A = \text{lm}(Y \sim X + i1 + i2) \end{array} \right]$$

$$(R): E(Y) = \beta_0 + \beta_1 X$$

$$\left[m_L = \text{lm}(Y \sim X) \right]$$

$$\underline{\text{anova}(m_L, m_A)}, \quad F^* = 68.81, \quad p = .000$$

[we accept the model which includes origin as a categorical predictor for turkey weight. (test for partial effect)]

Estimating categorical predictor effects: confint(m_A)

$$g-w: \quad \underline{CI \text{ for } \beta_2 = [-2.375, -1.462]}, \quad \underline{t_2 = -9.506}$$

data supports the hypothesis that Georgia turkeys weigh less than Wisconsin turkeys

$$v-w: \quad \underline{CI \text{ for } \beta_3 = [-2.67, -1.71]}, \quad \underline{t_3 = -10.367}$$

data supports the hypothesis that Virginia turkeys weigh less than Wisconsin turkeys

④

$$g-v : \text{CI for } \beta_2 - \beta_3 = [-0.22, 0.77] , \quad \underline{t_{23} = 1.25} \\ (\underline{p_{23} = .242})$$

[The data is compatible with a model that combines g, v into a single group.]

some quick distribution theory:

$$\underline{b} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} , \quad \underline{a}' = [0 \ 0 \ 1 \ -1]$$

$$\underline{b} \sim N(\underline{\beta}, \sigma^2 (X'X)^{-1}) , \quad E(\underline{a}'\underline{b}) = \underline{a}'\underline{\beta}$$

$$\text{Var}(\underline{a}'\underline{b}) = \sigma^2 \underline{a}'(X'X)^{-1}\underline{a}$$

so,

$$\text{SE}(\underline{a}'\underline{\beta}) = \sqrt{\underline{a}' \hat{\text{Cov}}(\underline{b}) \underline{a}} , \quad \underline{\text{CI for } \underline{a}'\underline{\beta} = \underline{a}'\underline{b} \pm \underline{t}_{\frac{\alpha}{2}}^{(n-p)} \text{SE}(\underline{a}'\underline{b})}$$

Estimating continuous predictor effects

$$\underline{\text{CI for } \beta_1 = [0.43, 0.545]} , \quad \underline{t = 18.908}$$

We accept the model which includes age as a continuous predictor for turkey weight. (test for partial effect)

Interaction Model :

$$E(Y) = \beta_0 + \beta_1 X + \beta_2 I_1 + \beta_3 I_2 + \beta_4 X \cdot I_1 + \beta_5 X \cdot I_2$$

<u>table:</u>	<u>group</u>	<u>intercept</u>	<u>slope</u>
	1: (1,0)	$\beta_0 + \beta_2$	$\beta_1 + \beta_4$
	2: (0,1)	$\beta_0 + \beta_3$	$\beta_1 + \beta_5$
	3: (0,0)	β_0	β_1

The interaction model implies that the regression functions will have unequal slopes.

(effect of continuous input depends on the level of the categorial input)

test for interaction effect : ($H_0: \beta_4 = \beta_5 = 0$).

(F): interaction model $mI = \text{lm}(y \sim x + i1 + i2 + I(x*i1) + I(x*i2))$

(R): additive model $mA = \text{lm}(y \sim x + i1 + i2)$

turkey data: $F^* = 0.5204$, $p = .6156$ $\text{anova}(mA, mI)$

~~the observed data is compatible with the~~

The observed data is compatible with the ~~observed data~~ additive model

see scatterplot