Math 152    1/27/09

## Simple Random Sampling (Cont.).

Last time we introduced the notion of biased and unbiased estimates and saw that

$$E(\bar{X}) = \mu$$
$$E(T) = \tau$$

so that $\bar{X}$ and $T$ are unbiased estimates of $\mu$ and $\tau$ respectively.

Recall also that on the way to the above we saw that

$$E(X_i) = \mu$$
$$\text{and } Var(X_i) = \sigma^2.$$

$\mu$ = true population mean

$\sigma^2$ = true population variance.

What is $\text{Var}(\bar{X})$ ?

$$= \text{Var}\left(\frac{1}{n}\sum_{i=1}^{n} X_i\right) = \frac{1}{n^2}\text{Var}\left(\sum_{i=1}^{n} X_i\right).$$

If the $X_i$ were independent

(e.g. if we sampled with replacement)

this would be

$$= \frac{1}{n^2}\sum_{i=1}^{n}\text{Var}(X_i) = \frac{\sigma^2}{n}.$$

and we would have, for the standard deviation of $\bar{X}$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

This _almost_ true. The $X_i, X_j$ $i \neq j$ are only very weakly dependent if $n$ is small compared to $N$.

$$\text{Var}(\bar{X}) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^{n} X_i\right)$$

$$= \frac{1}{n^2} \text{Cov}\left(\sum_{i=1}^{n} X_i, \sum_{i=1}^{n} X_i\right)$$

$$= \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \text{Cov}(X_i, X_j)$$

By the calculation made in the proof of Lemma B on pg 207 of your text

$$\text{Cov}(X_i, X_j) = -\sigma^2/N-1$$

$$\text{if} \quad i \neq j.$$

( See problems 25 and 26 ).

So

$$\text{Var}(\bar{X}) = \frac{1}{n^2}\left(\sum_{i=1}^{n} \text{Var}(X_i) + \sum\sum_{i \neq j} \text{Cov}(X_i, X_j)\right)$$

$$= \frac{1}{n^2}\left(n \cdot \sigma^2 + (n^2-n)\left(\frac{-\sigma^2}{N-1}\right)\right)$$

$$= \frac{\sigma^2}{n}\left(1 - \frac{n-1}{N-1}\right).$$

and

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \sqrt{1 - \frac{n-1}{N-1}}.$$

$$\approx \frac{\sigma}{\sqrt{n}} \quad \text{when} \quad \frac{n}{N} \text{ is small.}$$

e.g. hospitals (discharges).

$n = 32$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \sqrt{1 - \frac{n-1}{N-1}}$$

$$= \frac{589.7}{\sqrt{32}} \sqrt{1 - \frac{31}{392}}$$

$$\approx (104.2)(.96).$$

$$\approx 100.$$

On the figure ( produced in class using R or in your text )

most of the observations ( out of 500 ) were within two standard errors of the mean. (814)

i.e. in (614, 1014).

e.g.    estimating proportions

$$\theta_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} \; \sqrt{1 - \frac{n-1}{N-1}}$$

with $n = 32$, $N = 393$.

$$\theta_{\hat{p}} = \sqrt{\frac{.654 \times .346}{32}} \; \sqrt{1 - \frac{31}{392}}$$

$$\simeq .08.$$

Compare with the computer simulation.

Though $\theta_{\bar{X}} \simeq \frac{\theta}{\sqrt{n}}$ is

nearly independent of $N$,

we aren't as fortunate for $\theta_T$

$$T = N \bar{X}$$

$$Var(T) = N^2 \, Var(\bar{X}) = N^2 \left(\frac{\theta^2}{n}\right) \left(\frac{N-n}{N-1}\right)$$

In applications, we will not know the true value of the population variance and we need to estimate it.

$$\hat{\sigma}^2 = \frac{1}{b} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

is a natural candidate for an estimate ...

But, ...

$$E(\hat{\sigma}^2) = E\left(\frac{1}{n} \sum_{i=1}^{n} X_i^2 - \bar{X}^2\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} E(X_i^2) - E(\bar{X}^2)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left(Var(X_i) + E(X_i)^2\right) - E(\bar{X}^2)$$

$$= \frac{1}{n} \sum_{i=1}^{n} (\sigma^2 + \mu^2) - E(\bar{X}^2)$$

$$= \sigma^2 + \mu^2 - Var(\bar{X}) - (E(\bar{X}))^2$$

$$= \sigma^2 - \frac{\sigma^2}{n}\left(1 - \frac{n-1}{N-1}\right) = \sigma^2 \left(\frac{n-1}{n}\right)\frac{N}{N-1}$$

so $\hat{\sigma}^2$ is a biased estimate of $\sigma^2$.

In fact $\left(\frac{n-1}{n}\right)\frac{N}{N-1} < 1$

$$\text{if } n < N$$

so $E(\hat{\sigma}^2) < \sigma^2$.

An unbiased estimate of $\sigma^2$

is $\frac{N-1}{N} \cdot \frac{n}{n-1} \cdot \hat{\sigma}^2$.

$$= \frac{1}{n-1}\left(1-\frac{1}{N}\right)\sum_{i=1}^{n}(X_i-\bar{X})^2.$$

Now since

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n}\left(\frac{N-n}{N-1}\right)$$

we obtain an unbiased estimate of $\text{Var}(\bar{X})$ by replacing $\sigma^2$ with $\frac{N-1}{N}\frac{n}{n-1}\hat{\sigma}^2$

in the previous expression.

We get the unbiased estimate

$$S_{\bar{X}}^2 = \frac{\hat{\sigma}^2}{n}\left(\frac{n}{n-1}\right)\left(\frac{N-1}{N}\right)\left(\frac{N-n}{N-1}\right)$$

$$= \frac{S^2}{n}\left(1 - \frac{n}{N}\right)$$

where

$$S^2 = \frac{1}{n-1}\sum_{i=1}^{n}\left(X_i - \bar{X}\right)^2.$$

It follows that

$$S_T^2 = N^2\, S_{\bar{X}}^2$$

is an unbiased estimate of $\mathrm{Var}(T)$.
In the special case where $X_i = 0$ or $1$
for each $i$ we have.

$$S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

$$= \left(\frac{n}{n-1}\right) \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

$$= \left(\frac{n}{n-1}\right) \left[ \frac{1}{n} \left(\sum_{i=1}^{n} X_i^2\right) - \bar{X}^2 \right]$$

$$= \left(\frac{n}{n-1}\right) \left[ \hat{p} - \hat{p}^2 \right].$$

$$= \frac{n}{n-1} \hat{p}(1-\hat{p})$$

By the above, we have

$$S_{\hat{p}}^2 = \frac{1}{n-1} \hat{p}(1-\hat{p}) \left(1 - \frac{n}{N}\right)$$

as an   unbiased estimate of $Var(\hat{p})$.

# Normal approximation to the distribution of $\overline{X}$

If $X_1, \ldots, X_n$ are i.i.d.

with mean $\mu$ and variance $\sigma^2$

$\boxed{\text{Let } \overline{X_n} = \frac{1}{n} \sum_{i=1}^{n} X_i}$

then

$$E(\overline{X_n}) = \frac{1}{n} \sum_{i=1}^{n} E(X_i) = \mu.$$

$$\text{Var}(\overline{X_n}) = \frac{1}{n^2} \sum_{i=1}^{n} \text{Var}(X_i) = \frac{\sigma^2}{n}.$$

and C.L.T $\Rightarrow$

$$P\left( \frac{\overline{X_n} - \mu}{\sigma/\sqrt{n}} \leq x \right) \to \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-u^2/2} \, du.$$

as $n \to \infty$. $\left( \Phi(x) \right)$

In simple random sampling, the $X_i$ are not independent but $\text{Cov}(X_i, X_j) = \frac{-\sigma^2}{N-1}$ $i \neq j$

is small for large $N$

and when ("<<" means "much less")

$$1 << n \lessapprox < N.$$

an approximate C.L.T applies.

e.g. when $1 << n << N$

we may approximate

$$P(|\bar{X} - \mu| \le \delta)$$

$$= P(-\delta \le \bar{X} - \mu \le \delta)$$

$$= P\left(-\frac{\delta}{\sigma_{\bar{X}}} \le \frac{\bar{X} - \mu}{\sigma_{\bar{X}}} \le \frac{\delta}{\sigma_{\bar{X}}}\right)$$
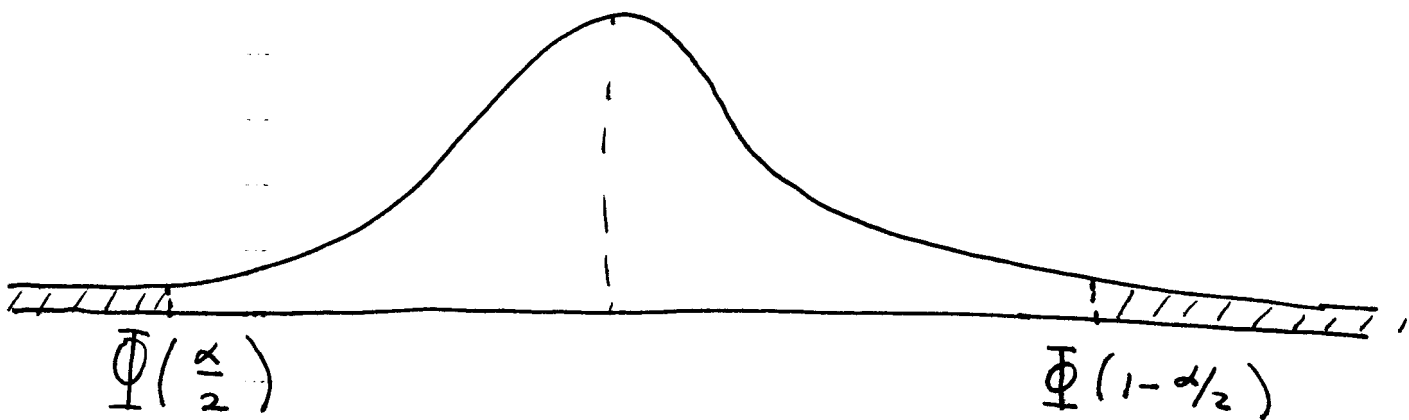
recall that
$$\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}\left(1 - \frac{n-1}{N-1}\right).$$

$$= \Phi\left(\frac{\delta}{\sigma_{\bar{X}}}\right) - \Phi\left(-\frac{\delta}{\sigma_{\bar{X}}}\right) = 2\Phi\left(\frac{\delta}{\sigma_{\bar{X}}}\right) - 1.$$

In practice, we will not know $\sigma^2$ (or $N$) and we will estimate $\sigma_{\bar{X}}^2$ by

$$S_{\bar{X}}^2 \left( or \quad \frac{S^2}{n} \right)$$

## Confidence Intervals



$$\Phi\left(\frac{\alpha}{2}\right) \qquad\qquad \Phi\left(1-\alpha/2\right)$$

If $X$ is $N(0,1)$ then

$$\text{Prob} \left\{ \Phi\left(\frac{\alpha}{2}\right) \leq X \leq \Phi\left(1-\frac{\alpha}{2}\right) \right\} = 1-\alpha.$$

Note that $\Phi\left(\frac{\alpha}{2}\right) = -\Phi\left(1-\frac{\alpha}{2}\right).$

Using our $\overset{\text{normal}}{\text{approximation}}$ to the distribution of $\overline{X}$, we assert that

$$P\left( \underset{(-\Phi(1-\frac{\alpha}{2}))}{\Phi\left(\frac{\alpha}{2}\right)} \leq \frac{\overline{X} - \mu}{S_{\overline{X}}} \leq \Phi\left(1-\frac{\alpha}{2}\right) \right)$$

$$\approx 1 - \alpha.$$

i.e.

$$P\left( \overline{X} - S_{\overline{X}}\,\Phi\left(1-\tfrac{\alpha}{2}\right) \leq \mu \leq \overline{X} + S_{\overline{X}}\,\Phi\left(1-\tfrac{\alpha}{2}\right) \right)$$

$$\approx 1 - \alpha.$$

We call

$$\left( \overline{X} - \Phi\left(1-\tfrac{\alpha}{2}\right)S_{\overline{X}} \;,\; \overline{X} + \Phi\left(1-\tfrac{\alpha}{2}\right)S_{\overline{X}} \right)$$

a $(1-\alpha)\%$ confidence interval for $\mu$.

Notice that it is the interval which is random and that $\mu$ is a fixed value.

---

Computer examples.
- confidence intervals
- example P pg 219.

---

## Estimation of a Ratio

If $X$ is a random variable, and we know something about its 1st two moments $E(X)$ and $E(X^2)$ then we can sometimes say something useful about the 1st two moments of

$$Y = g(X)$$

if the function $g$ is not badly behaved.

We expand $g$ in a Taylor series about $\mu_x$

$$Y = g(X) \simeq g(\mu_x) + (X - \mu_x) g'(\mu_x)$$
$$+ \frac{(X - \mu_x)^2}{2} g''(\mu_x)$$
$$+ \ldots$$

and from this we get.

$$E(Y) \simeq g(\mu_x) + \frac{g''(\mu_x)}{2} Var(X).$$

$$Var(Y) \simeq g'(\mu_x)^2 Var(X).$$

We know (by Tchebychev's inequality) that $X$ will tend to be near $\mu_x$.

~~If it is close enough +~~

If $X$ is usually (i.e. w/ large prob.)

in a nbhd of $\mu_x$ for which $g$ is well approximated by the 1st two terms of its Taylor expansion, then these approximations will be reasonably good.

Similarly, if $X, Y$ are r.v's about whose 1st two moments, we have some information then we may argue as follows.

$$Z = g(X, Y).$$

$$Z = g(X, Y)$$

$$\approx g(\mu_X, \mu_Y) + (X - \mu_X)\frac{\partial g}{\partial x}(\mu_X, \mu_Y)$$

$$+ (Y - \mu_Y)\frac{\partial g}{\partial y}(\mu_X, \mu_Y)$$

$$+ \frac{1}{2}(X - \mu_X)^2 \frac{\partial^2 g}{\partial x^2}(\mu_X, \mu_Y)$$

$$+ (X - \mu_X)(Y - \mu_Y)\frac{\partial^2 g}{\partial x \partial y}(\mu_X, \mu_Y)$$

$$\vec{\mu} = (\mu_X, \mu_Y) \qquad\qquad + \frac{1}{2}(Y - \mu_Y)^2 \frac{\partial^2 g}{\partial y^2}(\mu_X, \mu_Y)$$

$$+ \dots$$

to obtain

$$E(Z) \approx g(\mu) + \frac{1}{2}\sigma_X^2 \frac{\partial^2 g}{\partial x^2}(\vec{\mu}) + \frac{1}{2}\sigma_Y^2 \frac{\partial^2 g}{\partial y^2}(\vec{\mu})$$

$$+ \sigma_{XY}\frac{\partial^2 g(\vec{\mu})}{\partial x \partial y}$$

$$Var(Z) \approx \sigma_X^2\left(\frac{\partial g(\vec{\mu})}{\partial x}\right)^2 + \sigma_Y^2\left(\frac{\partial g(\vec{\mu})}{\partial y}\right)^2 + 2\sigma_{XY}\left(\frac{\partial g(\vec{\mu})}{\partial x}\right)\left(\frac{\partial g(\vec{\mu})}{\partial y}\right)$$

The application we have in mind is to.

$$Z = \frac{Y}{X} = g(X, Y).$$

for which we have.

$$\frac{\partial g}{\partial x} = -\frac{y}{x^2} \qquad \frac{\partial g}{\partial y} = \frac{1}{x}$$

$$\frac{\partial^2 g}{\partial x^2} = \frac{2y}{x^3} \qquad \frac{\partial^2 g}{\partial y^2} = 0.$$

$$\frac{\partial^2 g}{\partial x \partial y} = -\frac{1}{x^2}.$$

If $\mu_x \neq 0$ then we get

$$E(Z) \simeq \frac{\mu_Y}{\mu_x} + \sigma_X^2 \frac{\mu_Y}{\mu_x^3} - \frac{\sigma_{XY}}{\mu_x^2}$$

$$= \frac{\mu_Y}{\mu_x} + \frac{1}{\mu_x^2}\left( \sigma_X^2 \frac{\mu_Y}{\mu_x} - \rho \sigma_X \sigma_Y \right)$$

$$\rho = corr(X, Y).$$

also

$$Var(z) \approx \sigma_x^2 \frac{\mu_Y^2}{\mu_X^4} + \frac{\sigma_Y^2}{\mu_X^2} - 2\sigma_{XY}\frac{\mu_Y}{\mu_X^3}$$

$$= \frac{1}{\mu_X^2} \left( \sigma_x^2 \frac{\mu_Y^2}{\mu_X^2} + \sigma_y^2 - 2\rho\sigma_x\sigma_y\frac{\mu_Y}{\mu_X} \right)$$

Note that

$\left| E(z) - \frac{\mu_Y}{\mu_X} \right|$ can be large if

$\mu_X$ is small or if

$\sigma_X, \sigma_Y$ are large.

also.

$Var(z)$ is large if $\mu_X$ is small

but correlation between $X, Y$ which

is of the same sign as $\frac{\mu_Y}{\mu_X}$

decreases the variance.