

Summary of Week 6

- linear regression w/ > 1 predictor
- interpretation of regression coefficients
 - \hookrightarrow difference between single and multiple predictors

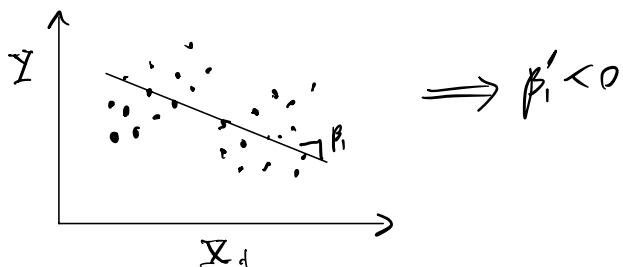
$$\beta_1' = \beta_1 + \beta_2 (E[Y|X_2=x+1] - E[Y|X_2=x])$$
- R^2 , p-values etc.
- sample distribution / collinearity \rightarrow want predictors to be uncorrelated

this week

- Simpson's paradox
- More on math behind multiple predictor regression
- Categorical predictors
- Nonlinear models:
 - * interactions between predictors
e.g. $Y = X_1 X_2 + \epsilon$
 - * Fourier modes e.g. $Y = \sin(x) + \epsilon$

Simpson's paradox

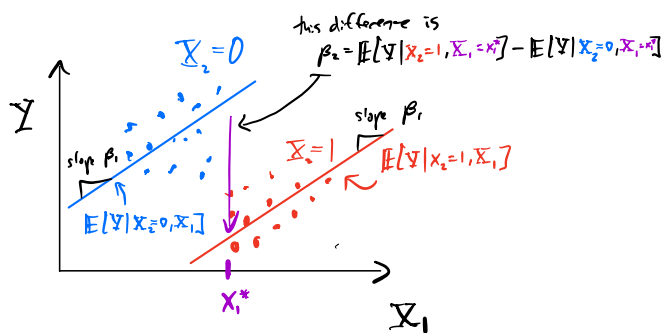
Look at outcome Y vs predictor X_1



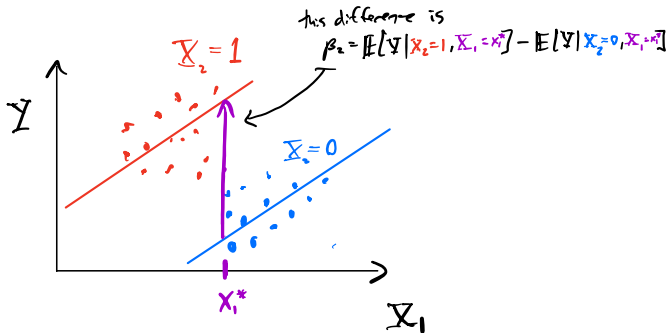
Now add predictor X_2 $\beta_1 > 0$. How can this happen?

$$0 > \beta_1' = \beta_1 + \beta_2 \underbrace{(E[X_2|X_1=x+1] - E[X_2|X_1=x])}_{=b} = \beta_1 + \beta_2 b$$

Case 1: $\beta_2 < 0$ AND $b > 0$



Case 2: $\beta_2 > 0$ AND $b < 0$



More Math (see page 6 in week 6 notes)

How are regression coefficients actually computed?

Use covariance like we did with one predictor.

to make things simpler, assume $E[X_i] = 0$, $\beta_0 = 0$

$$\begin{aligned} \text{Cov}(X_j, Y) &= \text{Cov}\left(X_j, \sum_{i=1}^K \beta_i X_i\right) \\ &= E\left[X_j, \sum_{i=1}^K \beta_i X_i\right] = \sum_{i=1}^K \beta_i E[X_i X_j] \\ &= \sum_{i=1}^K \beta_i \text{Cov}(X_i, X_j) \end{aligned}$$

get system of K equations w/ K unknowns β_1, \dots, β_K

Replace expectations w/ sample avg. using N data points

$$\begin{aligned} \mathbf{X} &= \begin{bmatrix} X_{1,1}, \dots, X_{1,K} \\ \vdots \\ X_{N,1}, \dots, X_{N,K} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix}, \quad \vec{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_K \end{bmatrix} \\ \mathbf{X}^T \mathbf{y} &= \begin{bmatrix} \sum_{i=1}^N Y_i X_{i,1} \\ \vdots \\ \sum_{i=1}^N Y_i X_{i,K} \end{bmatrix} = \underbrace{\mathbf{X}^T \mathbf{X}}_{\mathbf{K} \times \mathbf{K}} \vec{\beta} \implies \boxed{\vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}} \quad \star \end{aligned}$$

Optional material - not on exam

The matrix $X^T X$ is called the covariance matrix

Note that

$$\begin{bmatrix} E[Y|X=x_1] \\ \vdots \\ E[Y|X=x_n] \end{bmatrix} = X \hat{\beta}$$

(x_{1,1}, ..., x_{1,k})

is vector of predicted averages conditional on values of response variable

formula (8) also come from least squares i.e. $\hat{\beta} = \min \|y - X\beta\|_2^2$

Key point: ideally want $X^T X \approx I$
 $\Rightarrow \text{Cov}(X_i, X_j) = 0 \quad i \neq j$

Categorical predictors

Say we have a predictor like X = someone's race on a survey where you can only select one race

then $X \in \{\text{Black, white, asian, hispanic, other}\}$

↓	↓	↓	↓	↓
0	1	2	3	4

← Can't do this because it implies intrinsic ordering. E.g. difference between black \rightarrow asian less than white \rightarrow other

Instead we break up into binary variables: $X_{\text{black}}, X_{\text{white}}, \dots$

and consider

$$Y = \beta_0 + \beta_{\text{black}} X_{\text{black}} + \beta_{\text{white}} X_{\text{white}} + \beta_{\text{asian}} X_{\text{asian}} + \beta_{\text{hispanic}} X_{\text{hispanic}} + \beta_{\text{other}} X_{\text{other}} + \epsilon$$

← but this has new problem:

$$\beta_{\text{black}} = E[Y | X_{\text{black}}=1, \text{all others } 0] - E[Y | X_{\text{black}}=0, \text{all others } 0]$$

this doesn't have sense because one of the $X=1$

to resolve this we drop one of the predictors.

Python drops the first one in alphabetical order:

$$Y = \beta_0 + \beta_{\text{black}} X_{\text{black}} + \beta_{\text{white}} X_{\text{white}} + \beta_{\text{Asian}} X_{\text{Asian}} + \beta_{\text{hispanic}} X_{\text{hispanic}} + \beta_{\text{other}} X_{\text{other}} + \epsilon$$

Selected as baseline

$$Y = \beta_0 + \beta_{\text{black}} X_{\text{black}} + \beta_{\text{white}} X_{\text{white}} + \beta_{\text{hispanic}} X_{\text{hispanic}} + \beta_{\text{other}} X_{\text{other}} + \epsilon$$

in this model okay to have all predictors = 0 because this implies $X_{\text{Asian}} = 1$ which is not in the model

Now $\beta_{\text{black}} = E[Y | X_{\text{black}} = 1, \text{all others} = 0] - E[Y | X_{\text{black}} = 0, \text{all others} = 0]$

all $X=0 \Rightarrow$ this is avg. of Y among those samples who selected asian on survey

all regression coefficients should be interpreted as avg. differences between given category (e.g. black for black) and baseline category (e.g. asian)