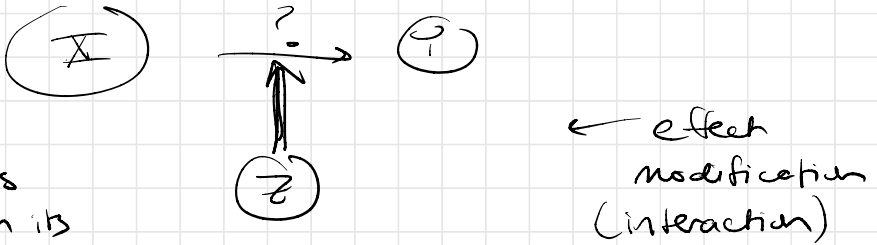
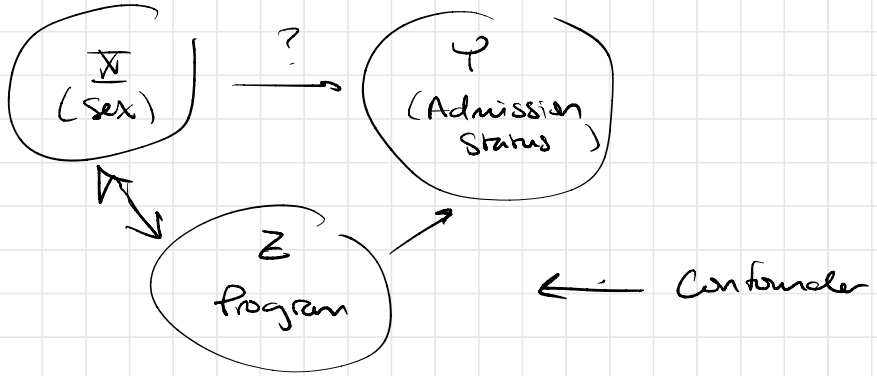


# Simpson's Paradox

2/15/22

Example: UC Berkeley Admissions



$Z$  interacts with  $X$  on its effect on  $Y$

Marginal Independence between  $X \rightarrow Y$  :

$$\rightarrow P(X=x | Y=y) = P(X=x) \quad \forall x, y$$

$$P(X=x \wedge Y=y) = P(X=x) P(Y=y)$$

Independent conditional on  $W$

$$P(X=x | Y=y, W=w) = P(X=x | W=w) \quad \forall x, y, w$$

PPT Ex:  $X \rightarrow Y$  indep given  $W=1$ ?

$$P(Y=1 | \underbrace{X=1}, W=1) \stackrel{?}{=} P(Y=1 | W=1)$$

or Odds ratio / LR / diff. prop.

$$\begin{array}{ccc} \swarrow & \searrow & \rightarrow \\ \theta = 1 & \theta = 1 & \theta = 0 \end{array}$$

Odds ratios: (conditional)

$$\hat{\theta}_{X \rightarrow Y | W=1} = \frac{4(9)}{6(6)} = 1$$

$$\hat{\theta}_{X \rightarrow Y | W=2} = \frac{16(1)}{4(4)} = 1$$

"independent"

$\downarrow$   
 $X \perp\!\!\!\perp Y | W$

Marginal independence: Sum over  $W$ :

	$Y=1$	$Y=2$
$X=1$	20	10
$X=2$	10	10

$$\hat{\theta}_{X \rightarrow Y} = 2$$

$\Rightarrow$   
 $X \not\perp\!\!\!\perp Y$

W an effect modifier? W interact w/ X on Y

⇒ Association between X-Y changes w/ values of W

$$\hat{\Delta}_{X \rightarrow Y | W=1} = \hat{\Delta}_{X \rightarrow Y | W=2} = 1$$

→ No

- Homogeneous association between X-Y given W

W a confounder?

- W associated w/ X?

- W associated w/ Y?

	W=1	W=2
X=1		
X=2		

$$\hat{\Delta}_{X \rightarrow W} = 1/6$$

→ Yes

	W=1	W=2
Y=1		
Y=2		

$$\hat{\Delta}_{Y \rightarrow W} = 1/6$$

	Y	
	1	0
X		
1		
2		
3		

$$\begin{array}{l} \text{ORs} \\ \text{odds } Y=1 | X=1 \\ \hline \text{odds } Y=1 | X=2 \end{array}$$

$$\ll \frac{X=1}{X=3}$$

$$\ll \frac{X=2}{X=3}$$

# Ch. 3 : Generalized Linear Models

Plan:

- ① review linear model
- ② What assumptions of linear models are not met w/ categorical response.
- ③ "Fix up" linear model to satisfy assumptions.

Linear Model :  $Y$  = response variable  
 $X_1, X_2, \dots, X_k$  = predictor variables

Model  $Y \mid X_1 = x_1, \dots, X_k = x_k$ .

$i=1, \dots, n$

$$Y_i = \underbrace{\beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki}}_{\text{fixed}} + \underbrace{\varepsilon_i}_{\text{random}}$$

Assumptions on  $\varepsilon$  :

- ① Normality
- ②  $E(\varepsilon_i) = 0$
- \* ③  $\text{Var}(\varepsilon_i) = \sigma^2 \rightarrow$  constant
- \* ④ Independence

Mathematically:

$$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

(independent  
identically  
distributed)

Alternate expression:  $Y_i \stackrel{\text{indep.}}{\sim} N(\mu_i, \sigma^2)$

$$\mu_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki}$$

Generalized linear models - response variable has some other distribution besides normal.

e.g. Binomial - Linear model assumptions violated  $\rightarrow$   
① + ③

Binomial:  $Y_i \stackrel{\text{indep}}{\sim} \text{Bin}(N_i, \pi_i) \quad i=1, \dots, n$

→ Model  $\pi_i$  using  $X_{1i} \rightarrow X_{ki}$

$$Y_i \neq \mu_i + \varepsilon_i \quad E(Y_i) = N_i \pi_i$$

$$N_i \pi_i \neq \varepsilon_i$$

GLM: 3 components:

① Systematic component (linear predictor):

$$\eta = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

② Link function:  $E(Y) = \mu \quad g(\mu) = \eta$

— how do we model  $E(Y)$  as a function of  $x_1, x_2, \dots, x_k$

③ Random component — probability distributional assumptions on  $Y$

Normal linear model:

$$Y \sim N(\mu, \sigma^2)$$

$$E(Y) = \mu$$

$$\mu = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

$$\text{Link: } g(\mu) = \mu$$

(identity)

Generalized LM:

$$Y \sim ??$$

$$E(Y) = \mu$$

$$g(\mu) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

Examples:  $\odot$  Bernoulli:

$$Y \sim \text{Bin}(1, \pi)$$

$$E(Y) = \pi$$

$$g(\pi) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

Logistic regression

$$\longrightarrow g(\pi) = \log\left(\frac{\pi}{1-\pi}\right)$$

$$\Rightarrow \log\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

Common error:  $\log\left(\frac{Y}{1-Y}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$

$Y \rightarrow$  random

fixed

$$\log\left(\frac{Y}{1-Y}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$$

\* Don't use an error term!

② Poisson:  $Y \sim \text{Pois}(\mu)$

→ counts

$Y = 0, 1, 2, 3, \dots$

Most common: "log-linear"

$$\frac{\log \mu}{g(\mu)} = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

Other common link functions:

① Identity:  $g(\mu) = \mu$

Normal

② Logit:  $g(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$

Binomial

③ Log:  $g(\mu) = \log(\mu)$

Poisson

④ Inverse:  $g(\mu) = \frac{1}{\mu}$

Gamma

⑤ Probit:  $g(\mu) = \Phi^{-1}(\mu)$

Binomial

Inverse cdf of  
Std. normal

$$\Phi(z) = P(Z \leq z)$$

⑥ Complementary log-log link:  $g(\mu) = \log(\log(1-\mu))$   
(Binomial)