

Mathematical Statistics v for Comptly

7 Limit Theorems

Motivation

For some new statistics, we may want to derive features of the distribution of the statistic.

When we can't do this analytically, we need to use statistical computing methods to *approximate* them.

We will return to some basic theory to motivate and evaluate the computational methods to follow.

7.1 Laws of Large Numbers

Limit theorems describe the behavior of sequences of random variables as the sample size increases ($n \rightarrow \infty$).

If X_1, \dots, X_n i.i.d

① What is the distribution of $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$?

② How big does n have to be for $\bar{X} \sim \text{Normal}$?

Often we describe these limits in terms of how close the sequence is to the truth.

How far is \bar{X} from μ ?
statistic true value we are estimating

How do we measure this distance? ex. $|\bar{X} - \mu|$ or $(\bar{X} - \mu)^2$ maybe?

We can evaluate this distance in several ways.

Some modes of convergence - e.g.

- almost surely ($P(\lim_{n \rightarrow \infty} X_n = x) = 1$)

- in probability ($\forall \varepsilon > 0, \lim_{n \rightarrow \infty} P(|X_n - x| > \varepsilon) = 0$)

- in distribution ($\lim_{n \rightarrow \infty} F_{X_n}(x) = F_x(x)$)

What happens to sequences of r.v.'s as n gets large.
(gives us useful approximations!)

e.g. Laws of large numbers -

Weak LLN: Sample mean \bar{X}_n converges in probability to pop. mean μ .

Strong LLN: Sample mean \bar{X}_n converges a.s. to pop mean μ

7.2 Central Limit Theorem

Theorem 7.1 (Central Limit Theorem (CLT)) Let X_1, \dots, X_n be a random sample from a distribution with mean μ and finite variance $\sigma^2 > 0$, then the limiting distribution of

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \text{ is } N(0, 1). \quad [\text{converges in distribution}]$$

i.e. $\bar{X}_n \xrightarrow{d} X$ where $X \sim N(\mu, \sigma^2/n)$.

Interpretation:

The sampling distribution of the sample mean approaches a normal distribution as the sample size increases.

Remember

Note that the CLT doesn't require the population distribution to be Normal.

8 Estimates and Estimators

Let X_1, \dots, X_n be a random sample from a population.

Let $T_n = T(X_1, \dots, X_n)$ be a function of the sample.

Then T_n is a "statistic"

and the pdf of T_n is called the "sampling distribution of T_n "

Statistics estimate parameters.

from sample from population

Example 8.1

\bar{X}_n estimates μ

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X}_n)^2 \text{ estimates } \sigma^2$$

$$s = \sqrt{s^2} \text{ estimates } \sigma$$

Definition 8.1 An estimator is a rule for calculating an estimate of a given quantity.

Definition 8.2 An estimate is the result of applying an estimator to observed data samples in order to estimate a given quantity.

A statistic is a point estimator

A CI is an interval estimator

(if based on observed data, they are estimates)

We need to be careful not to confuse the above ideas:

\bar{X}_n function of R.V.'s \rightarrow estimator (statistic)

\bar{x}_n function of observed data (an actual #) \rightarrow estimate (sample statistic)

μ fixed but unknown quantity \rightarrow parameter

We can make any number of estimators to estimate a given quantity. How do we know the "best" one?

What are some properties we can use to say an estimator is "better" than another one?

9 Evaluating Estimators

There are many ways we can describe how good or bad (evaluate) an estimator is.

9.1 Bias

Definition 9.1 Let X_1, \dots, X_n be a random sample from a population, θ a parameter of interest, and $\hat{\theta}_n = T(X_1, \dots, X_n)$ an estimator. Then the *bias* of $\hat{\theta}_n$ is defined as

$$\text{bias}(\hat{\theta}_n) = E[\hat{\theta}_n] - \theta. \quad \begin{matrix} \leftarrow E(T(x_1, \dots, x_n)) = \int_T(x_1, \dots, x_n) f_x(x) dx \\ \text{joint d}\sigma_{x_1, \dots, x_n} \end{matrix}$$

Definition 9.2 An *unbiased estimator* is defined to be an estimator $\hat{\theta}_n = T(X_1, \dots, X_n)$ where

$$\text{bias}(\hat{\theta}_n) = 0, \text{ i.e. } E[\hat{\theta}_n] = \theta.$$

Example 9.1

If you used $\text{Unif}(0,1)$ as your envelope for the Rayleigh dsn, your histogram of values would be biased (too many small values, no large values).

Example 9.2 Let X_1, \dots, X_n be a random sample from a population w/ mean μ and variance $\sigma^2 < \infty$.

$$E(\bar{X}) = E\left(\frac{1}{n} \sum X_i\right) = \frac{1}{n} \sum E(X_i) = \frac{1}{n} \cdot n\mu$$

$$\Rightarrow \text{bias}(\bar{X}) = E(\bar{X}) - \mu = 0 \Rightarrow \text{unbiased estimator for } \mu.$$

Example 9.3 Compare 2 estimators of σ^2 fr. Ex. 9.2:

Sample variance

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

MLE estimate of variance

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

Can show $E s^2 = \sigma^2$ but $\hat{\sigma}^2 = \frac{n-1}{n} s^2$, so

$$E(\hat{\sigma}^2) = \frac{n-1}{n} E s^2 = \frac{n-1}{n} \sigma^2 \Rightarrow \hat{\sigma}^2 \text{ is biased estimator.}$$

Note for large n , $s^2 \approx \hat{\sigma}^2$.

9.2 Mean Squared Error (MSE)

Definition 9.3 The *mean squared error (MSE)* of an estimator $\hat{\theta}_n$ for parameter θ is defined as

$$\begin{aligned} MSE(\hat{\theta}_n) &= E[(\theta - \hat{\theta}_n)^2] \\ &= Var(\hat{\theta}_n) + (bias(\hat{\theta}_n))^2. \end{aligned}$$

can show

Generally, we want estimators with

- ① small bias
 - ② small variance
- often there is a bias-variance trade-off
(can't get both)

Sometimes an unbiased estimator $\hat{\theta}_n$ can have a larger variance than a biased estimator $\tilde{\theta}_n$.

Example 9.4 Let's compare two estimators of σ^2 .

$$s^2 = \frac{1}{n-1} \sum (X_i - \bar{X}_n)^2 \quad \hat{\sigma}^2 = \frac{1}{n} \sum (X_i - \bar{X}_n)^2$$

$$E(s^2) = \sigma^2 \quad E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2$$

but $Var(s^2) > Var(\hat{\sigma}^2)$!

Can show:

$$MSE(s^2) = E(s^2 - \sigma^2)^2 = \frac{2}{n-1} \sigma^4$$

$$MSE(\hat{\sigma}^2) = E(\hat{\sigma}^2 - \sigma^2)^2 = \frac{2n-1}{n^2} \sigma^4$$

$$\Rightarrow MSE(s^2) > MSE(\hat{\sigma}^2).$$

See pg. 331 of Casella & Berger

9.3 Standard Error

Definition 9.4 The *standard error* of an estimator $\hat{\theta}_n$ of θ is defined as

$$se(\hat{\theta}_n) = \sqrt{Var(\hat{\theta}_n)}.$$

← standard error =
 st. dev. of sampling dsn
 of $\hat{\theta}_n$.

We seek estimators with small $se(\hat{\theta}_n)$.

Example 9.5

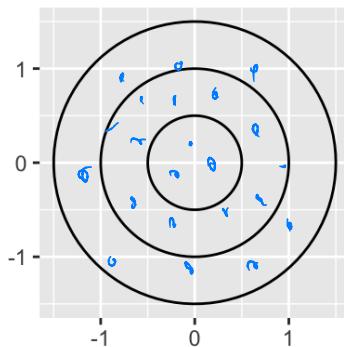
$$se(\bar{X}) = \sqrt{Var(\bar{X}_n)} = \sqrt{\frac{Var X}{n}} = \frac{\sigma_x}{\sqrt{n}}$$

10 Comparing Estimators

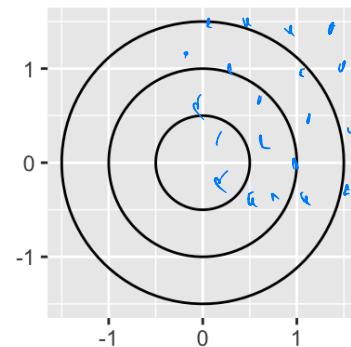
We typically compare statistical estimators based on the following basic properties:

1. *Consistency*: as $n \uparrow$ does estimator converge to parameter it's estimating?
(converges in probability)
2. *Bias*: Is the estimator unbiased? $E(\hat{\theta}_n) = \theta$
3. *Efficiency*: $\hat{\theta}_n$ is more efficient than $\tilde{\theta}_n$ if
 $\text{Var}(\hat{\theta}_n) < \text{Var}(\tilde{\theta}_n)$
4. *MSE*: Compare $\text{MSE}(\hat{\theta}_n)$ to $\text{MSE}(\tilde{\theta}_n)$ but remember the
bias/variance trade-off, $\text{MSE}(\hat{\theta}_n) = \text{Var}(\hat{\theta}_n) + [\text{Bias}(\hat{\theta}_n)]^2$

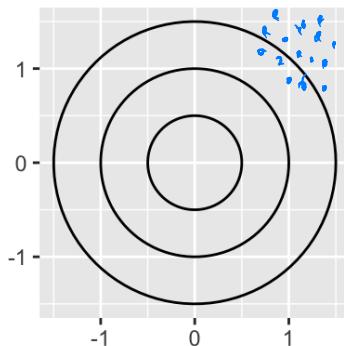
Unbiased and Inefficient



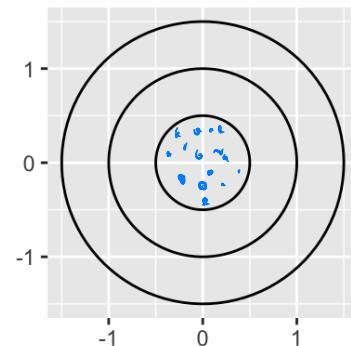
Biased and Inefficient



Biased and Efficient



Unbiased and Efficient



Example 10.1 Let us consider the efficiency of estimates of the center of a distribution. A **measure of central tendency** estimates the central or typical value for a probability distribution.

Mean and median are two measures of central tendency. They are both **unbiased**, which is more efficient?

↪ i.e. which has smaller variance?

```
set.seed(400)
```

```
times <- 10000 # number of times to make a sample
n <- 100 # size of the sample
uniform_results <- data.frame(mean = numeric(times), median =
  numeric(times))
normal_results <- data.frame(mean = numeric(times), median =
  numeric(times))

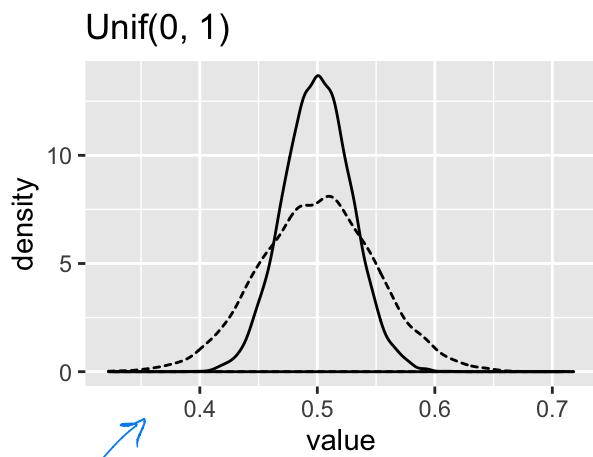
for(i in 1:times) {
  x <- runif(n) ← draw a uniform sample
  y <- rnorm(n) ← draw a normal sample
  uniform_results[i, "mean"] <- mean(x) ← store mean
  uniform_results[i, "median"] <- median(x) ← store median
  normal_results[i, "mean"] <- mean(y)
  normal_results[i, "median"] <- median(y)
}

uniform_results %>%
  gather(statistic, value, everything()) %>%
  ggplot() +
  geom_density(aes(value, lty = statistic)) +
  ggtitle("Unif(0, 1)") +
  theme(legend.position = "bottom") }
```

estimate the density sample →

plot results

```
normal_results %>%
  gather(statistic, value, everything()) %>%
  ggplot() +
  geom_density(aes(value, lty = statistic)) +
  ggtitle("Normal(0, 1)") +
  theme(legend.position = "bottom") }
```



statistic mean median

sampling dsn of

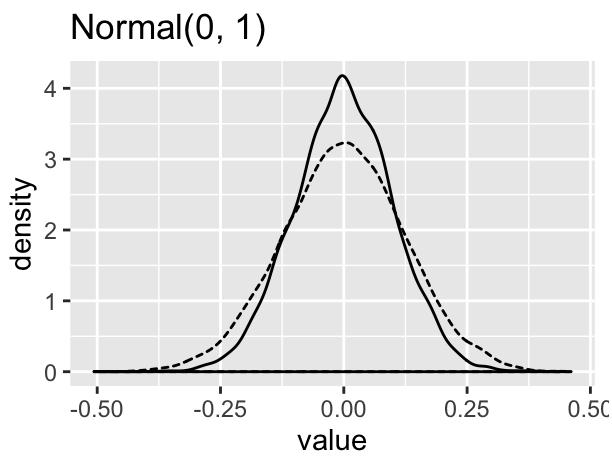
	\bar{X}	S
mean	.4999	.029
median	.4999	.0494

$$\begin{aligned} \text{true mean} &= \\ \text{true median} &= 0.5 \end{aligned}$$

For both $\text{Unif}(0,1)$ and $N(0,1)$,
BIAS: Both mean and median unbiased

EFFICIENCY: mean is more efficient $\hat{\text{Var}}(\text{mean}(x_1, \dots, x_{100})) < \hat{\text{Var}}(\text{median}(x_1, \dots, x_{100}))$

Next Up In Ch. 5, we'll look at a method that produces unbiased estimators of $E(g(X))$!



statistic mean median

sampling dsn of

	\bar{X}	S
mean	0.0001	0.10
median	-0.0009	0.12

$$\begin{aligned} \text{true mean} &= \\ = \text{true median} & \end{aligned}$$

$$= 0$$

NOTE: this is not the case
for all dsns. When a dsn

is heavy tailed, median
is more efficient than the mean.