

Biased and Unbiased Samples

James J. Heckman
Econ 312, Spring 2019

May 13, 2019

Definitions and Some Examples of Biased Samples

- All sampling models can be described by the following set-up.
- Let \mathbf{Y} be a vector of outcomes of interest and let \mathbf{X} be a vector of “control” or “explanatory” variables.
- The population distribution of (\mathbf{Y}, \mathbf{X}) is $F(y, x)$.
- Assume that the density is well defined and write it as $f(y, x)$.

- Any **sampling rule** can be interpreted as producing a non-negative weighting function of $\omega(\mathbf{y}, \mathbf{x})$ that alters the population density.
- Let $(\mathbf{Y}^*, \mathbf{X}^*)$ denote the sampled random variables.
- The density of the sampled data $g(\mathbf{y}^*, \mathbf{x}^*)$ may be written as

$$g(\mathbf{y}^*, \mathbf{x}^*) = \frac{\omega(\mathbf{y}^*, \mathbf{x}^*) f(\mathbf{y}^*, \mathbf{x}^*)}{\int \omega(\mathbf{y}^*, \mathbf{x}^*) f(\mathbf{y}^*, \mathbf{x}^*) d\mathbf{y}^* d\mathbf{x}^*} \quad (1)$$

- The denominator of the expression introduced to make the density $g(\mathbf{y}^*, \mathbf{x}^*)$ integrate to one.

- Alternatively, the weight may be defined as

$$\omega^*(y^*x^*) = \frac{\omega(y^*, x^*)}{\int \omega(y^*, x^*) f(y^*, x^*) dy^* dx^*}$$

so that

$$g(y^*, x^*) = \omega^*(y^*, x^*)f(y^*, x^*). \quad (2)$$

-

$$\Pr(\Delta = 0) = 1 - \Pr(\Delta = 1).$$

- A *truncated sample* is one for which $\Pr(\Delta = 1)$ is not known and cannot be identified.
- A *censored sample* is one for which $\Pr(\Delta = 1)$ is known or can be identified.
- Sampling rule in this case is such that frequency of \mathbf{y}, \mathbf{x} for which $\omega(\mathbf{y}, \mathbf{x}) = 0$ are not known.
- It is known whether or not $i(\mathbf{y}, \mathbf{x}) = 0$ for all values of \mathbf{Y}, \mathbf{X} .

Two Cases

- A *truncated sample* is one for which $\Pr(\Delta = 1)$ is not known and cannot be identified.
- A *censored sample* is one for which $\Pr(\Delta = 1)$ is known or can be identified.
- Sampling rule in this case is such that frequency of \mathbf{y}, \mathbf{x} for which $\omega(\mathbf{y}, \mathbf{x}) = 0$ are not known.
- It is known whether or not $i(\mathbf{y}, \mathbf{x}) = 0$ for all values of \mathbf{Y}, \mathbf{X} .

- Thus

$$g(\mathbf{y}^*, \mathbf{x}^*, \delta) = \left[\frac{\omega(\mathbf{y}^*, \mathbf{x}^*) f(\mathbf{y}^*, \mathbf{x}^*)}{\int \omega(\mathbf{y}^*, \mathbf{x}^*) f(\mathbf{y}^*, \mathbf{x}^*) d\mathbf{y}^* d\mathbf{x}^*} \right]^\delta \quad (4)$$

$$\times \left[\int i(\mathbf{y}, \mathbf{x}) f(\mathbf{y}, \mathbf{x}) d\mathbf{y} d\mathbf{x} \right]^\delta$$

$$\times [1]^{1-\delta} \left[\int (1 - i(\mathbf{y}, \mathbf{x})) f(\mathbf{y}, \mathbf{x}) d\mathbf{y} d\mathbf{x} \right]^{1-\delta}.$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Selection on the exogenous variables:

$$g(y^*, x^*) = f(y^*|x^*) \frac{\omega(x^*)f(x^*)}{\int \omega(x^*)f(x^*)dx}$$

and

$$g(x^*) = \frac{\omega(x^*)f(x^*)}{\int \omega(x^*)f(x^*)dx^*}.$$

- $$g(y^*|x^*) = \frac{g(y^*, x^*)}{g(x^*)} = f(y^*|x^*).$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

General Stratified Sampling

- Sampling on both \mathbf{y} and \mathbf{x} .

-

(5)

- Integrating the left-hand side of (5) it is possible to determine $\int \omega(\mathbf{y}^*, \mathbf{x}^*) f(\mathbf{y}^*, \mathbf{x}^*) d\mathbf{y}^* d\mathbf{x}^*$.
- Hence can use (5) to recover the population density of the data.

$$\int f(y^*, x^*) dy^* dx^* = 1.$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Distribution of Y^* is

$$G(y^*|Y > c) = F(y^*|Y > c) = F(y^*|\Delta = 1) \quad (6a)$$

$$= \frac{F(y^*)}{1 - F(c)}, y^* > c.$$

Point mass at $Y^* = 0$ (Convention) for $Y^* = 0$ ($\Delta = 0$). (6b)

- Observe that (6a) is obtained from (1) by setting $\omega(y^*) = 1$ if $y > c$, and $\omega(y^*) = 0$ otherwise, and integrating up with respect to y^* .
- The distribution of Δ is

$$\Pr(\Delta = \delta) = [1 - F(c)]^\delta [F(c)]^{1-\delta}, \delta \in \{0, 1\}.$$

- The joint distribution of (Y^*, Δ) for a censored sample:

$$\begin{aligned} F(y^*, \delta) &= F(y^*|\delta)\Pr(\delta) \\ &= \left\{ \frac{F(y^*)}{(1 - F(c))} \right\}^\delta [1 - F(c)]^\delta (1)^{1-\delta} [F(c)]^{1-\delta} \\ &= [F(y^*)]^\delta [F(c)]^{1-\delta}. \end{aligned} \tag{7}$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- More information in a censored sample than in a truncated sample because one can obtain (6a) from (7) (by conditioning on $\Delta = 1$) but not vice versa.

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

- Let

$$Y = \mathbf{X}\boldsymbol{\beta} + U \quad (8)$$

be the population earnings function where Y is earnings.

- “ β ”: suitably dimensioned parameter vector.
- X is a regressor vector assumed to be distributed independently of mean zero disturbance U .
- $U \perp\!\!\!\perp X$; $E(XX')$ full rank, $E(U) = 0$.

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Invoke independence between U and \mathbf{X} and letting F_u denote the distribution of U ,

$$\Pr(\Delta = 1 | \mathbf{X} = \mathbf{x}) = 1 - F_u(c - \mathbf{x}\beta) \quad (9a)$$

and

$$\Pr(\Delta = 0 | \mathbf{X} = \mathbf{x}) = F_u(c - \mathbf{x}\boldsymbol{\beta}). \quad (9b)$$

$$\begin{aligned} G(y^* | Y > \mathbf{0}, \mathbf{X} = x) &= F(y^* | X = x, Y > c) \\ &= F(y^* | \mathbf{X} = x, \Delta = 1) \\ &= \frac{F_u(y^* - x\beta)}{1 - F_u(c - x\beta)}, \quad y^* > c. \end{aligned} \quad (10a)$$

$$G(y^* | Y \leq 0) = 1 \text{ for } Y^* = 0 \ (\Delta = 0). \quad (10b)$$

- In particular,

$$\begin{aligned} F(y^*, \delta | \mathbf{X} = \mathbf{x}) &= F(y^* | \delta, \mathbf{x}) \Pr(\delta | \mathbf{x}) \\ &= \{F_u(y^* - \mathbf{x}\beta)\}^\delta \{F_u(c - \mathbf{x}\beta)\}^{1-\delta}. \end{aligned} \quad (11)$$

$$\begin{aligned} E(Y^* \mid \mathbf{X} = \mathbf{x}, \Delta = 1) &= \mathbf{x}\beta + E(U \mid \mathbf{X} = \mathbf{x}, \delta = 1) \quad (12) \\ &= \mathbf{x}\beta + \int_{c - \mathbf{x}\beta}^{\infty} \frac{z \, dF_u(z)}{(1 - F_u(c - \mathbf{x}\beta))} \end{aligned}$$

- z : dummy variable of integration.

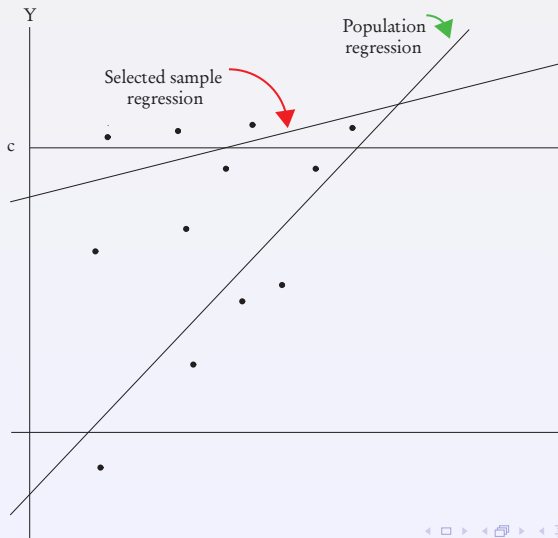
- Contrast between (12) and (13) illuminating.
- When theoretical model is estimated on a selected sample ($\Delta = 1$), the true conditional expectation is (12) not (13).

$$E(Y \mid \mathbf{X} = x) = x\beta. \quad (13)$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Illustrate the nature of the bias, it is useful to draw on the work of Cain and Watts (1973).
- Suppose that X is a scalar random variable (e.g., education) and that its associated coefficient is positive ($\beta > 0$).
- Under conventional assumptions about U (e.g., mean zero, independently and identically distributed and distributed independently of X), the population regression of Y on X is a straight line.
- The scatter about the regression line and the regression line are given in Figure 1.

Figure 1:



- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Fruitful to distinguish between the case of a truncated sample and the case of a censored sample.
- In the truncated sample case, no information is available about the fraction of the population that would be allocated to the truncated sample $[\Pr(\Delta = 1)]$.
- In the censored sample case, this fraction is known or can be consistently estimated.
- Fruitful to distinguish two further cases:
- Case (a), the case in which \mathbf{X} is not observed when $\Delta = 0$.
- Case (b) is the one most fully developed in the literature: \mathbf{X} observed when $D = 0$.

- Conditional mean $E(U \mid \mathbf{X} = \mathbf{x}, \Delta = 1)$ is a function of $c - \mathbf{x}\beta$ solely through $\Pr(\Delta = 1 \mid \mathbf{x})$.
- Since $\Pr(\Delta = 1 \mid \mathbf{x})$ is monotonic in $c - \mathbf{x}\beta$.
- The conditional mean depends solely on $\Pr(\Delta = 1 \mid \mathbf{x})$ and the parameters F_u i.e., since

$$F_u^{-1}(1 - \Pr(\Delta = 1 \mid \mathbf{x})) = c - \mathbf{x}\beta$$

$$\begin{aligned} E(U \mid X = x, \Delta = 1) &= \int_{F_u^{-1}[1 - \Pr(\Delta=1|x)]}^{\infty} \frac{z dF_u(z)}{\Pr(\Delta = 1 \mid x)} \\ &= K(P(\Delta = 1|x)) \\ \ln P(\Delta = 1|x) \rightarrow 1, K(P(\Delta = 1|x)) &= 0. \end{aligned}$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- If $F(y_1, y_2)$ is the population distribution of (Y_1, Y_2) , the distribution of Δ is

$$\Pr(\Delta = \delta) = [1 - F_2(c)]^{1-\delta} [F_2(c)]^\delta, \quad \delta = 0, 1,$$

- F_2 is the marginal distribution of Y_2 .

- The distribution of Y_1^* is

$$G(y_1^*) = F(y_1^*; \delta = 1) = \frac{F(y_1^*; c)}{F_2(c)}, \quad \Delta = 1, \quad (14a)$$

$$G(y_1^* = 0) = 1, \quad \Delta = 0. \quad (14b)$$

- (14a): the distribution function corresponding to the density in (1) when $\omega(y_1, y_2) = 1$ if $y_2 \leq c$ and $\omega(y_1, y_2) = 0$ otherwise.

- This is the distribution function corresponding to density (4) for the special weighting rule of this example.
- In a censored sample, under general conditions it is possible to consistently estimate $\Pr(\Delta = \delta)$ and $G(y_1^*)$.

$$G(y_1^*, \delta) = [F(y_1^*; c)]^\delta [1 - F_2(c)]^{1-\delta}. \quad (15)$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- $Y_1^* = Y_1$ if $Y_1 - Y_2 < 0$ while $Y_2^* = Y_2$ if $Y_1 - Y_2 \geq 0$.

- In example 3 set

$$Y_1 = \mathbf{X}_1\boldsymbol{\beta}_1 + U_1 \quad (16a)$$

$$Y_2 = \mathbf{X}_2\boldsymbol{\beta}_2 + U_2 \quad (16b)$$

where $(\mathbf{X}_1, \mathbf{X}_2)$ are distributed independently of (U_1, U_2) , a mean zero, finite variance random vector.

- Conventional assumptions are invoked to ensure that if Y_1 and Y_2 can be observed, least squares applied to a random sample of data on $(Y_1, Y_2, \mathbf{X}_1, \mathbf{X}_2)$ would consistently estimate β_1 and β_2 .
- $Y_1^* = Y_1$ if $Y_2 < 0$.
- If $Y_2 < 0, \Delta = 1$.
- Regression function for the selected sample is

$$E(Y_1^* | \mathbf{X}_1 = \mathbf{x}_1, Y_2 < 0) = E(Y_1^* | \mathbf{X}_1 = \mathbf{x}_1, \Delta = 1) = \mathbf{x}_1\beta_1 + E(U_1 | \mathbf{X}_1 = \mathbf{x}_1, \Delta = 1) \quad (17)$$

- Regression function for the population is

$$E(Y_1 | \mathbf{X}_1 = \mathbf{x}_1) = \mathbf{x}_1\beta_1. \quad (18)$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- *Example 5. Length biased sampling.*
- Let T be the duration of an event such as a completed unemployment spell or a completed duration of a job with an employer.
- The population distribution of T is $F(t)$ with density $f(t)$.
- The sampling rule is such that *individuals* are sampled at random.
- Data are recorded on a completed spell *provided that at the time of the interview the individual is experiencing the event.*
- Such sampling rules are in wide use in many national surveys of employment and unemployment.

- This is the **hazard rate**.

- This is the **hazard rate**.

- Suppose that the environment is stationary.
- The population entry rate into the state at each instant of time is k .
- From each vintage of entrants into the state distinguished by their distance from the survey date t_b , only $1 - F(t_b) = \Pr(T > t_b)$ survive.
- Aggregating over all cohorts of entrants, the population proportion in the state at the date of the interview is P where

$$P = \int_0^{\infty} k(1 - F(t_b))dt_b \quad (20)$$

which is assumed to exist.

- In a duration of unemployment example, P is the unemployment rate.

- The density of sampled **completed** durations is thus

$$g(t_b^* | t_b^* > 0) = \frac{k(1 - F(t_b^*))}{\rho}. \quad (21)$$

$$\begin{aligned} g(t^*) &= \int_0^{t^*} f(t^*|t_b^*)g(t_b^*|t_b^* > 0)dt_b^* \\ &= k \frac{f(t^*)}{1 - F(t_b^*)} \frac{1 - F(t_b^*)}{P} \int_0^{t^*} dt_b^* \\ &= k \frac{t^* f(t^*)}{P}. \end{aligned}$$

- **Length biased sampling.**

- Integration by parts:

$$P = k \int_0^{\infty} (1 - F(z)) dz = k \int_0^{\infty} z dF(z) = kE(T).$$

- Note that

$$g(t^*) = \frac{t^* f(t^*)}{E(T)}. \quad (22)$$

- We know that $g(t^*)$.
- Can form $\frac{g(t^*)}{t^*}$, $t^* > 0$.
- \therefore we know $\frac{f(t^*)}{E(T)}$.

- Apply analysis of (5): $\int_0^{\infty} \frac{g(t^*)}{t^*} dt^* = \frac{\overbrace{\int_0^{\infty} f(t^*) dt^*}^{=1}}{E(T)}.$

- \therefore know $f(t^*)$.

- In this form (22) is equivalent to (1) with $\omega(t) = t$.
- $E(T)$.
- **Length biased sampling.**
- Intuitively, longer spells are oversampled when the requirement is imposed that a spell be in progress at the time the survey is conducted ($T_b > 0$).
- Suppose, instead, that individuals are randomly sampled and data are recorded on the **next** spell of the event (after the survey date).

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

- $$f(t_c|t_b) = \frac{f(t_c)}{1 - F(t_b)}$$

- $$\begin{aligned} g(t_c) &= \int_0^{t_c} f(t_c|t_b)g(t_b)dt_b \\ &= \int_0^{t_c} \frac{f(t_c)}{m}dt_b = \frac{f(t_c)t_c}{m} \end{aligned}$$

- $$\begin{aligned} g(t_a) &= \int_0^\infty f(t_a + t_b | t_b) g(t_b) dt_b \\ &= \int_0^\infty \frac{f(t_a + t_b)}{m} dt_b \\ &= \frac{1}{m} \int_{t_a}^\infty f(t_b) dt_b \\ &= \frac{1 - F(t_a)}{m} \end{aligned}$$

- 1 If $f(t) = \theta e^{-t\theta}$, then $g(t_b) = \theta e^{-t_b\theta}$ and $g(t_a) = \theta e^{-t_a\theta}$.
- 2 **Proof:**

$$f(t) = \theta e^{-t\theta} \rightarrow m = \frac{1}{\theta},$$

$$F(t) = 1 - e^{-t\theta} \rightarrow g(t_a) = \frac{1 - F(t)}{m} = \theta e^{-t\theta}$$

$$\textcircled{1} \quad E(T_a) = \frac{m}{2} \left(1 + \frac{\sigma^2}{m^2} \right).$$

Proof:

$$\begin{aligned} E(T_a) &= \int t_a f(t_a) dt_a = \int t_a \frac{1 - F(t_a)}{m} dt_a \\ &= \frac{1}{m} \left[\frac{1}{2} t_a^2 (1 - F(t_a)) \Big|_0^\infty - \int \frac{1}{2} t_a^2 d(1 - F(t_a)) \right] \\ &= \frac{1}{m} \int \frac{1}{2} t_a^2 f(t_a) dt_a = \frac{1}{2m} [\text{var}(t_a) + E^2(t_a)] \\ &= \frac{1}{2m} [\sigma^2 + m^2] \end{aligned}$$

- $$\textcircled{1} \quad E(T_b) = \frac{m}{2} \left(1 + \frac{\sigma^2}{m^2} \right).$$

2 **Proof:** See proof of Proposition 2.

$$\textcircled{3} \quad E(T_c) = m(1 + \frac{\sigma^2}{m^2}).$$

4 Proof:

$$E(T_c) = \int \frac{t_c^2 f(t_c)}{m} dt_c = \frac{1}{m}(\text{var}(t_c) + E^2(t_c))$$

$$\rightarrow E(T_c) = 2E(T_a) = 2E(T_b), E(T_c) > m \text{ unless } \sigma^2 = 0$$

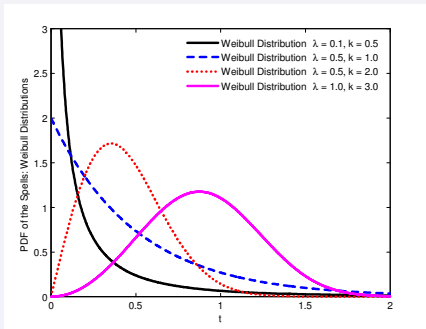
Examples

1

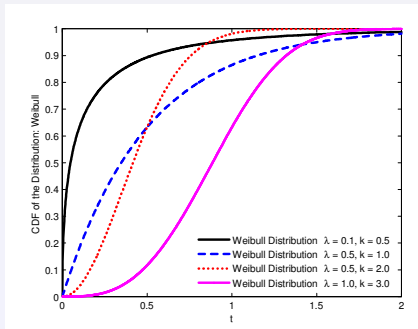
Figure 1

Basic Distribution Graphs

PDF for Weibull Distribution



CDF of Weibull Distribution



1000



[illegible]

1. *Journal of Management Studies*, 1997, 34(1), 1-15.



1. *Journal of Management Studies*, 1996, 33, 1, 1-14.



1. *Journal of Management Studies*, 1997, 34, 1, 1-14.



— *Journal of the American Medical Association*, 1997



1. *Journal of the American Medical Association*, 2000; 284: 2689-2695.



- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- The population density of (D, \mathbf{X}) is

$$f(d, \mathbf{x}) = \Pr(D = d | \mathbf{X} = \mathbf{x})h(\mathbf{x}) \quad (24)$$

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

1. *Journal of the American Medical Association*, 1997; 277: 1039-1043.

- Notice that the dominator can be simplified to

$$\sum_{i=1}^I \omega(i) f(i)$$

- $f(d^*)$ is the marginal distribution of D^* so that

$$g(d^*, \mathbf{x}^*) = \frac{\omega(d^*) f(d^*, \mathbf{x}^*)}{\sum_{i=1}^I \omega(i) f(i)}. \quad (26)$$

$$g(d^*) = \frac{\omega(d^*)f(d^*)}{\sum_{i=1}^I \omega(i)f(i)} \quad (27)$$

- Sampling rule causes the sampled proportions to deviate from the population proportions.

- $$h(\mathbf{x}^*|d^*) = \frac{f(d^*, \mathbf{x}^*)}{f(d^*)} \quad (28)$$

- The density of x in the sample is thus

$$g(x^*) = \sum_{i=1}^I h(x^* | i) g(i). \quad (29)$$

- Then using (26)-(29) we reach

$$g(d^*|x^*) = f(d^*|x^*) \times \left\{ \left[\frac{\omega(d^*)}{\sum_{i=1}^I \omega(i)f(i)} \right] \left[\frac{1}{\sum_{i=1}^I f(i|x^*) \frac{g(i)}{f(i)}} \right] \right\}. \quad (30)$$

- The bias that results from using choice based samples to make inference about $f(d^*|x^*)$ is a consequence of neglecting the terms in braces on the right-hand side of (30).

- Notice that if the data are generated by a random sampling rule, $\omega(d^*) = 1$, $g(d^*) = f(d^*)$ and the term in braces is one.

Further Discussion of Choice Based Samples

- Pick D first (e.g. travel mode).
- Probability of selecting D is $C(D)$.
- $f(D, X)$ is the joint density of D and X in the population.

$$f(D, X | \theta) = g(D | X, \theta)h(X) = \varphi(X | D)f(D | \theta)$$

$$f(D | \theta) = \int g(D | X, \theta)h(X)dX$$

- Given D we observe X (the implicit assumption is that we are sampling only on D , not on D and X).
- Probability of *sampled* X, D is $\varphi(X | D)C(D)$.

- A fact we use later is

$$\begin{aligned}\varphi(X | D)C(D) &= \left\{ \frac{g(D | X)h(X)}{f(D)} \right\} C(D) \\ &= \frac{g(D | X)h(X)C(D)}{\left[\int g(D | X)h(X)dX \right]}.\end{aligned}$$

- When $C(D) = f(D) = \int g(D | X)h(X)dX$, choice based sampling is random sampling.

- Note, the likelihood function in an exogenous sampling scheme is

$$\begin{aligned}\mathcal{L} &= \prod_{i=1}^I f(D_i, X_i) = \prod_{i=1}^I f(D_i \mid X_i, \theta) h(X_i) \\ \ln \mathcal{L} &= \sum_{i=1}^I \ln f(D_i \mid X_i) + \sum \ln h(X_i).\end{aligned}$$

- By exogeneity, we get the lack of dependence of distribution of X on θ .

- Likelihood function for a choice-based sampling scheme is

$$\ln \mathcal{L} = \sum_{i=1}^I [\ln g(D_i | X_i) + \ln h(X_i) - \ln f(D_i) + \ln C(D_i)] .$$

- Suppose $f(D)$ depends on parameters θ . \therefore Max with θ .

$$\frac{\partial \ln \mathcal{L}}{\partial \theta} = \sum_{i=1}^I \frac{\partial \ln g(D_i | X_i)}{\partial \theta} - \underbrace{\sum_{i=1}^I \frac{\partial \ln f(D_i)}{\partial \theta}}_{\text{source of bias}} .$$

- We neglect the second term in forming the usual estimators using only the first term.
- That is the source of the inconsistency.

Further Analysis of Choice Based Samples:

- An example in discrete choice.
- (c) Draw d by $\varphi(d)$.
- (d) Draw X by $f(X \mid d = 1)$.
- Joint density of data:

$$= \varphi(d = 1) \left[\frac{\Pr(d = 1 \mid X, \theta) f(X)}{\Pr(d = 1 \mid \theta)} \right]$$

- Now in a choice-based sample

$$\Pr^*(d = 1 \mid X) = \frac{f(X \mid d = 1, \theta)\varphi(d = 1)}{h^*(X)}$$

where $g^*(X)$ is the sampled X data.

- Joint density of *data* X is given by:

$$h^*(X) = f(X \mid d = 1, \theta)\varphi(d = 1) + f(X \mid d = 0, \theta)\varphi(d = 1)$$

and

$$\Pr(D = 1 \mid X) = \frac{f(X \mid d = 1) \Pr(d = 1)}{f(X)}$$

- $\Pr^*(D = 1 \mid X) =$

$$\frac{\frac{\Pr(D = 1 | X, \theta)f(X)}{\Pr(D = 1 | \theta)}\varphi(D = 1)}{\frac{\Pr(D = 1 | X, \theta)f(X)}{\Pr(D = 1 | \theta)}\varphi(D = 1) + \frac{\Pr(D = 0 | X, \theta)f(X)}{\Pr(D = 0 | \theta)}\varphi(D = 0)}$$

$$= \frac{\Pr(D = 1 | X, \theta)\varphi(D = 1) / \Pr(D = 1 | \theta)}{\Pr(D = 1 | X, \theta)\frac{\varphi(D = 1)}{\Pr(D = 1 | \theta)} + \Pr(D = 0 | X, \theta)\frac{\varphi(D = 0)}{\Pr(D = 0 | \theta)}}.$$

- Now we missample the population with density $f(X \mid D = 1)$ in a choice based sample:

$$\begin{aligned}
 \Pr^*(D = 1 \mid X) &= \frac{f(X \mid D = 1, \theta) \varphi(D = 1)}{f(X \mid D = 1, \theta) \varphi(D = 1) + f(X \mid D = 0, \theta) \varphi(D = 0)} \\
 &= \frac{\frac{f(X) \Pr(D = 1 \mid X)}{\Pr(D = 1)} \varphi(D = 1)}{\frac{f(X) \Pr(D = 1 \mid X)}{\Pr(D = 1)} \varphi(D = 1) + \frac{f(X) \Pr(D = 0 \mid X)}{\Pr(D = 0)} \varphi(D = 0)} \\
 &= \frac{\Pr(D = 1 \mid X)}{\Pr(D = 1 \mid X) + \Pr(D = 0 \mid X) \frac{\varphi(D = 0)}{\varphi(D = 1)} \cdot \frac{\Pr(D = 1)}{\Pr(D = 0)}} \\
 &= \frac{1}{1 + \left[\frac{\Pr(D = 0 \mid X)}{\Pr(D = 1 \mid X)} \right] \cdot \frac{\varphi(D = 0)}{\varphi(D = 1)} \cdot \frac{\Pr(D = 1)}{\Pr(D = 0)}}
 \end{aligned}$$

- With logit we get

$$\Pr^*(D = 1 \mid X) = \frac{1}{1 + e^{-(\alpha_0 + X\beta) + \ln \left[\frac{\varphi(D = 0) \Pr(D = 1)}{\varphi(D = 1) \Pr(D = 0)} \right]}}.$$

This goes into an intercept term:

$$= \frac{e^{\alpha^* + X\beta}}{1 + e^{\alpha^* + X\beta}}$$

$$\alpha^* = \alpha_0 - \ln \left[\frac{\varphi(D = 0)}{\varphi(D = 1)} \cdot \frac{\Pr(D = 1)}{\Pr(D = 0)} \right].$$

- How to solve problem: Reweight data by relative frequency in population.
- (Idea due to C.R. Rao, 1965, 1986.)
- Joint density of the data is

$$f(X \mid D = 1)\varphi(D = 1).$$

Use Bayes' rule to obtain

$$\frac{P(D = 1 \mid X)f(X)}{P(D = 1)}\varphi(D = 1).$$

- Now weight by

$$\frac{P(D = 1)}{\varphi(D = 1)}.$$

- Solution: Reweight the data to form the following weighted likelihood:

$$\frac{1}{N} \sum_{i=1}^N \left[\frac{\Pr(D_i = 1)}{\varphi(D_i = 1)} (D_i^*) \ln \Pr(D_i = 1 \mid X, \theta) + \frac{\Pr(D_i = 0)}{\varphi(D_i = 0)} (1 - D_i^*) \ln \Pr(D_i = 0 \mid X, \theta) \right]$$

$$P \int \{[\Pr(D = 1 | X, \theta_0)f(X | \theta_0)] \ln \Pr(D = 1 | X, \theta) +$$

$$\int [\Pr(D = 0 | X, \theta_0)f(X | \theta_0)] \ln \Pr(D = 0 | X, \theta)\} f(X | D)DX$$

- $$\frac{f(X \mid D = 1)\varphi(D = 1)}{g^*(X)} = \frac{\Pr(D = 1 \mid X)f(X)}{g^*(X)} \frac{\varphi(D = 1)}{\Pr(D = 1)}.$$

- $$f(X) = f(X \mid D = 1)\varphi(D = 1) \left[\frac{P(D = 1)}{\varphi(D = 1)} \right] + f(X \mid D = 0)\varphi(D = 0) \frac{\Pr(D = 0)}{\varphi(D = 0)}.$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- The probability of sampled family size of $N^* = n^*$ is

$$g(n^*) = \frac{\omega(n^*)f(n^*)}{E[\omega(N^*)]} \quad (31)$$

where $\omega(n^*) = 1 - (1 - \beta)^{n^*}$ (the probability that at least one child from a family of size n^* will be sampled).

- Note $(1 - \beta)$ = probability of sampling a child (assumed the same across all n^*).



$$E[\omega(N^*)] = \sum_{n^*} (1 - (1 - \beta)^{n^*}) f(n^*)$$

is the probability of observing a family.

- In a large population $\beta \rightarrow 0$ with increasing population size.

- Thus the limit form of (31) is identical to (22).
- Larger families tend to be oversampled and hence a misleading estimate of family size will be produced from such samples

$$\lim_{\beta \rightarrow 0} g(n^*) = \frac{n^* f(n^*)}{E(N^*)}. \quad (32)$$

- Since the model is formally equivalent to the length biased sampling model, all references and statements about identification given in Example 6 apply with full force to this example.
- See the discussion in Rao (1965).

Appendix

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

- Max $U(\mathbf{Z}, E)$ subject to $\mathbf{P}'\mathbf{Z} \leq M$.
- In the population \mathbf{P} and M are distributed independently of E .
- First order conditions for this problem are

$$\frac{\partial U(\mathbf{Z}, E)}{\partial \mathbf{Z}} \leq \lambda \mathbf{P} \quad (33)$$

where λ is the Lagrange multiplier associated with the budget constraint.

- Focusing on the demand for the first good, Z_1 , none of it is purchased if at zero consumption of Z_1

$$\frac{\partial U(\mathbf{Z}, E)}{\partial Z_1} \Big|_{Z_1=0} \leq \lambda P_1 \quad (34)$$

i.e., marginal valuation is less than marginal cost in utility terms.

- Conventional interior solution demand functions for Z_1 are defined for a given \mathbf{P} , M only for values of E such that

$$\frac{\partial U(\mathbf{Z}, E)}{\partial Z_1} \Big|_{Z_1=0} \geq \lambda P_1. \quad (35)$$

- $$\Pr(\Delta_1 = 0 \mid \mathbf{P}, M) = 1 - \int_E dF(\varepsilon).$$

- $$Z_1 = Z_1(\mathbf{P}, M, E) \quad (36)$$
- is well defined and $Z_1 = Z_1^*$.
- When $\Delta_1 = 0$, observed $Z_1 = Z_1^* = 0$.

$$Z_1 = Z_1(\mathbf{P}, M, E) \quad (36)$$

is well defined and $Z_1 = Z_1^*$.

- When $\Delta_1 = 0$, observed $Z_1 = Z_1^* = 0$.

- Equation (36) is the conventional object of interest in consumer theory.
- Partial derivatives of that function *holding E and the other arguments constant* have well defined economic interpretations.
- Suppose that some non-negligible proportion of the population buys none of the good Z_1 .
- Regression estimates of the parameters of (36) using Z_1^* approximate the conditional expectation

$$E(Z_1 \mid \Delta_1 = 1, \mathbf{P}, M) = \int_{\underline{E}} Z_1(\mathbf{P}, M, \varepsilon) dF(\varepsilon) \quad (37)$$

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

- Then the derivatives of (37) are, for the j th price

$$\frac{\partial E(Z_1 \mid \Delta = 1, \mathbf{P}, M)}{\partial P_j} = \int_E \frac{\partial Z_1(\mathbf{P}, M, \varepsilon)}{\partial P_j} dF(\varepsilon) \quad (38)$$

$$+ \lim_{\Delta P_j \rightarrow 0} \int_E \frac{[(I_{\underline{\underline{E}} + \Delta \underline{\underline{E}}_{P_j}}(\varepsilon) - I_{\underline{\underline{E}}}(\varepsilon))]Z(\mathbf{P}, M, \varepsilon)}{\Delta P_j} dF(\varepsilon).$$

- When the limit in the second term does not exist, the derivative does not exist.
- We assume for expositional convenience that the limit is well defined.

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Neither term is the same as the price derivative of (36) for an arbitrary value of $E = \varepsilon$ although the first term on the right-hand side of (38) approximates the price derivative of (36) for some value of $E = \varepsilon$.

- A set of navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and other slide controls.

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻

- Cain, G. G. and H. W. Watts (1973). Summary and overview. In *Income Maintenance and Labor Supply: Econometric Studies*. Chicago: Rand McNally College Publishing Company.
- Domencich, T. and D. L. McFadden (1975). *Urban Travel Demand: A Behavioral Analysis*. Amsterdam: North-Holland. Reprinted 1996.
- Flinn, C. and J. J. Heckman (1982a, January). New methods for analyzing structural models of labor force dynamics. *Journal of Econometrics* 18(1), 115–168.
- Flinn, C. J. and J. J. Heckman (1982b). New methods for analyzing individual event histories. *Sociological Methodology* 13, 99–140.
- Heckman, J. J. (1976, December). The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. *Annals of Economic and Social Measurement* 5(4), 475–492.

- Heckman, J. J. (1979, January). Sample selection bias as a specification error. *Econometrica* 47(1), 153–162.
- Heckman, J. J. and B. S. Singer (1986). Econometric analysis of longitudinal data. In Z. Griliches and M. D. Intriligator (Eds.), *Handbook of Econometrics*, Volume 3, Chapter 29, pp. 1690–1763. Amsterdam: North-Holland.
- Pearson, K. (1900). Mathematical contributions to the theory of evolution, VII: On the correlation of characters not quantitatively measureable. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 195(262–273), 1–47.
- Rao, C. R. (1965). On discrete distributions arising out of methods of ascertainment. In G. Patil (Ed.), *Classical and Contagious Discrete Distributions; Proceedings*. New York: Pergamon Press.