

Definition 8 (sample mean)

The sample mean of X_1, X_2, \dots, X_n is $\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$.

\uparrow cf. population mean: μ of F_0 \uparrow data

a function of data

$$\begin{aligned} X &\sim F_0 \\ \mu &= E(X) \\ &= \frac{1}{N} \sum_{i=1}^N x_i \\ &= \sum_{j=1}^m \left(\frac{n_j}{N} \right) c_j \end{aligned}$$

Note 4 (Some notes about sample mean)

- \bar{X} is clearly a statistic, and hence a random variable.

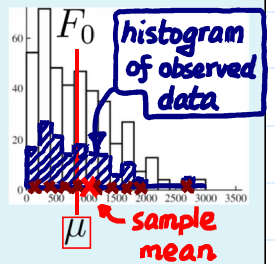
- \bar{X} is an intuitive estimator of μ .

- In the dichotomous case, we have $\mu = p$ and $\sum_{k=1}^n X_k = \text{\# of 1's in the sample}$. $\hat{p} \equiv \bar{X}$ is the sample proportion.

In X_1, \dots, X_n , observe c_1 about $(n \cdot \frac{n_1}{N})$ times c_m about $(n \cdot \frac{n_m}{N})$ times

$$\Rightarrow \sum_{k=1}^n X_k \approx n \times \sum_{j=1}^m \left(\frac{n_j}{N} \right) c_j$$

population proportion



The statistic that contributes the most to the world.

Example 3 (sampling distribution of sample mean, cont. Ex.2 in LNp.4)

- Consider the population of $N=393$ hospitals.

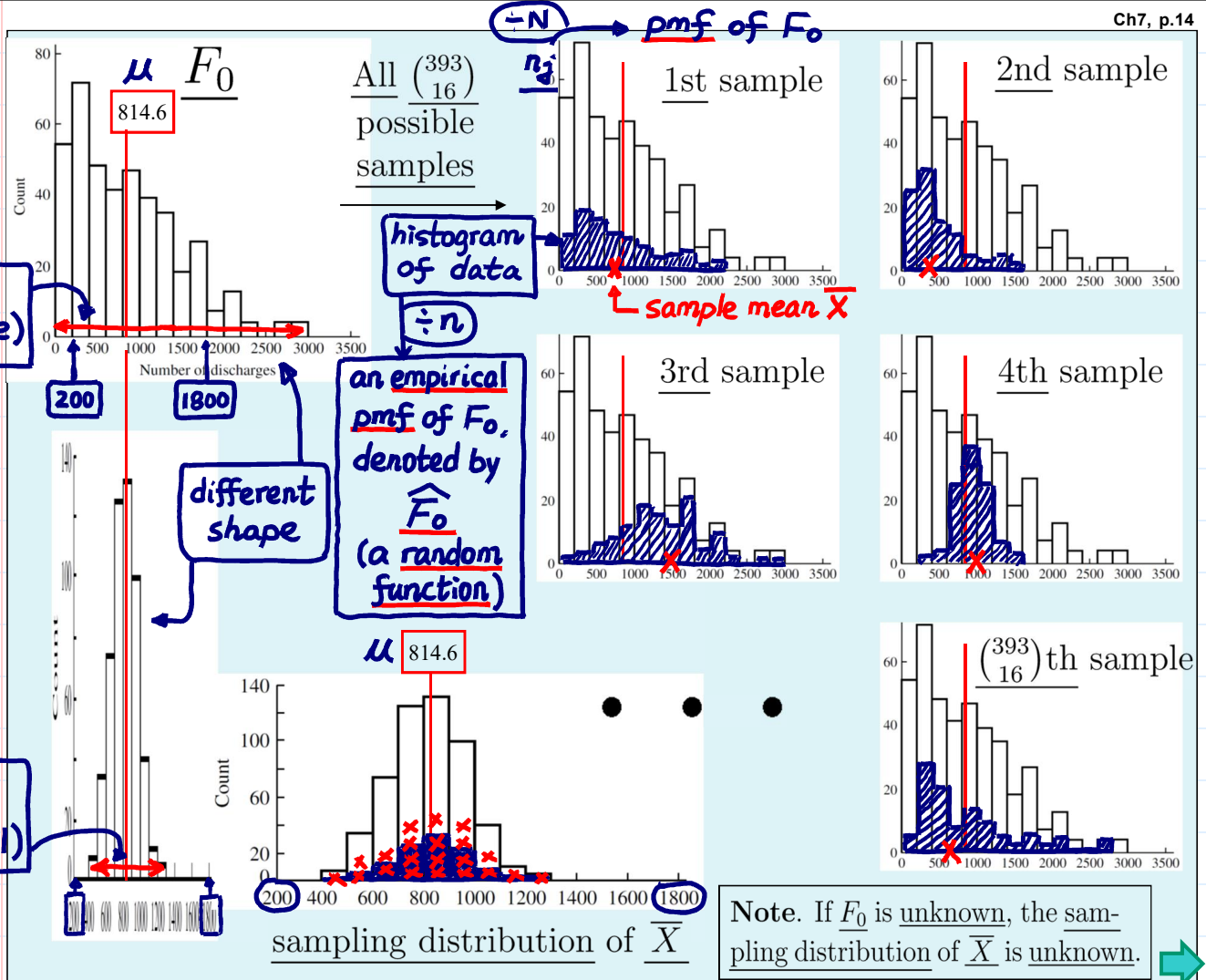
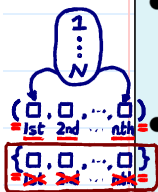
- Suppose we want to know the sampling distribution of \bar{X} of a s.r.s. without replacement of sample size $n=16$.

- There are $\binom{393}{16}$ possible samples. Note that $\binom{393}{16}$ is of order 10^{28} !

- Sampling distribution of \bar{X} is formed by the (sample) mean of each of the possible samples along with their probabilities.

Q: What is the source of randomness in \bar{X} ?

Ans: random sampling



- $\binom{393}{16} = 10^{28}$ is too large
- To reduce computation, we can use the technique of **simulation** to understand the sampling distribution of \bar{X} .

perform an S.R.S on the $\binom{393}{16}$ samples

- randomly draw many (say, 500) s.r.s. of size n
- compute the mean of each sample
- form a histogram of the collection of these sample means

Why?

This histogram will be an approximation to the sampling distribution of \bar{X} .

- Figure 7.2 (textbook) shows the results for sample size $n=8, 16, 32$, or 64 .

Thm 1 LNp.16

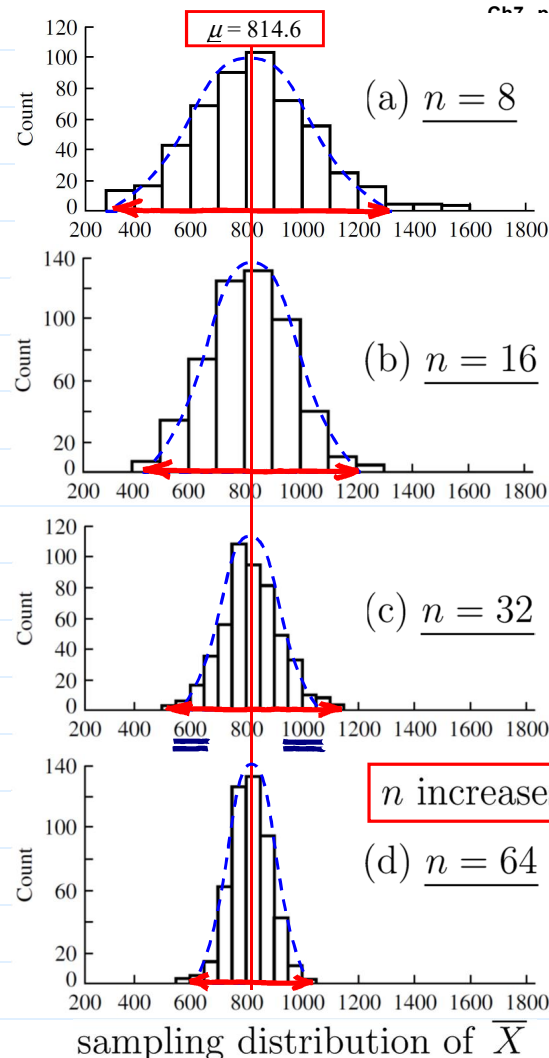
→ All the four histograms are centered at $\mu=814.6$.

Thm 2-3 LNp.17-18

→ As n increases, the histograms become less spread out.

Thm 12-13 LNp.29-30

→ Although shape of F_0 (population distribution) is not symmetric about μ , these histograms are nearly so.



Theorem 1 (expectation of sample mean)

(1) Under simple random sampling, with or without replacement,

The k th observation
Note. $X_k \sim F_0$

$$E(\underline{X}_k) = \underline{\mu} \quad \text{and} \quad \text{Var}(\underline{X}_k) = \underline{\sigma}^2.$$

population variance

population mean

(2) Under simple random sampling, with or without replacement,

$$E(\underline{\bar{X}}) = \underline{\mu}.$$

parameter

So, \bar{X} is an unbiased estimator of μ , i.e., the sampling distribution of \bar{X} is centered at μ . Recall graphs in LNp.15

Proof: Under simple random sampling, no matter with or without replacement, the marginal distribution of \underline{X}_k is F_0 . Thus, we have

LNp.11

$$E(\underline{X}_k) = \sum_{j=1}^m \zeta_j P(\underline{X}_k = \zeta_j) = \sum_{j=1}^m \zeta_j (n_j/N) = \frac{1}{N} \sum_{j=1}^m n_j \zeta_j = \underline{\mu}.$$

$$\text{Var}(\underline{X}_k) = E(\underline{X}_k^2) - [E(\underline{X}_k)]^2 = \frac{1}{N} \left(\sum_{j=1}^m n_j \zeta_j^2 \right) - \underline{\mu}^2 = \underline{\sigma}^2,$$

and

$$E(\underline{\bar{X}}) = E\left(\frac{1}{n} \sum_{k=1}^n \underline{X}_k\right) = \frac{1}{n} \sum_{k=1}^n E(\underline{X}_k) = \frac{1}{n} (n \underline{\mu}) = \underline{\mu}.$$

Definition 3 in LNp.5

Theorem 2 (variance of sample mean, s.r.s. with replacement)

Under simple random sampling with replacement, we have

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n},$$

parameter

and the standard error (st.e.) of \bar{X} , denoted by $\sigma_{\bar{X}}^*$, is σ/\sqrt{n} .

$$n \rightarrow 4n$$

$$\sigma_{\bar{X}}^* \rightarrow \frac{1}{2} \sigma_{\bar{X}}^*$$

Proof: Under simple random sampling with replacement,

we have

 \therefore independent $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} F_0$. \leftarrow LNP. IIThus, $\text{Cov}(X_k, X_l) = 0$ for any $1 \leq k < l \leq n$, and

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{k=1}^n X_k\right) = \frac{1}{n^2} \sum_{k=1}^n \text{Var}(X_k) = \frac{1}{n^2} (n \sigma^2) = \frac{\sigma^2}{n}.$$

~~$+2 \sum_{k,l} \text{Cov}(X_k, X_l)$~~

	n	$\sigma_{\bar{X}}^*$
100	32	-16
400	16	-8
1600	8	-8

Note 5 (Some notes about the st.e. of sample mean, with replacement)① $\sigma_{\bar{X}}^* = \sigma/\sqrt{n}$ (a measure of how spread out \bar{X} is)measures the precision of the estimator \bar{X} .• $\sigma_{\bar{X}}^*$ is determined by n and σ , but not N .• $\sigma_{\bar{X}}^*$ is inversely proportional to \sqrt{n} , i.e., inorder to double the accuracy, n must be quadrupled (the contribution of each observation to the accuracy of \bar{X} decays with the increase of n)larger n ,
more
accurate
estimator① same $n \rightarrow$ same precision \rightarrow irrelevant to N ② No need to sample a certain proportion (i.e., n/N) of the population to reach same precision.**Theorem 3** (variance of sample mean, s.r.s. without replacement)

Under simple random sampling without replacement, we have

What if
 $n=N$?

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

parameter

and the standard error of \bar{X} , denoted by $\sigma_{\bar{X}}^*$, is $(\sigma/\sqrt{n})\sqrt{1 - \frac{n-1}{N-1}}$.

a value between 0 and 1

Proof: First, for $1 \leq k < l \leq n$,

$$\text{Cov}(X_k, X_l) = E(X_k X_l) - E(X_k) E(X_l)$$

$$= \left(\sum_{s=1}^m \sum_{t=1}^m \zeta_s \zeta_t P(X_k = \zeta_s, X_l = \zeta_t) \right) - \mu^2$$

LNP. II

$$= \sum_{s=1}^m \zeta_s^2 \left(\frac{n_s(n_s-1)}{N(N-1)} \right) + \sum_{s=1}^m \sum_{t \neq s}^m \zeta_s \zeta_t \left(\frac{n_s n_t}{N(N-1)} \right) - \mu^2$$

$\zeta_s = \zeta_t$ iff $s=t$

$$= \zeta_s^2 \left(\frac{n_s^2}{N(N-1)} \right)$$

when $\zeta_s = \zeta_t$

$$= \frac{N}{N-1} \sum_{s=1}^m \sum_{t=1}^m \zeta_s \zeta_t \left(\frac{n_s n_t}{N \cdot N} \right) - \frac{1}{N-1} \sum_{s=1}^m \zeta_s^2 \left(\frac{n_s}{N} \right) - \mu^2$$

$$= \frac{N}{N-1} E(X_k) E(X_l) - \frac{1}{N-1} E(X_k^2) - \mu^2$$

$$= \frac{N}{N-1} \cancel{\mu^2} - \frac{1}{N-1} (\sigma^2 + \cancel{\mu^2}) - \cancel{\mu^2} = \frac{-\sigma^2}{N-1}$$

$$\text{Var}(X_k) = E(X_k^2) - [E(X_k)]^2$$

$$\sigma^2 = \mu^2$$

 \leftarrow why negative?

$$E(X_k)$$

$$= \left(\sum_{s=1}^m \zeta_s \frac{n_s}{N} \right) \times$$

$$\left(\sum_{t=1}^m \zeta_t \frac{n_t}{N} \right)$$

$$= E(X_l)$$

When
 $\zeta_s = \zeta_t$ When
 $\zeta_s \neq \zeta_t$

Then,

$$Var(\bar{X}) = Var\left(\frac{1}{n} \sum_{k=1}^n X_k\right)$$

Note. $1 \leq k < l \leq n$ = 0 in with repl.

$$= \frac{1}{n^2} \sum_{k=1}^n Var(X_k) + \frac{2}{n^2} \sum_{k=1}^{n-1} \sum_{l=k+1}^n Cov(X_k, X_l)$$

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

$\frac{\sigma^2}{n} - \sigma^2 \times \frac{n-1}{n} \times \frac{1}{N-1}$
 $\approx \sigma^2/n - \sigma^2/N$

Why

$\sigma_{\bar{X}}^2$ (without)
 $< \sigma_{\bar{X}}^{*2}$ (with)?

Note 6 (Some notes about the st.e. of sample mean, without replacement)

- The variance of \bar{X} in s.r.s. without replacement differs from that in s.r.s. with replacement by the factor $(1 - \frac{n-1}{N-1})$, which is called the

N

finite population correction. (Note. $1 - \frac{n-1}{N-1} \rightarrow 1$ when $N \rightarrow \infty$)

- n/N : sampling fraction ($\approx \frac{n-1}{N-1}$ in most cases) $\Rightarrow 1 - \frac{n-1}{N-1} \approx$ unsampled fraction

- $\sigma_{\bar{X}} \approx \sigma_{\bar{X}}^* = \sigma/\sqrt{n}$ if the sampling fraction is very small (i.e., $n \ll N$).

When $n \ll N$,
 without repl.
 \approx with repl.

- $\sigma_{\bar{X}}$ also depends on n and σ , i.e.,

$$\sigma_{\bar{X}} \downarrow \text{ as } n \uparrow \quad \text{and} \quad \sigma_{\bar{X}} \uparrow \text{ as } \sigma \uparrow,$$

and $\sigma_{\bar{X}}$ depends on N only through the sampling fraction.

Example 4 (st.e. of sample mean, cont. Ex.2 in LNp.4)

- $N = 393$ hospitals. Consider s.r.s. without replacement of size $n = 32$.

- Because $\sigma = 589.7$ (of 393 hospitals), we have

σ : st.d. of F_0 , usually unknown,
need to estimate using data

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \sqrt{1 - \frac{n-1}{N-1}} \approx 100,$$

\times : unknown
 in sampling
 survey

where finite population correction $1 - \frac{31}{392} \approx 0.92$ makes little difference.

- Most of sample means differ from the population mean $\times 814$ by less than $2 \times \sigma_{\bar{X}} = 200$ (see graph (c) of Figure 7.2 in LNp.15).

$$814 \pm 200$$

$$= (614, 1014)$$

$$Z \sim \text{normal}(\mu_Z, \sigma_Z^2)$$

$$P(Z \in \mu_Z \pm 2 \times \sigma_Z)$$

$$\approx 0.95$$

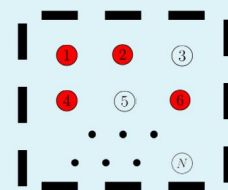
Theorem 4 (mean and variance of sample mean for dichotomous x_i 's)In the dichotomous case, $\bar{X} = \hat{p}$ (sample proportion), and

check
 Note 1
 in
 LNp.6

- under s.r.s. with or without replacement, $E(\hat{p}) = p$

- under s.r.s. with replacement, $Var(\hat{p}) = \frac{p(1-p)}{n}$

variance
 is a
 function
 of mean

and $n\hat{p} = \sum_{k=1}^n X_k$ follows binominal(n, p) distribution

mean = np
 variance = $np(1-p)$

of 1's in
 X_1, \dots, X_n

$\therefore X_1, \dots, X_n$ i.i.d. Bernoulli(p)

- under s.r.s. without replacement, $Var(\hat{p}) = \frac{p(1-p)}{n} \left(1 - \frac{n-1}{N-1}\right)$

and $n\hat{p} = \sum_{k=1}^n X_k$ follows hypergeometric($n, Np, N(1-p)$) distribution

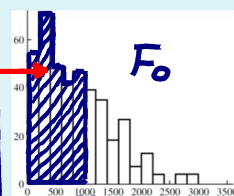
Example 5 (st.e. of sample mean, dichotomous case, cont. Ex.2 in LNp.4)

- In the population of 393 hospitals, a proportion of $\overset{\times}{p} = 0.654$ had fewer than 1000 discharges. parameter \rightarrow

- $y_i = 1$ if $x_i < 1000$ and $y_i = 0$ if $x_i \geq 1000$

Data: X_1, \dots, X_n For $k = 1, \dots, n$, $\underline{Y_k} = \underline{I_{[0, 1000)}(X_k)}$ Indicator functionFor a set A ,

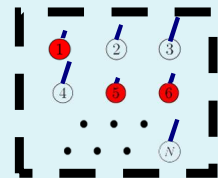
$$\underline{I_A(x)} = \begin{cases} 1, & x \in A \\ 0, & x \notin A \end{cases}$$



- For an s.r.s. without replacement sample Y_1, \dots, Y_n of size $n = 32$, the estimator of p is $\hat{p} = \bar{Y}$ and

 \times : unknown in sampling survey

$$\underline{\sigma_{\hat{p}}} = \sqrt{\frac{p(1-p)}{n}} \sqrt{1 - \frac{n-1}{N-1}} = \sqrt{\frac{.654 \times .346}{32}} \sqrt{1 - \frac{31}{392}} = \underline{0.08}.$$

**Definition 9** (estimator of population total)

Because $\tau = \sum_{i=1}^N x_i$ (population total) equals $N\mu$, known value μ
 an intuitive estimator of τ is $\underline{T} = \underline{N \bar{X}}$. parameter μ

Note. \underline{T} is not $\sum_{k=1}^n X_k = n \bar{X}$. cf. \rightarrow estimate

Theorem 5 (mean of population total estimator)

Under simple random sampling, with or without replacement, we have

$$\underline{E(\underline{T})} = \underline{\tau}. \quad \underline{E(T)} = \underline{E(N\bar{X})} = \underline{NE(\bar{X})} = \underline{N\mu} = \underline{\tau}$$

That is, \underline{T} is an unbiased estimator of $\underline{\tau}$.

Theorem 6 (variance of population total estimator)

- Under simple random sampling with replacement, $\underline{Var(\underline{T})} = \underline{N^2 \left(\frac{\sigma^2}{n} \right)}$. $\leftarrow \underline{Var(\bar{X})}$
- Under simple random sampling without replacement,

9/8

cf. the precision of \bar{X}

• Note 5 in LNp.17

• Note 6 in LNp.19

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

$$\underline{Var(T)} = \underline{Var(N\bar{X})} = \underline{N^2 Var(\bar{X})}$$

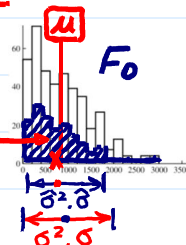
Note. The precision of the estimator \underline{T} does depend on population size \underline{N} .

• Estimation of population variance

$$\sigma^2 = \sum_{j=1}^m \frac{n_j}{N} (c_j - \mu)^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

Recall. When $\underline{F_0}$ is unknown, the $\underline{\sigma}$ in the standard error of $\underline{\bar{X}}$ is a parameter, i.e., it is unknown.

Q: how to estimate $\underline{\sigma}$ or $\underline{\sigma^2}$?

histogram of data**Definition 10** (sample variance)

The sample variance of X_1, X_2, \dots, X_n is defined as $\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})^2$.
 \uparrow • a function of data • a r.v. • an estimator

Theorem 7 (expectation of sample variance, s.r.s. with replacement)

Under s.r.s. with replacement, we have $\underline{E(\hat{\sigma}^2)} = \underline{\sigma^2} \left(\frac{n-1}{n} \right)$ \leftarrow not unbiased