

Jan 9th, 2013

Instructor: Zhou Zhou  
zhou@utstat.toronto.edu

TA: Alexander Shestopaloff alexander@utstat.utoronto.ca

Final 55% + Mid-term 30% + Quiz 15% (3 out of 4)

Today, we are gonna start from Ch. 7 and finish it.

### Chapter 7 Sample Survey

Goal: To make inference about a large population via sampling just a small fraction of it.  
Sample surveys are random in nature.

Should be  
↓  
objective  
↓  
doctors giving medications

#### 1. Population Parameters

Think of a finite population of size  $N$ . We are interested in one characteristic of the population the values of which are ordered as  $x_1, x_2, \dots, x_N$

↓  
Not Random

Population mean  $\frac{1}{N} \sum_{i=1}^N x_i = \mu$

← not random  
← parameter

Population Variance  $\sigma^2 = \frac{1}{N} \sum_{i=1}^N [x_i - \mu]^2$

← parameter

#### Simple Random Sample

(S.R.S) For each sample of size  $n$  in the population the chance of being selected are the same.

\* Same as choosing randomly without replacement.  $X_i$ 's are not independent

The selected sample is denoted by

$$x_1, x_2, \dots, x_n$$

Sample Mean  $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$

Sample variance  $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$

Thm:  $E[x_i] = \mu$  for any  $i$ .

Sampling without replacement is equivalent to grabbing  $n$  items at the same time from an urn of  $N$  items.

$E[x_i] = E[X_1] = E[X_2] = \dots$  (because random labeling)

Now suppose that characteristic we are interested in has  $k$  different values  $\beta_1, \beta_2, \dots, \beta_k$ . Suppose in the population there are  $N_1 \beta_1$ 's,  $N_2 \beta_2$ 's, ...,  $N_k \beta_k$ 's.

$$\begin{aligned} E[X_i] &= \sum_{i=1}^k \beta_i p[X_i = \beta_i] = \sum_{i=1}^k \beta_i \frac{N_i}{N} = \frac{1}{N} \sum_{i=1}^k \beta_i N_i \\ &= \frac{1}{N} \sum_{i=1}^N X_i \\ &= \mu \end{aligned}$$

**Corollary**  $E[\bar{X}] = \mu$

Proof:  $E[\bar{X}] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \mu$   
 $\bar{X}$  is an unbiased estimate of  $\mu$ .

### Quality of Estimators

Commonly used measure of quality is the mean square error: (MSE)

$$\boxed{\text{MSE} = E[\hat{\theta} - \theta]^2}$$
 small MSE means good (accurate) estimate

$$\begin{aligned} &= E[\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta]^2 \\ &= E[\hat{\theta} - E(\hat{\theta})]^2 + [E[\hat{\theta}] - \theta]^2 + 2E[(\hat{\theta} - E[\hat{\theta}]) \times \underbrace{(E[\hat{\theta}] - \theta)}_{\text{not random}}] \\ &= \underbrace{E[\hat{\theta} - E(\hat{\theta})]^2}_{\text{variance of estimator}} + \underbrace{[E[\hat{\theta}] - \theta]^2}_{\text{bias of estimator}} \quad \text{two non-negative parts} \end{aligned}$$



B is better because of smaller variance

$$\begin{aligned} \text{Var}[X_i] &= E[X_i^2] - (E[X_i])^2 \\ &= E[X_i^2] - (\mu)^2 \\ &= \sum_{i=1}^k \beta_i^2 p[X_i = \beta_i] - \mu^2 \\ &= \sum_{i=1}^k \beta_i^2 \frac{N_i}{N} - \mu^2 = \frac{1}{N} \sum_{i=1}^k \beta_i^2 N_i - \mu^2 \end{aligned}$$

For each  $i$ ,  $\beta_i^2 N_i$  = total sum of squares in class,  $\Rightarrow \sum_{i=1}^k \beta_i^2 N_i = \sum_{i=1}^N X_i^2$

$$\text{Hence } \text{Var}[X_i] = \frac{1}{N} \sum_{i=1}^N X_i^2 - \mu^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N} = \sigma^2$$

But we don't want  $\text{Var}[X_i]$  but  $\text{Var}[\bar{X}]$

$$\text{Var}[\bar{X}] = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \left( \sum_{i,j} 2\text{cov}[X_i, X_j] + \underbrace{\sum_{i=1}^n \text{Var}(X_i)}_{\text{! ! !}} \right) \neq n \sigma^2$$

Because of exchangeability,  $\text{cov}[X_i, X_j] = \text{cov}[X_1, X_i]$  for  $\forall i \neq j$

$$\begin{aligned}\text{cov}[X_1, X_2] &= E[X_1 X_2] - E[X_1]E[X_2] \\ &= E[X_1 X_2] - \mu^2\end{aligned}$$

$$E[X_1 X_2] = E[E[X_1 X_2 | X_2]]$$

$$* E[X] = E[E[X | i]]$$

$$E[X_1 X_2 | X_2 = \beta_t] = \beta_t E[X_1 | X_2 = \beta_t]$$

$$\text{For } r \neq t, P[X_1 = \beta_r | X_2 = \beta_t] = \frac{N_r}{N-1}$$

$$P[X_1 = \beta_t | X_2 = \beta_t] = \frac{N_{t-1}}{N-1}$$

$$\begin{aligned}E[X_1 | X_2 = \beta_t] &= \beta_t \frac{N_{t-1}}{N-1} + \sum_{i \neq t} \beta_i \frac{N_i}{N-1} = \frac{1}{N-1} \sum_{i=1}^k \beta_i N_i - \frac{\beta_t}{N-1} \\ &= \frac{1}{N-1} \mu - \frac{\beta_t}{N-1}\end{aligned}$$

$$E[X_1 X_2 | X_2 = \beta_t] = \frac{\beta_t [\mu - \beta_t]}{N-1}$$

$$E[X_1 X_2] = \sum_{t=1}^k \frac{\beta_t [N\mu - \beta_t]}{N-1} \quad P[X_2 = \beta_t] = \sum_{t=1}^k \frac{\beta_t [N\mu - \beta_t]}{N-1} \cdot \frac{N_t}{N}$$

$$\begin{aligned}&= \frac{1}{(N-1)N} \left[ \sum_{t=1}^k (N\mu) \beta_t N_t - \sum_{t=1}^k \beta_t^2 N_t \right] \\ &= \frac{1}{(N-1)N} [N^2 \mu^2 - \sum_{t=1}^k \chi_t^2]\end{aligned}$$

$$\text{known: } \sigma^2 = \frac{1}{N} \sum_{i=1}^N \chi_i^2 - \mu^2$$

$$\Rightarrow \sum_{i=1}^N \chi_i^2 = (\sigma^2 + \mu^2)N$$

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

$$= \mu^2 - \frac{1}{N-1} \sigma^2$$

$$\text{cov}[X_i, X_j] = \mu^2 - \frac{1}{N-1} \sigma^2 - \mu^2$$

$$= -\frac{1}{N-1} \sigma^2$$

*final answer*

Note  $\nabla \bar{X} = \frac{1}{n^2} [2(-\frac{1}{N-1} \sigma^2) n(n-1) + n \sigma^2]$

$$= \frac{1}{n} [2(-\frac{1}{N-1}) \frac{n-1}{2} + 1] \sigma^2$$

$$= \frac{1}{n} (1 - \frac{n-1}{N-1}) \sigma^2 \quad \leftarrow \text{MSE of } \bar{X}$$

Frequently  $n \ll N$   
 $\text{MSE of } \bar{X} \approx \frac{\sigma^2}{n}$

to make  $\bar{X}$  more accurate :

① increase  $n$

② or decrease  $\sigma^2$ .

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

$$\begin{aligned}
E[\hat{\sigma}^2] &= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \frac{1}{n} \sum_{i=1}^n E[X_i^2] - E[\bar{X}^2] = [\sigma^2 + \mu^2] - [Var[\bar{X}] + E[\bar{X}]^2] \\
&= \sigma^2 + \mu^2 - \left[\frac{1}{n}(1 - \frac{n-1}{N-1})\sigma^2 + \mu^2\right] \\
&= \sigma^2 - \frac{1}{n}(1 - \frac{n-1}{N-1})\sigma^2 \\
&= \sigma^2 \left[1 - \frac{1}{n} + \frac{n-1}{(N-1)n}\right]
\end{aligned}$$

expectation  
of sample  
variance

If  $N \gg n$ ,  $E[\hat{\sigma}^2] \approx \sigma^2(1 - \frac{1}{n})$

Define  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  If  $n \ll N$ ,  $E[S^2] \approx \sigma^2$

$$E[\hat{\sigma}^2] \approx S^2 \left(\frac{N}{N-1}\right)$$