# I. Probability

Population of interest = US adults in labor force in 1992.

Relationship of interest = Risk of unemployment and Education attainment

Let, Y be the random variable representing the unemployment status

$$Y = \begin{cases} 1, & \text{if unemployed} \\ 0, & \text{if employed} \end{cases}$$

Let X be the random variable representing education attainment.

$$X = \begin{cases} 1, & \text{if no degree} \\ 2, & \text{if GED} \\ 3, & \text{if High School degree} \\ 4, & \text{if Associates degree} \\ 5, & \text{if Bachelors degree} \\ 6, & \text{if Masters degree or higher} \end{cases}$$

# 1. Constructing a theoretical contingency table.

| | X = 1 | X = 2 | X = 3 | X = 4 | X = 5 | X = 6 | |
|---|---|---|---|---|---|---|---|
| Y = 1 | $\pi_{11}$ | $\pi_{12}$ | $\pi_{13}$ | $\pi_{14}$ | $\pi_{15}$ | $\pi_{16}$ | $\pi_{1.}$ |
| Y = 0 | $\pi_{01}$ | $\pi_{02}$ | $\pi_{03}$ | $\pi_{04}$ | $\pi_{05}$ | $\pi_{06}$ | $\pi_{0.}$ |
| | $\pi_{.1}$ | $\pi_{.2}$ | $\pi_{.3}$ | $\pi_{.4}$ | $\pi_{.5}$ | $\pi_{.6}$ | $\pi$ |

Let $i = \{1, 0\}$ be the indices for Y

$j = \{1, 2, 3, 4, 5, 6\}$ be the indices for X

$\pi_{ij}$ = Joint probability of the event when
$y = i$ and $x = j$

$\pi_{1.}$ = marginal probability of unemployment
= $P(Y = 1)$

$\pi_{0.}$ = marginal probability of employment
= $P(Y = 0)$

$\pi_{.1}$ = marginal probability of having no degree
= $P(X = 1)$

$\pi_{.2}$ = marginal probability of having a GED

$\quad$ = $P(x = 2)$

$\pi_{.3}$ = marginal probability of having a high School degree

$\quad$ = $P(x = 3)$

$\pi_{.4}$ = marginal probability of having an Associate's Degree

$\quad$ = $P(x = 4)$

$\pi_{.5}$ = marginal probability of having a Bachelor's Degree

$\quad$ = $P(x = 5)$

$\pi_{.6}$ = marginal probability of having a Master's degree or higher

$\quad$ = $P(x = 6)$

Based on the definitions

$$\pi_{1.} = \sum_{j=1}^{6} \pi_{1j}$$

$$\pi_{0.} = \sum_{j=1}^{6} \pi_{0j}$$

Similarly,

$$\Pi_{.1} = \sum_{i=0}^{1} \Pi_{i1}$$

$$\Pi_{.2} = \sum_{i=0}^{1} \Pi_{i2}$$

$$\Pi_{.3} = \sum_{i=0}^{1} \Pi_{i3}$$

$$\Pi_{.4} = \sum_{i=0}^{1} \Pi_{i4}$$

$$\Pi_{.5} = \sum_{i=0}^{1} \Pi_{i5}$$

$$\Pi_{.6} = \sum_{i=0}^{1} \Pi_{i6}$$

2. Marginal probability of unemployment

$$P(Y=1) = ?$$

$$P(Y=1) = \Pi_{1.} = \sum_{j=1}^{6} \Pi_{ij} = \Pi_{11} + \Pi_{12} + \Pi_{13} + \Pi_{14} + \Pi_{15} + \Pi_{16}$$

$$= P(Y=1, X=1) + P(Y=1, X=2) + P(Y=1, X=3) +$$

$$P(Y=1, X=4) + P(Y=1, X=5) + P(Y=1, X=6)$$

3. Conditional Probability of unemployment for each possible level of education

(i) $X = 1$

$$P(Y=1 \mid X=1) = \frac{P(Y=1, X=1)}{P(X=1)} = \frac{\pi_{11}}{\pi_{\cdot 1}}$$

This relationship is true due to the property,

$$P(Y=y \mid X=x) = \frac{P(Y=y, X=x)}{P(X=x)}$$

Similarly,

(ii) $X = 2$

$$P(Y=1 \mid X=2) = \frac{P(Y=1, X=2)}{P(X=2)} = \frac{\pi_{12}}{\pi_{\cdot 2}}$$

(iii) $X = 3$

$$P(Y=1 \mid X=3) = \frac{P(Y=1, X=3)}{P(X=3)} = \frac{\pi_{13}}{\pi_{\cdot 3}}$$

(iv) $X = 4$

$$P(Y=1 \mid X=4) = \frac{P(Y=1, X=4)}{P(X=4)} = \frac{\pi_{14}}{\pi_{\cdot 4}}$$

(v) $X = 5$

$$P(Y=1 \mid X=5) = \frac{P(Y=1, X=5)}{P(X=5)} = \frac{\pi_{15}}{\pi_{\cdot 5}}$$

(vi)  X = 6

$$P(Y=1 \mid X=6) = \frac{P(Y=1, X=6)}{P(X=6)} = \frac{\pi_{16}}{\pi_{.6}}$$

Substantive Interpretation.

$P(Y=1 \mid X=x) =$ Probability of being unemployed if the education level is at $X = x$.


4.  If unemployed, $Y=1$

   if no degree, $X=1$

   $P(Y=1, X=1) =$ Joint probability of being unemployed and having no degree

   $$P(Y=1, X=1) = \underbrace{P(Y=1 \mid X=1)}_{\substack{\text{conditional} \\ \text{probability} \\ \text{component}}} \cdot \underbrace{P(X=1)}_{\substack{\text{marginal} \\ \text{probability} \\ \text{component}}}$$

5. Marginal probability of being unemployed **/2**

$$P(Y=1) = \sum_{x=1}^{6} P(Y=1|X=x) P(X=x)$$

USING DATA FROM NALS

**/4**

6.

$P(Y=1|X=1) = 0.16$

$P(Y=1|X=2) = 0.13$

$P(Y=1|X=3) = 0.88$

$P(Y=1|X=4) = 0.07$

$P(Y=1|X=5) = 0.05$

$P(Y=1|X=6) = 0.03$

We observe that the risk of unemployment decreases as education level increases.

7. Let's assume that education and unemployment are independent,

$$X \perp\!\!\!\perp Y$$

$$\therefore P(Y=1 \mid X) = P(Y=1) = 0.084$$

$\therefore 8.4 \%$ of individuals in each level of X are expected to be unemployed.

In $X=1$, we have 1579 individuals

$\therefore$ Expected number of unemployed

$$= 0.084 * 1579$$
$$\approx 126$$

But the number of those with no degree and unemployed from the data

Count for $(X=1, Y=1) = 253$

The actual value is greater than that estimated under the assumption of x and y being independent.

Thus, the assumption may not hold true.

# II. EXPECTATION

Variables of interest:

$X$ = Parent years of education

$Y$ = Adult literacy (outcome of interest)

$Z$ = Respondent years of education

1. Theoretical Linear model for $Y = f(X, Z)$

True model:

$$y_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i + \varepsilon_i$$

$i = \{1, \dots, n\}$ = indice for sample subjects

$y_i$ = adult literacy of $i^{th}$ individual in the sample

$\beta_0$ = intercept (no substantive meaning)

$X_i$ = parent years of education for $i^{th}$ individual in the sample

$\beta_1$ = For one unit change in parent years of education, holding 'Z' constant, $\beta_1$ is the change in adult literacy outcome

$z_i$ = years of education of individual (respondent) "i" in the sample

$\beta_2$ = For one unit change in Z (respondent years of education), holding X constant, $\beta_2$ is the change in adult literacy outcome.

$\varepsilon_i$ = error term

= random deviation of $i^{th}$ individual's outcome from that predicted by the true model.

$$\varepsilon_i = y_i - (\beta_0 + \beta_1 x_i + \beta_2 z_i)$$

2. (a) $Z_i = \gamma_0 + \gamma_1 x_i + e_i$

(b) $X_i = \delta_0 + \delta_1 z_i + u_i$

3. (a) $E(Y|Z) = ?$ if we only use 'z' as a predictor.

$$Y_i = \beta_0 + \beta_2 z_i + \beta_1 x_i + \varepsilon_i$$

$$E(Y|Z) = \beta_0 + \beta_2 z + \beta_1 E(X|Z) + E(\varepsilon|Z)$$

(b) Bias due to omission of $X$.

If we estimate the true model,

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i$$

then from OLS estimates, $\{\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2\}$

$$E(\hat{\beta}_0) = \beta_0 \; ; \; E(\hat{\beta}_1) = \beta_1 \; ; \; E(\hat{\beta}_2) = \beta_2$$

But if we instead believe that,

$$y_i = \beta_0^* + \beta_2^* z_i + \varepsilon_i^*$$

From what we studied in class for a single predictor linear model,

$$\hat{\beta}_2^* = \frac{\widehat{cov}(Y, z)}{\widehat{var}(z)} = \frac{\widehat{cov}(\beta_0 + \beta_1 x + \beta_2 z + \varepsilon, z)}{\widehat{var}(z)}$$

$$= \frac{\overbrace{\widehat{cov}(z, \beta_0)}^{0} + \widehat{cov}(z, \beta_1 x) + cov(z, \beta_2 z) + \overbrace{cov(z, \varepsilon)}^{0}}{\widehat{var}(z)}$$

$$= \frac{\beta_1 \widehat{cov}(z, x) + \beta_2 \widehat{var}(z)}{\widehat{var}(z)}$$

$$= \beta_2 + \beta_1 \cdot \frac{\widehat{cov}(z, x)}{\widehat{var}(z)}$$

$$\text{Bias} = \beta_1 \cdot \frac{\hat{\text{cov}}(z,x)}{\hat{\text{var}}(z)}$$

Bias = 0    if,

(i)  $\beta_1 = 0$,  x does not belong in the true model, so its omission has no effect on X.

(ii)  $\hat{\text{cov}}(z,x) = 0$, if z and x are not correlated then omitting one doesn't affect the other.

4. $y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i$ — ①

$x_i = \delta_0 + \delta_1 z_i + u_i$ — ②

$y_i = \theta_0 + \theta_1 x_i + v_i$ — ③

$z_i = \gamma_0 + \gamma_1 x_i + e_i$ — ④

a) $\theta_1 = $ Total effect of $x$ on $y$

$$E(y|x) = \beta_0 + \beta_1 x + \beta_2 E(z|x) + \varepsilon$$

$$\approx \theta_0 + \theta_1 x + u_i$$

S.t. $u_i = \beta_2 E(z|x) + \varepsilon_i$

b) Based on ①

Direct effect of $x$ on $y = \beta_1$

c) $y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i$

$$= \beta_0 + \beta_1 x_i + \beta_2 (\gamma_0 + \gamma_1 x_i) + \varepsilon_i + \beta_2 e_i$$

$$= \beta_0 + (\beta_1 + \beta_2 \gamma_1) x_i + \varepsilon_i + \beta_2 e_i$$

earlier we saw that
Direct effect of $x$ on $y = \beta_1$
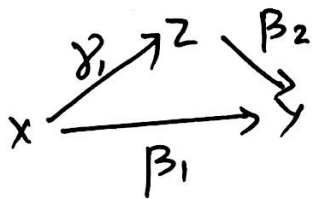
Indirect effect $= \beta_2 \gamma_1$

(d) if
$$y_i = \beta_0 + (\beta_1 + \beta_2 \gamma_1) x_i + \varepsilon_i + \beta_2 e_i$$

and

$$y_i = \theta_0 + \theta_1 x_i + v_i$$

we see that,

$$\theta_1 = \underbrace{\beta_1}_{\substack{\text{Direct} \\ \text{effect}}} + \underbrace{\beta_2 \gamma_1}_{\substack{\text{Indirect} \\ \text{effect}}}$$



5. Using OLS in R.

Direct effect $= \hat{\beta}_1 = 3.99$

Indirect effect $= \hat{\beta}_2 \cdot \hat{\gamma}_1 = 4.75$

Total effect $= \hat{\theta}_1 = \hat{\beta}_1 + \hat{\beta}_2 \hat{\gamma}_1$

$= 3.99 + 4.75$

$= 8.74$