

1. 5-1, 5-2, 5-3

Large n normal approximations (Central Limit Theorem).

$\bar{x} \sim N[\mu, \sigma^2 / n]$
(sketch a normal with mean μ and $sd = \sigma / \sqrt{n}$).

$\hat{p} \sim N[p, pq / n]$
(sketch a normal with mean p and $sd = \sqrt{pq} / \sqrt{n}$).

Use FPC = $\sqrt{(N-n) / (N-1)}$ as appropriate.

ME

$$\bar{x} \pm 1.96 s / \sqrt{n}$$

$$\hat{p} \pm 1.96 \sqrt{\hat{p} \hat{q}} / \sqrt{n}$$

2. 6-3, 6-4, 6-6

Large n confidence intervals

$$\bar{x} \pm z s / \sqrt{n}$$

$$\hat{p} \pm z \sqrt{\hat{p} \hat{q}} / \sqrt{n}$$

Use FPC = $\sqrt{(N-n) / (N-1)}$ as appropriate.

Claim: $P(\mu \text{ in } z\text{-CI}) \sim P(Z \text{ in } [-z, z])$ n large

$P(p \text{ in } z\text{-CI}) \sim P(Z \text{ in } [-z, z])$ n large

3. Any $n > 1$ confidence intervals if population is normal.

$$\bar{x} \pm t s / \sqrt{n}$$

Example problem. A population distribution is known to be “in control,” that is “approximately normal.” A sample of 12 finds $\bar{x} = 16.3$ with $s = 7.1$. Give a 95% confidence interval for population mean μ .

ans. We need the t-score (replaces z-score) for degrees of freedom $n - 1 = 12 - 1 = 11$). Consult the t-table, inside front cover to the right, for degrees of freedom = $n - 1$). Here is how it looks,

Critical Values of the t Distribution

Degrees of Freedom	
11	2.201
infinity	1.96
C.I.	95%

So the desired t-score for 11 degrees of freedom and confidence level 0.95 is $t = 2.201$. The 95% CI is therefore

$$\bar{x} \pm t s / \sqrt{n}$$

$$16.3 \pm 2.201 \cdot 7.1 / \sqrt{12}$$

Claim: For sample of any size $n > 1$ from a **normal** population :
 $P(\mu \text{ in } \bar{x} \pm t s / \sqrt{n}) = 0.95$ (not just ~ 0.95).

4. Supplement to the readings: Stratified sampling produces a better estimate for μ (general x scores).

Example question. A population of accounts is divided into two types, those having good credit and those not. We are interested in estimating the population average balance μ of all accounts. However, we will need to audit each account we include in our

sample and this is costly so we need to do things as efficiently as possible on a per-sample basis. Suppose it is known that

22% of accounts do not have good credit

We decide to invest in 100 samples. A statistically trained employee suggests that we draw a random sample of 22 accounts from those that do not have good credit and 78 from those that do. This is done, from which we find

$$\bar{x} \text{ from 22} = 147.21 \quad \bar{x} \text{ from 78} = 194.49$$

$$s \text{ from 22} = 75.38 \quad s \text{ from 78} = 110.33$$

What is the stratified estimator of the overall population mean μ based upon this stratified sample and what is a 95% confidence interval for μ based on this stratified estimator?

ans. The stratified estimate of μ is

$$\bar{x}_{\text{STRAT}} = 0.22 \cdot 147.21 + 0.78 \cdot 194.49$$

and the margin of error is

$$\pm 1.96 \sqrt{0.22^2 \cdot 75.38^2 / 22 + 0.78^2 \cdot 110.33^2 / 78}$$

The sample of 22 would be regarded as rather too small for application of this method based on large n z -approximations.

Defining the stratified estimator.

- Divide the population into two or more disjoint sub-populations (strata).
- From each of stratum $_i$ select a with-replacement sample of size n_i . The samples are to be independent between strata also.
- Form the stratified estimate of the population mean μ

$$\bar{x}_{\text{strat}} = \text{sum, over all strata } i, \text{ of } (W_i \cdot \bar{x}_{\text{BAR}_i})$$

where

$W_i = N_i / N$ is the relative size of stratum i in the population

\bar{x}_{BAR_i} is the mean of stratum i sample scores.

Using E and Var rules to evaluate the expectation and variance of the stratified estimator.

$$\begin{aligned}
E \bar{x}_{\text{STRAT}} &= \text{sum of } W_i E \bar{x}_{i} \quad (\text{linearity of } E) \\
&= \text{sum of } W_i \mu_i \quad (\text{each } \bar{x}_{i} \text{ is unbiased}) \\
&= \mu \quad (\mu \text{ is the wtd avg of subpopulation means})
\end{aligned}$$

That is, the stratified estimator \bar{x}_{strat} is **unbiased** as an estimator of the population mean μ .

$$\begin{aligned}
\text{Var } \bar{x}_{\text{STRAT}} &= \text{sum of } \text{Var} (W_i \bar{x}_{i}) \quad (\bar{x}_{i} \text{ are indep}) \\
&= \text{sum of } (W_i^2 \text{Var } \bar{x}_{i}) \\
&= \text{sum of } (W_i^2 \sigma_i^2 / n_i)
\end{aligned}$$

ME

This leads to the following definition of ME for the stratified estimator of μ

ME for $\bar{x}_{\text{STRAT}} = \pm 1.96 \text{ root}(\text{sum of } (W_i^2 s_i^2 / n_i))$
where s_i denotes the sample standard deviation obtained from the sample of stratum i .

Important fact: In the special case $n_i = W_i n$ (called proportionally stratified sampling and used in the example problem above) the theoretical sd of \bar{x}_{STRAT} (obtained just above) is never larger than the sd of \bar{x} .

In practice it is usually the case that the estimator \bar{x} (without stratification) and its competitor \bar{x}_{STRAT} (with stratification) are each approximately normally distributed. This makes it easy to compare their performances. Since they have the same expectation μ (each is unbiased for the population mean μ), and the **proportionally stratified** estimator \bar{x}_{STRAT} has the smaller sd, it is the better estimator. To see it just think of comparing two normal curves, both having the same mean μ , but one having a smaller sd.

5. Supplement to the readings: Outline of a ME applicable to stratified sampling for p (0-1 scores).

Example question. A population of accounts is divided into two types, those having good credit and those not. We are interested in estimating the population fraction p of all accounts that will respond to an offer we plan to make them. However, it is costly to make this offer properly so we need to do things as efficiently as possible on a per-sample basis. Suppose it is known that

30% of accounts are held by women

We decide to invest in 50 samples. A statistically trained employee suggests that we draw a random sample of 50 $0.3 = 15$ accounts from those held by women and 35 from those held by men. This is done, from which we find

$$\text{pHAT from 15} = 6/15 \quad \text{pHAT from 35} = 20/35$$

What is the stratified estimator of the overall population proportion p based upon this stratified sample and what is a 95% confidence interval for p based on this stratified estimator?

ans. The stratified estimate of p is

$$\text{pHAT}_{\text{STRAT}} = 0.3 \cdot 6/15 + 0.7 \cdot 20/35$$

and the margin of error is

$$\pm 1.96 \sqrt{0.3^2 (6/15 - 9/15) / 15 + 0.7^2 (20/35 - 15/35) / 35}$$

The sample size of 15 would be regarded as rather too small and insufficient to support the use of this large n method.

Defining the stratified estimator.

- Divide the population into two or more disjoint sub-populations (strata).
- From each of stratum $_i$ select a with-replacement sample of size n_i . The samples are to be independent between strata also.

c. Form the stratified estimate of the population proportion p having some particular characteristic of interest

$$\text{pHAT}_{\text{strat}} = \text{sum, over all strata } i, \text{ of } (W_i \text{ pHAT}_i)$$

where

$W_i = N_i / N$ is the relative size of stratum i in the population

pHAT_i is the stratum i proportion for sample scores.

Then

$$\begin{aligned} E \text{pHAT}_{\text{strat}} &= \text{sum of } W_i E \text{pHAT}_i \quad (\text{linearity of } E) \\ &= \text{sum of } W_i p_i \quad (\text{each } \text{pHAT}_i \text{ is unbiased for } p_i) \\ &= p \quad (p \text{ is the wtd avg of subpopulation proportions}) \end{aligned}$$

That is, the stratified estimator $\text{pHAT}_{\text{strat}}$ is also **unbiased** as an estimator of the population proportion p .

$$\begin{aligned} \text{Var } \text{pHAT}_{\text{strat}} &= \text{sum of } \text{Var} (W_i \text{ pHAT}_i) \quad (\text{pHAT}_i \text{ are indep}) \\ &= \text{sum of } (W_i^2 \text{ Var } \text{pHAT}_i) \\ &= \text{sum of } (W_i^2 p_i q_i / n_i) \end{aligned}$$

ME

This leads to the following definition of ME for the stratified estimator of μ

$$\begin{aligned} \text{ME for } \text{pHAT}_{\text{STRAT}} \\ &= \pm 1.96 \text{ root}(\text{sum of } (W_i^2 \text{ pHAT}_i \text{ qHAT}_i / n_i)) \end{aligned}$$

where $\text{pHAT}_i \text{ qHAT}_i$ is the estimated standard deviation of i -th stratum (0-1 scores), obtained from the sample of stratum i .

Important fact: In the special case $n_i = W_i n$ (called proportionally stratified sampling and used in the example problem above) the theoretical sd of $\text{pHAT}_{\text{STRAT}}$ (obtained just above) is never larger than the sd of pHAT .

In practice it is usually the case that the estimator pHAT (without stratification) and its competitor $\text{pHAT}_{\text{STRAT}}$ (with stratification)

are each approximately normally distributed. This makes it easy to compare their performances. Since they have the same expectation p , and the **proportionally** stratified estimator \hat{p}_{strat} has the smaller sd, we prefer to use the stratified estimator. To see it just think of comparing two normal curves, both having the same mean μ , but one having a smaller sd.