

Limited Dependent Variables: Part II

Truncated and Censored Samples

$$y_i = x_i' \beta + \epsilon_i \quad i = 1, \dots, N$$

In usual case, we have all observations on y_i & x_i with no irregularities

- Censoring occurs when we observe x_i for the entire sample but for some observations we only have limited information about the dependent variable

Example (Tobin's study of household expenditures)

Consumer maximizes utility by purchasing durable goods under constraint that total expenditures do not exceed income

expenditure on durables $>$, cost of least expensive durable good

If available income $<$ least expensive durable good
no expenditure is observed

Don't know how much a household would have spent if a durable good could be purchased for less than the least expensive item.

~~y_i = observed consumption of durables
 c = least expensive durable good
 y_i^* = optimal consumption
 $y_i = y_i^*$ if $y_i^* >$~~

Example (Lang's study of Prestige of 1st academic job & of Biochemists)

y_i = prestige of 1st academic job

Prestige is measured on a continuous scale from 1-5 with ~~schools~~ schools from

1-1.99 = adequate

2-2.99 = good

3-3.99 = strong

3.99+ = distinguished

Prestige is ~~not measured for~~ coded as 1 for
A) Graduate programs \oplus rates below adequate

Truncation

excludes all observations based on characteristics of the dependent variable

Example

(Hausman & Wise's 1977 analysis of New Jersey negative income tax experiment)

~~Example~~ Goal: Estimate earnings function for low income individuals

Individuals with ~~income~~^{earnings} less than 1.5 * poverty level were excluded from sample

Inference problems

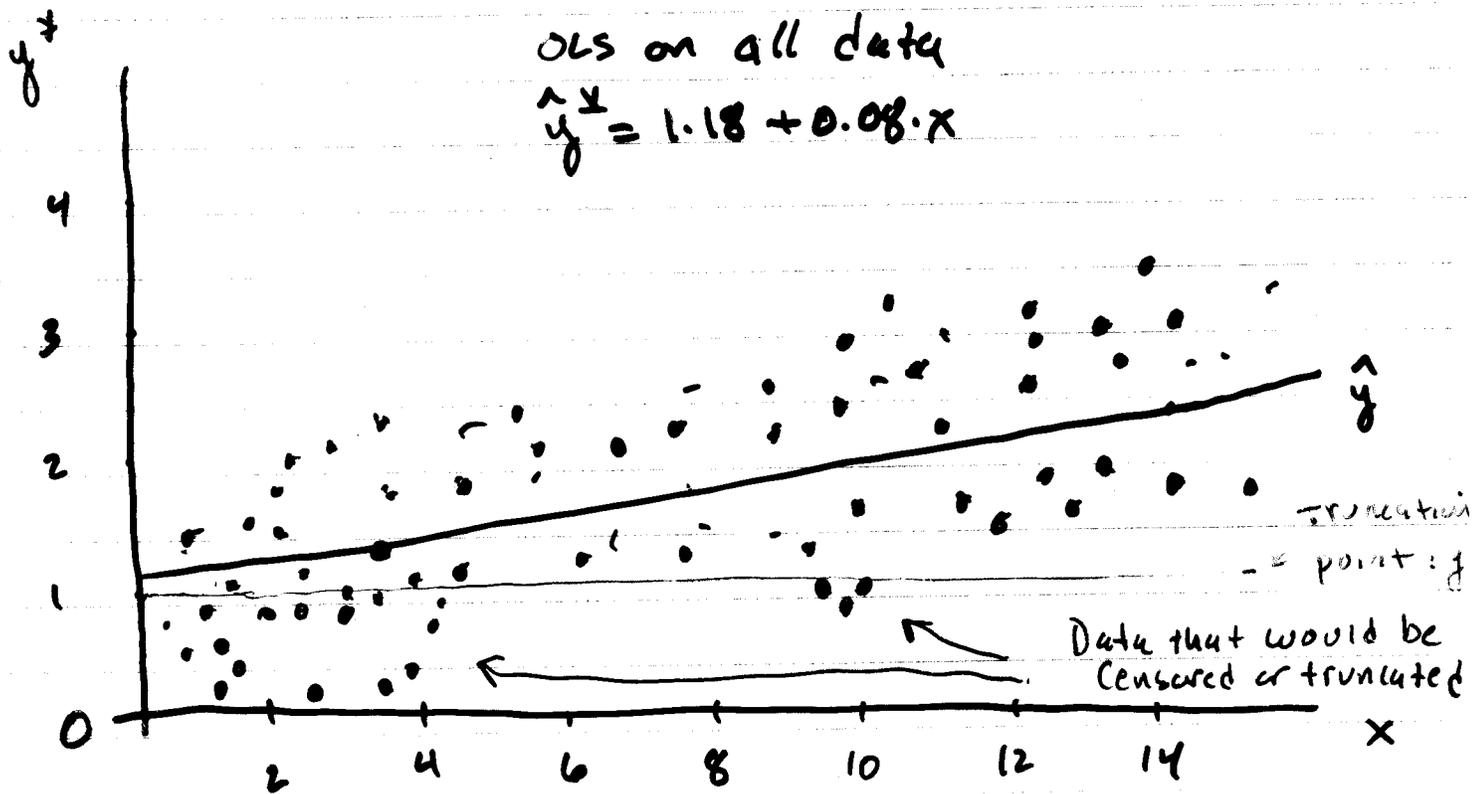
- Inference about entire population in presence of censoring / truncation
- Inference about subpopulation

Results

OLS has severe problems

MLE ~~can give~~ is appropriate alternative

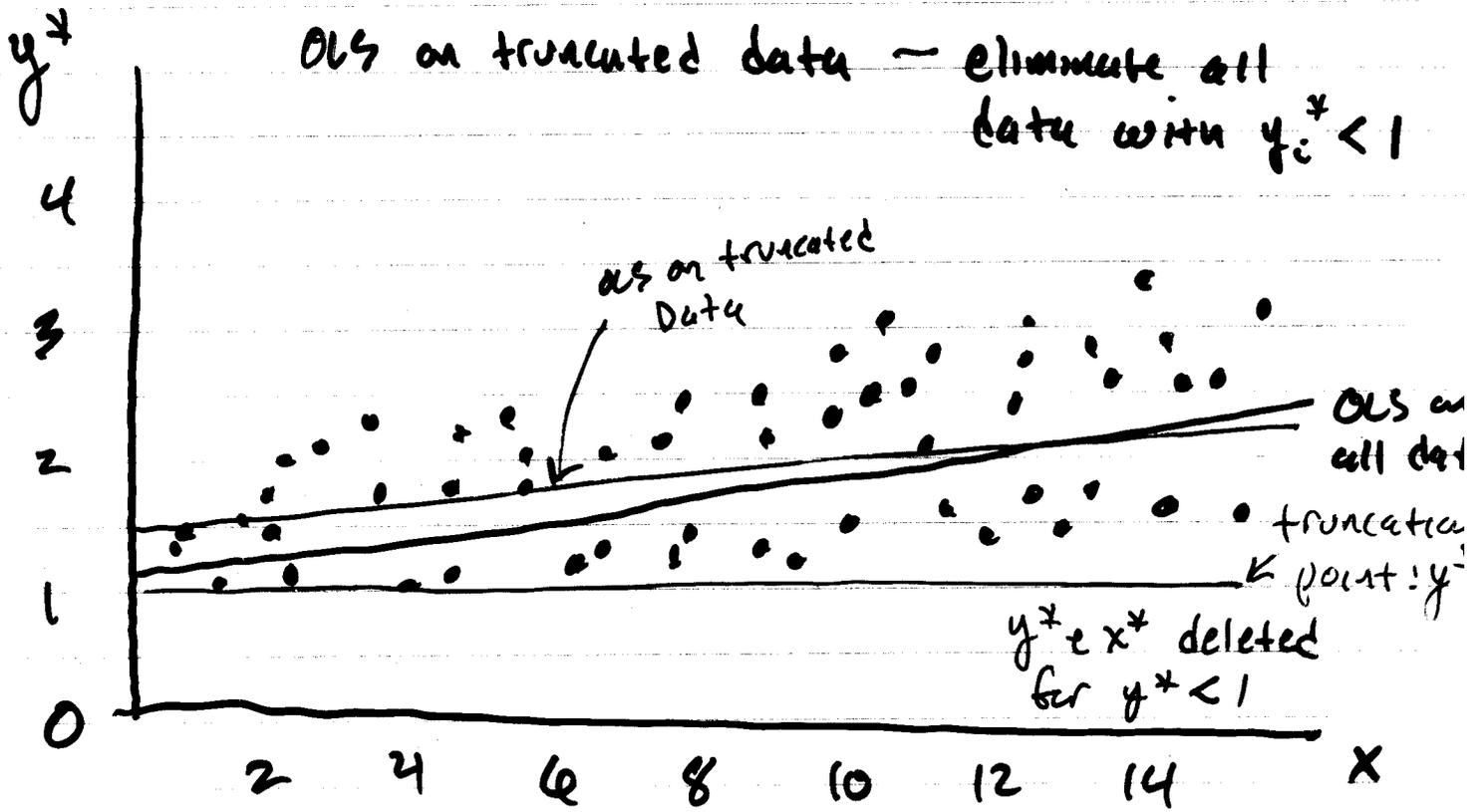
True Regression: $y^* = 1.2 + 0.08x + \epsilon$



Censored Data: $y_i = y_i^*$ if $y_i^* > 1$
 $= 0$ if $y_i^* \leq 1$

Truncated Data: Only observe $y_i = y_i^*$ and x_i
 if $y_i^* > 1$

Truncated Regression: $\hat{y}^* = 1.41 + 0.06kx$



- OLS on truncated Data is biased and inconsistent
- Truncation causes correlation between x and ϵ in general

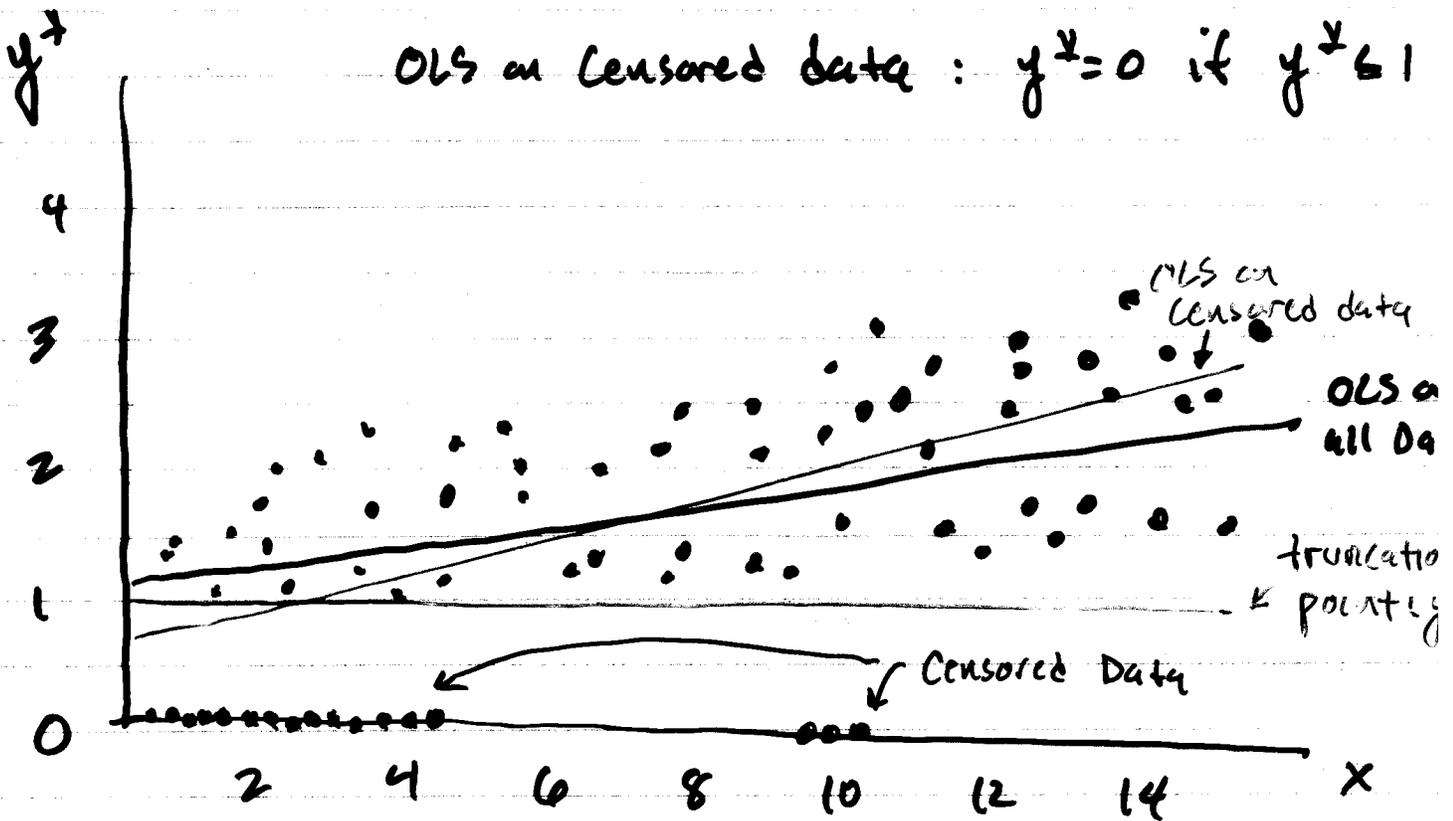
Here, observations with large negative errors have been deleted

intercept is overestimated

slope is underestimated

Censored Regression: $\hat{y}^x = 0.95 + 0.11x$

OLS on Censored data: $y^x = 0$ if $y^x \leq 1$



- OLS on censored data creates biased and inconsistent estimates of intercept & slope

Here

intercept is underestimated

slope is over estimated

Truncated Regression Model

truncation from below

$$y_i = \beta' x_i + \epsilon_i$$

where we observe x_i and y_i only if $y_i > a$

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

$$\Rightarrow y_i | x_i \sim N(\beta' x_i, \sigma^2)$$

Q: What is the pdf of $y_i | x_i, y_i > a$?

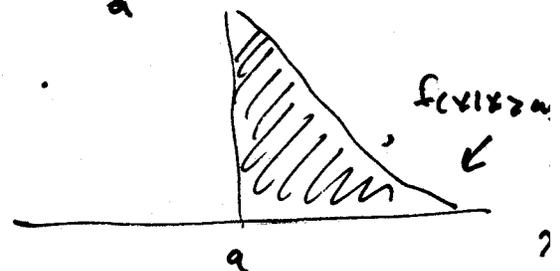
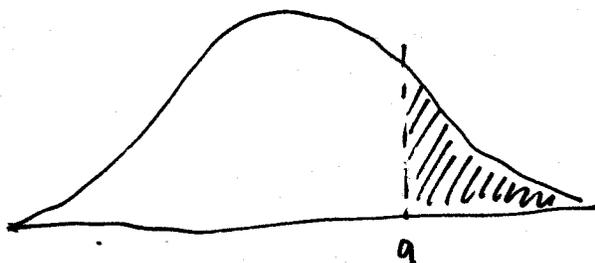
In order to MLE we don't use the unconditional pdf of y_i , we must use the pdf conditional on $y_i > a$

Truncated Distributions

Thm (see pg 449)

Let x be a cts rv with density $f(x)$ and let a be a constant. Then

$$f(x | x > a) = \frac{f(x)}{P(x > a)} = \frac{f(x)}{\int_a^{\infty} f(x) dx} \quad x > a.$$



Truncated Normal Dista (TN(μ, σ^2, a)) \uparrow $a = \text{truncation point}$

let $X \sim N(\mu, \sigma^2)$. Then

$$f(z; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left(\frac{z-\mu}{\sigma}\right)^2}$$

$$= \frac{1}{\sigma} \varphi\left(\frac{z-\mu}{\sigma}\right)$$

where $\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} = \text{pdf for Std. Normal}$

Further,

$$\begin{aligned} \Pr(X > a) &= \Pr\left(\frac{X-\mu}{\sigma} > \frac{a-\mu}{\sigma}\right) \\ &= \Pr\left(Z > \frac{a-\mu}{\sigma}\right), \quad Z \sim N(0,1) \\ &= 1 - \Pr\left(Z \leq \frac{a-\mu}{\sigma}\right) \\ &= 1 - \Phi\left(\frac{a-\mu}{\sigma}\right) \end{aligned}$$

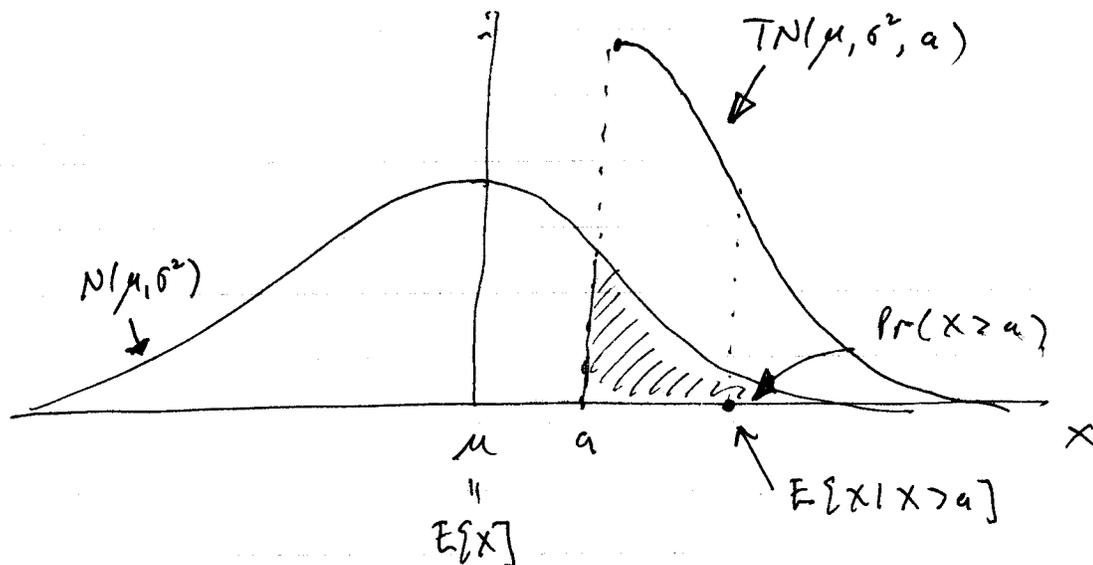
where $\Phi(z) = \text{CDF of } N(0,1) \text{ evaluated at } z$.

Result: ~~pdf~~ $f(x | x > a, \mu, \sigma^2) = \frac{f(x; \mu, \sigma^2)}{\Pr(X > a)}$

$$= \frac{\frac{1}{\sigma} \varphi\left(\frac{x-\mu}{\sigma}\right)}{1 - \Phi\left(\frac{a-\mu}{\sigma}\right)}$$

Note:

$$\int_a^{\infty} f(x | x > a, \mu, \sigma^2) dx = 1$$



Remarks

- (1) The mean of the truncated distribution is always larger than the mean of the untruncated distribution
- (2) The variance of the truncated distribution is always smaller than the untruncated distribution

Moments of a truncated Normal

let $X \sim N(\mu, \sigma^2)$ and let a be a constant

Then

$$(i) E[X | X > a] = \mu + \sigma \cdot \lambda\left(\frac{a-\mu}{\sigma}\right)$$

$$(ii) \text{Var}(X | X > a) = \sigma^2 \left(1 - \delta\left(\frac{a-\mu}{\sigma}\right)\right)$$

where

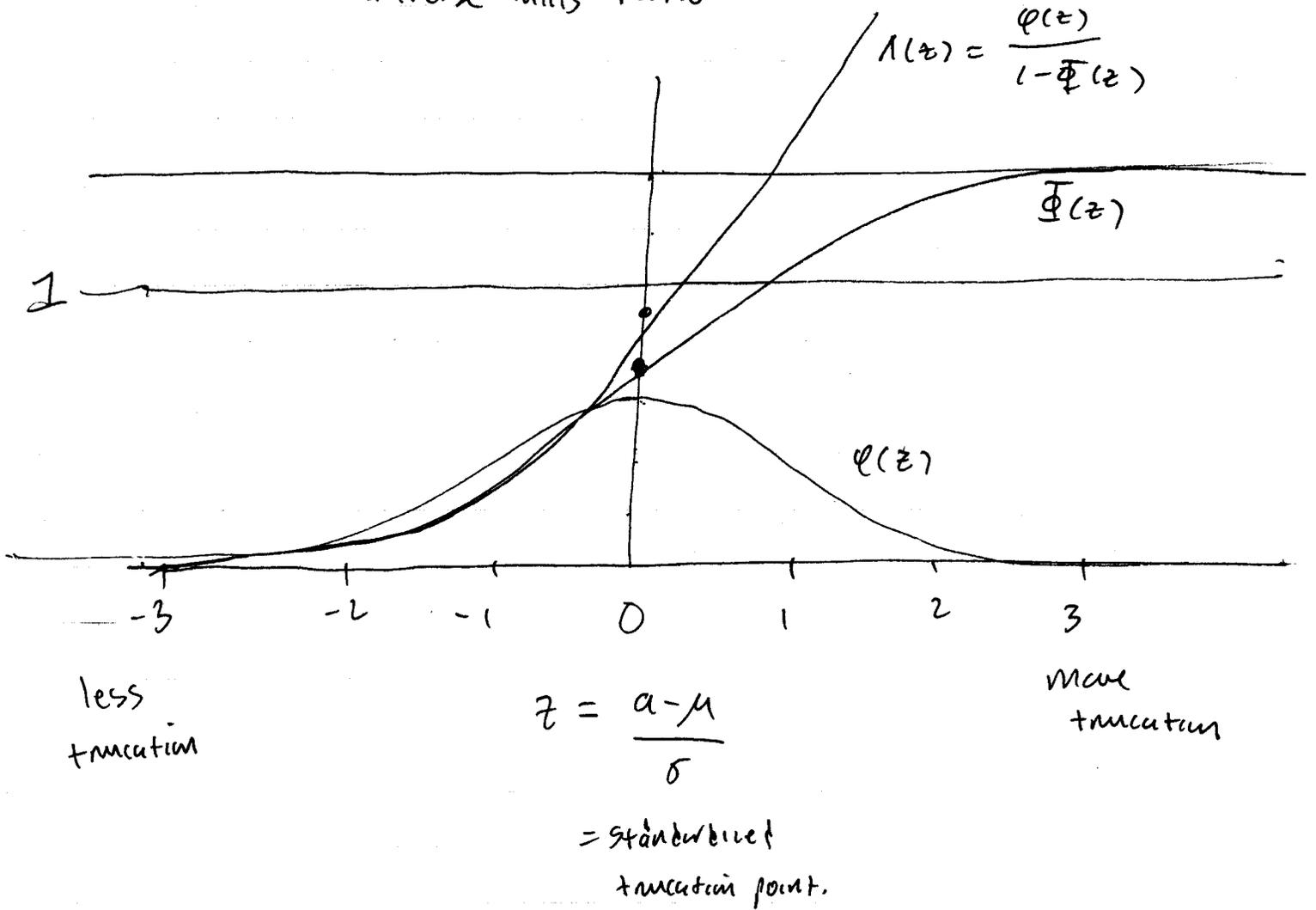
$$\lambda(x) = \frac{\varphi(x)}{1 - \Phi(x)} = \text{inverse mills ratio}$$
$$= \text{Normal hazard function}$$

$$\delta(x) = \lambda(x)(\lambda(x) - x), \quad 0 < \delta(x) < 1$$

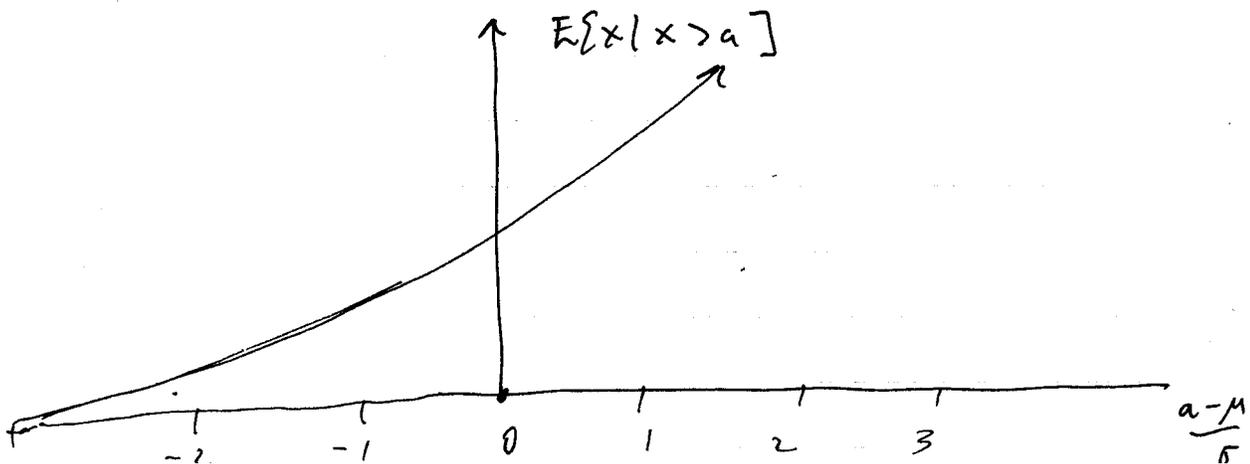
Note: prove using brute force manipulations; e.g.,

$$E[X | X > a] = \int_a^{\infty} x f(x | X > a, \mu, \sigma^2) dx$$
$$= \int_a^{\infty} x \frac{\frac{1}{\sigma} \varphi\left(\frac{x-\mu}{\sigma}\right)}{1 - \Phi\left(\frac{a-\mu}{\sigma}\right)} dx$$

Relationship b/w degree of truncation &
Inverse Mills Ratio



\Rightarrow Inverse Mills ratio increases with the amount of truncation!



The truncated Regression Model

$$y_i = x_i' \beta + \epsilon_i, \quad \epsilon_i \sim \text{iid } N(0, \sigma^2)$$

x_i is treated as fixed

$$\Rightarrow y_i | x_i \sim N(x_i' \beta, \sigma^2)$$

Truncation: only observe y_i, x_i for $y_i > a$
where a is some constant

Using the results now

$$\bullet E\{y_i | x_i\} = x_i' \beta$$

$$\frac{\partial E\{y_i | x_i\}}{\partial x_i} = \beta = \text{marginal effect for entire population}$$

Using the results on truncated distributions

$$E\{y_i | x_i, y_i > a\} = x_i' \beta + \sigma \lambda \left(\frac{a - x_i' \beta}{\sigma} \right)$$

and the marginal effects for the subpopulation are

$$\frac{\partial E\{y_i | x_i, y_i > a\}}{\partial x_i} = \beta + \sigma \lambda \left(\frac{a - x_i' \beta}{\sigma} \right)$$

which, by the chain rule, reduces to

$$= \beta \left(1 - \delta \left(\frac{a - x_i' \beta}{\sigma} \right) \right)$$

$$< \beta$$

since $\delta(x)$ lies b/w 0 and 1.

Maximum Likelihood Estimation

$y_i | x_i, y_i > a \sim$ truncated Normal so that

$$f(y_i; x_i, y_i > a, \beta, \sigma^2) = \frac{\frac{1}{\sigma} \varphi\left(\frac{y_i - x_i' \beta}{\sigma}\right)}{1 - \Phi\left(\frac{a - x_i' \beta}{\sigma}\right)}$$

Given an iid sample, the joint density is

$$L(\beta, \sigma^2 | y, X) = \prod_{i=1}^n \frac{\frac{1}{\sigma} \varphi\left(\frac{y_i - x_i' \beta}{\sigma}\right)}{1 - \Phi\left(\frac{a - x_i' \beta}{\sigma}\right)}$$

and the log-likelihood is

$$\begin{aligned} \ln L(\beta, \sigma^2 | y, X) &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i' \beta)^2 \\ &\quad - \sum_{i=1}^n \ln \left(1 - \Phi \left(\frac{a - x_i' \beta}{\sigma} \right) \right) \end{aligned}$$

Example: Hausman & Wise (1977) "Social Experimentation, truncated distributions, and efficient Estimation, Econometrica

Earnings regression $y_i = x_i' \beta + \epsilon_i$, $\epsilon_i \sim iid N(0, \sigma^2)$, $i = 1, \dots, 484$
 ↑
 log earnings

Truncation: y_i, x_i observed if $y_i < 1.5 \times \text{poverty level} = \5002

Variables	$\hat{\beta}_{OLS}$	$\hat{\beta}_{MLE}$	$2E[y_i y_i < 5002] / \sigma^2$
Const	8.20 (0.09)	9.10 (0.03)	
Education	0.010 (0.006)	0.015 (0.007)	0.0036
IQ	0.002 (0.002)	0.006 (0.005)	0.0014
Training	0.002 (0.002)	0.007 (0.003)	0.0014
Union	0.090 (0.031)	0.246 (0.090)	
Illness	-0.076 (0.040)	-0.226 (0.107)	
Age	-0.003 (0.002)	-0.016 (0.005)	-0.0039

Remark

The bias & inconsistency of OLS in the truncated regression model can be explained in terms of omitted variables bias. For the truncated regression model, the regression function takes the form

$$y_i = E\{y_i | x_i, y_i > a\} + \epsilon_i$$

where $\epsilon_i = y_i - E\{y_i | x_i, y_i > a\}$ is the regression error. Using the results in the mean of a

truncated normal r.v.

$$E\{y_i | x_i, y_i > a\} = x_i' \beta + \sigma \cdot \lambda \left(\frac{a - x_i' \beta}{\sigma} \right)$$

inverse mills ratio



So that the correct regression function is

$$y_i = x_i' \beta + \sigma \cdot \lambda \left(\frac{a - x_i' \beta}{\sigma} \right) + \epsilon_i$$

Standard OLS estimates

$$y_i = x_i' \beta + u_i$$

and which omits the term $\sigma \cdot \Lambda \left(\frac{a - x_i' \beta}{\sigma} \right)$

So that we may think of u_i as

$$u_i = \sigma \cdot \Lambda \left(\frac{a - x_i' \beta}{\sigma} \right) + \epsilon_i$$

Consequently, x_i in the OLS regression is correlated

with u_i (x_i is correlated with $\Lambda \left(\frac{a - x_i' \beta}{\sigma} \right)$)

due to the correlation between x_i and the
omitted variable $\Lambda \left(\frac{a - x_i' \beta}{\sigma} \right)$.