

## Factor + numeric covariate

one factor with  $J$  levels

one numeric covariate,  $x_i$

$$\text{model: } Y_i = b_0 + b_{j(i)} + c_1 x_i + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$E(Y_i) = b_0 + b_{j(i)} + c_1 x_i$$

interpretation:  $c_1$  = effect (slope) of  $x_i$

$$b_0 + b_{j(i)} = \text{intercept} = E(Y_i) \text{ when } x_i = 0$$

Example:  $y_i$  = 3pm temperature on one day in 1 of 3 cities

$x_i$  = 7am temperature in same city on same day

$j(i) = 1, 2, \text{ or } 3$  (city)

$b_0 + b_1$  = expected 3pm temp when 7am temp is 0 in city 1

$b_0 + b_2$  = expected 3pm temp when 7am temp is 0 in city 2

$b_0 + b_3$  = expected 3pm temp when 7am temp is 0 in city 3

increasing 7am temp by 1 degree means that we expect 3pm temp to change by  $c_1$  degrees

Note that the slope  $c_1$  does not depend on city.

Additive Model: effect of one variable does not depend on the value of the other

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & x_1 \\ 1 & 1 & 0 & 0 & x_2 \\ 1 & 0 & 1 & 0 & x_3 \\ 1 & 0 & 1 & 0 & x_4 \\ 1 & 0 & 0 & 1 & x_5 \\ 1 & 0 & 0 & 1 & x_6 \end{bmatrix} \quad \text{rank}(X) = 4 \text{ (probably)}$$

Drop intercept column (set  $b_0 = 0$ )

$b_0 = 0$  (no interpretation)

$b_1 = \text{intercept for city 1}$

$b_2 = \text{intercept for city 2}$

$b_3 = \text{intercept for city 3}$

Drop level 1 column (set  $b_1 = 0$ , R default)

$b_0 = \text{level 1 intercept} \quad (b_0 + b_1)$

$b_1 = 0$

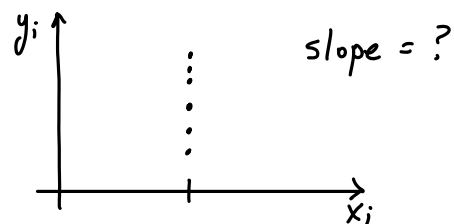
$b_2 = \text{level 2 intercept} - \text{level 1 intercept} \quad (b_2 - b_1)$

$b_3 = \text{level 3 intercept} - \text{level 1 intercept} \quad (b_3 - b_1)$

### Digression on multicollinearity

In the multiple linear model, there is no assumption about the relationships between  $x_{i0}$ ,  $x_{i1}$ ,  $x_{i2}$ , etc. However, correlation in the design matrix can affect our ability to estimate things

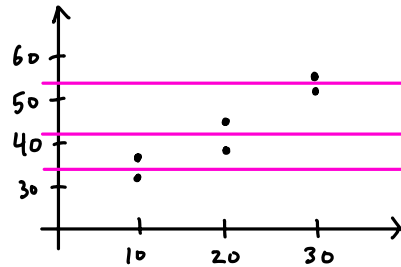
simple linear model:  $X = \begin{bmatrix} 1 & 9.2 \\ 1 & 9.2 \\ 1 & 9.2 \\ 1 & 9.2 \\ 1 & 9.2 \\ 1 & 9.2 \end{bmatrix}$



in that example  $\text{rank}(X) = 1$ ,  $(X^T X)^{-1}$  does not exist,  $\text{Var}(\hat{\beta}_1) = \infty$   
 we need variation in the covariate to estimate a slope.

city temperature example:

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 10 \\ 1 & 1 & 0 & 0 & 10 \\ 1 & 0 & 1 & 0 & 20 \\ 1 & 0 & 1 & 0 & 20 \\ 1 & 0 & 0 & 1 & 30 \\ 1 & 0 & 0 & 1 & 30 \end{bmatrix} \quad y = \begin{bmatrix} 32 \\ 37 \\ 45 \\ 39 \\ 52 \\ 55 \end{bmatrix}$$

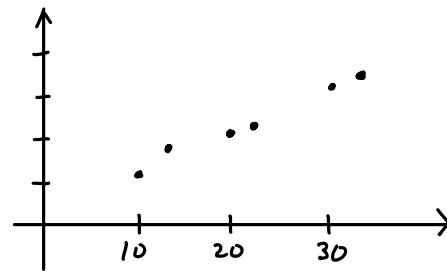


$$\text{rank}(X) = 3$$

Can't tell how much of the changes in  $y$  are due to  
 the intercepts versus the slope

Need variation in 7am temp that's not collinear with city.

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 10 \\ 1 & 1 & 0 & 0 & 12 \\ 1 & 0 & 1 & 0 & 20 \\ 1 & 0 & 1 & 0 & 22 \\ 1 & 0 & 0 & 1 & 30 \\ 1 & 0 & 0 & 1 & 32 \end{bmatrix} \quad y = \begin{bmatrix} 32 \\ 39 \\ 41 \\ 42 \\ 52 \\ 55 \end{bmatrix}$$



$$\text{rank}(X) = 4$$

A small positive slope will explain the data better than  $c_1 = 0$

In fact,  $\hat{c}_1 = \frac{1}{3} \left( \frac{y_2 - y_1}{2} + \frac{y_4 - y_3}{2} + \frac{y_6 - y_5}{2} \right)$

## Adding another factor

factor 1 has  $J$  levels: (repub., dem, other)

factor 2 has  $K$  levels: (female, male)

$y_i$  numeric response: (income)

model:  $Y_i = b_0 + b_{j(i)} + c_{k(i)} + \epsilon_i, \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

Example: 6 subjects.  $X \underline{b} = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \\ c_1 \\ c_2 \end{bmatrix}$

		$k$	
		1 (F)	2 (M)
$j$	1 (R)	$E(Y_1)$	$E(Y_2)$
	2 (D)	$E(Y_3)$	$E(Y_4)$
	3 (I)	$E(Y_5)$	$E(Y_6)$

		1	2
$j$	1	$b_0 + b_1 + c_1$	$b_0 + b_1 + c_2$
	2	$b_0 + b_2 + c_1$	$b_0 + b_2 + c_2$
	3	$b_0 + b_3 + c_1$	$b_0 + b_3 + c_2$

can interpret quantities that are expected values of linear combinations of observations.

$$E(Y_3) - E(Y_1) = (b_0 + b_2 + c_1) - (b_0 + b_1 + c_1) = b_2 - b_1$$

$$E(Y_4) - E(Y_2) = (b_0 + b_2 + c_2) - (b_0 + b_1 + c_2) = b_2 - b_1$$

$$E(Y_6) - E(Y_5) = (b_0 + b_3 + c_2) - (b_0 + b_3 + c_1) = c_2 - c_1$$

$$E(Y_2) - E(Y_1) = (b_0 + b_1 + c_2) - (b_0 + b_1 + c_1) = c_2 - c_1$$

Note that the effect of changing from  $j=1$  to  $j=2$  is  $b_2 - b_1$ , regardless of the level of factor 2.

Same for changing from  $k=1$  to  $k=2$  ( $c_2 - c_1$ )

Effect of changing one of the factor levels does not depend on the level of the other factor.

Another example of an additive model

$\text{rank}(X) = 4$ , yet we have 6 columns and 6 parameters need to drop 2 columns (e.g. set 2 parameters to 0)

R default:  $b_1 = 0$ ,  $c_1 = 0$

		k	
		1	2
j	1	$b_0$	$b_0 + c_2$
	2	$b_0 + b_2$	$b_0 + b_2 + c_2$
	3	$b_0 + b_3$	$b_0 + b_3 + c_2$

$b_0 = \text{EV for } j=1, k=1$

$b_2 = \text{EV for } (2,1) - (1,1) \quad (\text{or } (2,2) - (1,2))$

$b_3 = \text{EV for } (3,1) - (1,1) \quad (\text{or } (3,2) - (1,2))$

$c_2 = \text{EV for } (1,2) - (1,1) \quad (\text{or } (2,2) - (2,1))$   
 $(\text{or } (3,2) - (3,1))$

$k$

1      2

$j$

1	*	*
2	*	
3	*	

Model dof = 4 < 6 = # of cells  
 consequence: can complete the  
 table with incomplete information.

1    2    3    4    5    6

1	*	*	*	*	*	*
2	*					
3	*					
4	*					

In general, model has

$$J + K - 1$$

model d.o.f.

## Multiple subscript notation

So far:  $i \in \{1, \dots, n\}$  = identifier for individual observation

$j(i) \in \{1, \dots, J\}$  = factor 1 level for obs.  $i$

$k(i) \in \{1, \dots, K\}$  = factor 2 level for obs.  $i$

$$\text{Model: } Y_i = b_0 + \underbrace{\sum_{\ell=1}^p b_{\ell} x_{i\ell}}_{\text{numeric}} + \underbrace{c_{j(i)}}_{\text{fac 1}} + \underbrace{d_{k(i)}}_{\text{fac 2}} + \underbrace{\varepsilon_i}_{\text{error}}$$

We can also write it like this:

$i$  = identifier for obs. with levels  $j$  and  $k$  of fac 1 and 2

$i$	$j$	$k$	$y_{ijk}$
1	1	1	..
2	1	1	..
1	1	2	..
2	1	2	..
1	2	1	..
2	2	1	..
1	2	2	..
2	2	2	..

$$\text{Model: } Y_{ijk} = b_0 + c_j + d_k + \varepsilon_{ijk}, \quad \varepsilon_{ijk} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

The model is the same, but notation is different

$i$ ,  $j$ , and  $k$  together uniquely identify an observation

In single subscript notation ( $y_i$ ),  $i$  uniquely identifies an observation.