


Statistics and Sampling Distributions

- Statistics are the numerical descriptive measures that you calculate from the random sample of a population (eg. sample mean, sample variance)
- these statistics can vary or change for each different random sample \Rightarrow they are random variables
- the probability distribution of these statistics, of this random variable, is called sampling distributions

The **sampling distribution of a statistic** is the probability distribution for the possible values of the statistic that results when random samples of size n are repeatedly drawn from the population.

- 3 ways of finding sampling distribution
 - ① mathematically from laws of probability
 - ② Simulation
calc a lot of stats from a lot of samples and plot the relative frequency histogram
 - ③ Use statistical theorem

= relative frequency histogram
size $\rightarrow 0$

bell curve's
importance

↖ The central limit theorem v. imp
+ Sample mean

- it is a statistical theorem that describes the sampling distribution of statistics that are sums or averages
- theorem :
sums and means of random samples of measurements drawn from a population tend to have an approximately normal distribution
- as n increases, the distribution of \bar{x} becomes more and more normal + skew decreases
- The CLT can be restated to apply to the sum of sample measurements $\sum x_i$; \rightarrow as n becomes large has an approximately normal distribution with mean $n\mu$ and standard deviation $\sigma\sqrt{n}$
- to apply to the mean of sample measurements \bar{x} , as $n \rightarrow$ large, has an approximately normal distribution with mean μ and standard deviation σ/\sqrt{n}

WHEN TO USE CLT

- if sampled population is normal, sampling distribution of \bar{x} will be normal $\forall n$
- if sampled population is approximately symmetric, CLT holds for relatively small n
- if sampled population is skewed, for $n \geq 30$, \bar{x} will be approximately normal

Standard Error of Sample Mean

· standard deviation of a statistic = standard error of the estimator (SE)

· \therefore standard deviation of $\bar{x} = \frac{\sigma}{\sqrt{n}}$
= standard error of the mean ($SE(\bar{x})$ SEM)

· finding probabilities for the sample mean \bar{x}
↑
assuming $\bar{x} \approx$ normal

① $\mu, \bar{x}, \sigma \leftarrow \text{calc}$
 $SE(\bar{x}) = \frac{\sigma}{\sqrt{n}} \leftarrow \text{calc}$

② calc z value $z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$

③ use normal table to compute the probability of \bar{x}

· note: population mean should be the peak of the sample mean's normal distribution

greater the n (sample size), less the deviation for mean $\therefore \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

$np \geq 10, n(1-p) \geq 10$ then normal

mean and std of the normal dist of sample mean and σ

Properties of sampling distribution of sample proportion

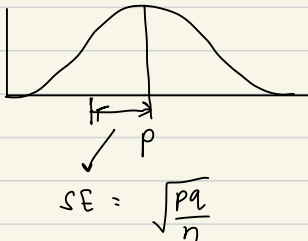
binomial population : success or failure
 ↓
 's parameter 'p'
 p = probability of success for whole population
 ↓
 \hat{p} : $\frac{\text{no of success}}{\text{no of trials}}$
 ↓
 prob success in sample
 sd = $\sigma = \sqrt{pq}$
 $\mu = p$ } of random variable $Y = \begin{cases} a & \text{success} \\ b & \text{failure} \end{cases}$

when plotting sampling distribution of \hat{p}

μ
 σ/\sqrt{n} mean = p
 SEM = $\sqrt{\frac{pq}{n}}$ } because its of \hat{p}
 mean = $\frac{\text{mean of sample}}{n} = np/n$
 sd = $\frac{\text{sd sample}}{n} = \frac{\sqrt{npq}}{n}$
 ↖ no of samples / observations

can be approximated for $np > 5$ AND $nq > 5$

calculate probabilities for \hat{p} the same way.



calculating skew :

less than 5 success → right
 less than 5 failure → left
 at least 5 s + f → normal
 neither → uniform

8.

Point Estimation

point estimator \rightarrow a statistic used to estimate the value (pe) of an unknown parameter of a population

practically, we can have many pe's \leftarrow how to pick best?

characteristics of pe.

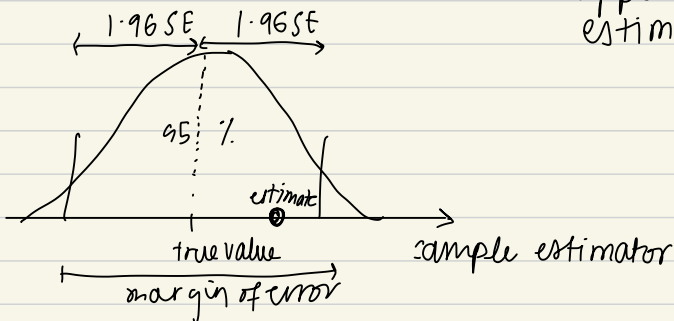
① sampling dist of pe should be centered over true value of parameter ($\mu_{pe} - \mu_p$)

\uparrow
should be unbiased

② variance of pe's sampling dist should be as small as possible

③ error of estimation : difference b/w an estimate & true value of the parameter measured in SE of estimator's sampling dist

margin of error = 95% margin of error = practical upper bound for error of estimation.



est.

pop

mean, $n \geq 30$, margin of error = $\pm 1.96 \left(\frac{s}{\sqrt{n}} \right)$

sample
s.d
↓

est.
population
proportion, $n\hat{p} > 5$ & $n\hat{q} > 5$, margin of error = $\pm 1.96 \left(\sqrt{\frac{\hat{p}\hat{q}}{n}} \right)$

eg.

sample size $n = 50$

sample mean = 980 \bar{x}

sample dev = 105 s

pop mean?

$$1.96 SE = 1.96 \left(\frac{s}{\sqrt{n}} \right) \\ = 1.96 \left(\frac{105}{\sqrt{50}} \right)$$

$$\leftarrow = 29.10 \approx 29 \text{ pounds}$$

sample estimate 980 is within
 ± 29 of population mean

eg.

$n = 100$

$\hat{p} = 0.73$

$\sqrt{\hat{p}\hat{q}} = 0.44$

estimate p

$$1.96 \left(\frac{\sqrt{\hat{p}\hat{q}}}{\sqrt{n}} \right) = 0.09.$$

$$0.73 \text{ is with } p \pm 0.09 \\ \Rightarrow p = 0.73 \pm 0.09$$

• ————— •
true

interval estimation

- interval estimator : rule for calculating 2 numbers (a,b) which contains the parameter of interest might
- probability that it will = confidence coefficient = $1 - \alpha$ → contain the parameter
- constructing a confidence interval for a sampling distribution of a point estimator
- 95% \rightsquigarrow $PE \pm 1.96 SE$
↙
variable center

confidence limits for confidence coeff $(1 - \alpha)$

↙ $PE \pm Z_{\alpha/2} (SE)$
↓
↗ upper and lower confidence limits

random
quantities

sample sd
use s ← if we don't know σ → $\frac{\sigma}{\sqrt{n}}$

good interval is as narrow as possible, w/ large $(1 - \alpha)$

but increasing confidence w/o increasing width can only happen by changing sample size n

for population proportion p , \hat{p} ← sample proportion is the best estimator

↓

$$\hat{p} \pm Z_{\alpha/2} \left(\sqrt{\frac{pq}{n}} \right)$$

assuming $np > 5$ & $nq > 5$
& independent, p constant

population mean : $\mu = \sum x_i p_i$

population sd : σ

sample size : n

there are ${}^N C_n$ possible samples

population size : N

CLT applies to random variable x_i

$n\mu, \sigma\sqrt{n} \leftarrow$ sum of x_i (which is also a random var)

$\mu, \sigma/\sqrt{n} \leftarrow$ mean : $\sum x_i p_i$ (" ")

pop normal \Rightarrow sampling dist normal $\forall n$

pop approx symmetric \Rightarrow for small n

pop is skewed $\Rightarrow n > 30$.

standard error = standard deviation of a statistic = SEM

$SE(\text{mean}) = \frac{\sigma}{\sqrt{n}}$ mean : μ

$SE(\text{sum}) = \sigma \cdot \sqrt{n}$ mean : $n\mu$

calculating probabilities : very simple \rightarrow you know what to do.

standard error for sampling proportion

$\mu = p, \sigma = \sqrt{\frac{pq}{n}}$

p = population proportion

$np > 5, nq > 5 \leftarrow$ can apply CLT

less than 5 successes \rightarrow right

less than 5 failure \rightarrow left

at least 5 s + f \rightarrow normal

neither \rightarrow uniform

population parameter estimation :

· mean : $\bar{x} \pm z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$ $\xrightarrow{\text{can be } s}$

· proportion : $\hat{p} \pm z_{\alpha/2} \left(\sqrt{\frac{\hat{p}\hat{q}}{n}} \right)$

other key points

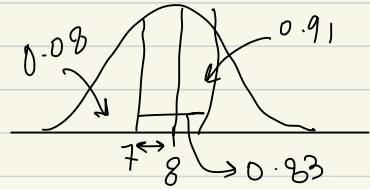
· $E(x^2) = \mu_x^2 + \sigma_x^2$

· $\text{var}(\sum X_i) = \sum (\text{var}(X_i)) \rightsquigarrow \text{var}(\text{sum of } x_i) = n \cdot \sigma^2$
 $\therefore \text{sd} = \sqrt{n} \cdot \sigma$

some sums

pg 30 ①

$\mu = 8$ pop mean
 $\sigma = 4$ pop sd
 $n = 30$ sample size
↑ enough to approx
∴ it is skewed.



records avg duration.

mid = 8

$$SEM = 4/\sqrt{30} < 1.$$

KA lesson 2 1. $\mu = 210$

2. $\mu = 30$
 $\sigma = 1.5$

$n = 3$

mean : $\mu = 30$

$$\sigma = 1.5/\sqrt{3}$$

3. $\mu = 60$

$\sigma = 0.5$

$n = 36$

$$0.5/6 =$$

pg 42 7.10

$N = 500$

$$p = 0.6 \quad \sigma = \sqrt{0.6 \times 0.4} = \sqrt{0.24}$$

sd :

$$\mu_{\hat{p}} = 0.6$$

$$\sigma_{\hat{p}} = \frac{\sqrt{0.24}}{\sqrt{500}} = 0.022$$