

Central Limit Theorem

The Most Important Result in Statistics

Functions of RVs

Linear Combinations

Sample Mean

CLT $n \geq 25$

Central limit theorem

One of the most frequently used statistics in engineering applications is the sample mean. We can relate the mean and variance of the probability distribution of the sample mean to those of a single observation when an iid model is appropriate.

In the case of the sample mean, if the sample size (n) is large enough, we can also approximate the shape of the probability distribution function of the sample mean!

Functions of RVs

Linear Combinations

Sample Mean

CLT

(CLT)

Central limit theorem

If X_1, \dots, X_n are **independent** and **identically** distributed (iid) random variable (with mean μ and variance σ^2), then for large n , the variable \bar{X} is approximately normally distributed. That is,

Sample mean $\bar{X} \sim \text{Normal} \left(\mu, \frac{\sigma^2}{n} \right)$
approximately

* This is one of the **most important** results in statistics.

Functions of RVs

Linear Combinations

Sample Mean

CLT

W_i discrete r.v.
 $i=1, 2$

Example: [Tool serial numbers]

Consider selecting the last digit of randomly selected serial numbers of pneumatic tools. Let

- * W_1 = the last digit of the serial number observed next Monday at 9am
- * W_2 = the last digit of the serial number observed the following Monday at 9am

A plausible model for the pair of random variables W_1, W_2 is that they are independent, each with the marginal probability function

$$f(w) = \begin{cases} .1 & w = 0, 1, 2, \dots, 9 \\ 0 & \text{otherwise} \end{cases}$$

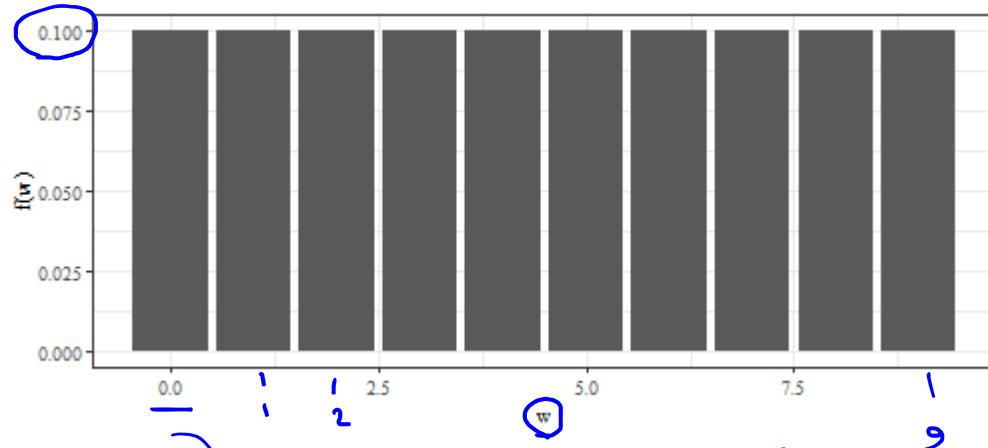
Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Tool serial numbers]



With $\underline{EW} = 4.5$ and $\underline{\text{Var}}W = 8.25$.

Using such a distribution, it is possible to see that

→ $\underline{\overline{W}} = \frac{1}{2}(\underline{W_1} + \underline{W_2})$ has probability distribution

$\begin{matrix} w_1 & \leftarrow w_2 \\ (0,0) & \rightarrow \end{matrix}$

\overline{w}	$f(\overline{w})$								
0.00	0.01	2.00	0.05	4.00	0.09	6.00	0.07	8	0.03
0.50	0.02	2.50	0.06	4.50	0.10	6.50	0.06	8.5	0.02
1.00	0.03	3.00	0.07	5.00	0.09	7.00	0.05	9	0.01
1.50	0.04	3.50	0.08	5.50	0.08	7.50	0.04		

$(w_1=1, w_2=2)$

$(\overline{w}_1=9, \overline{w}_2=9)$

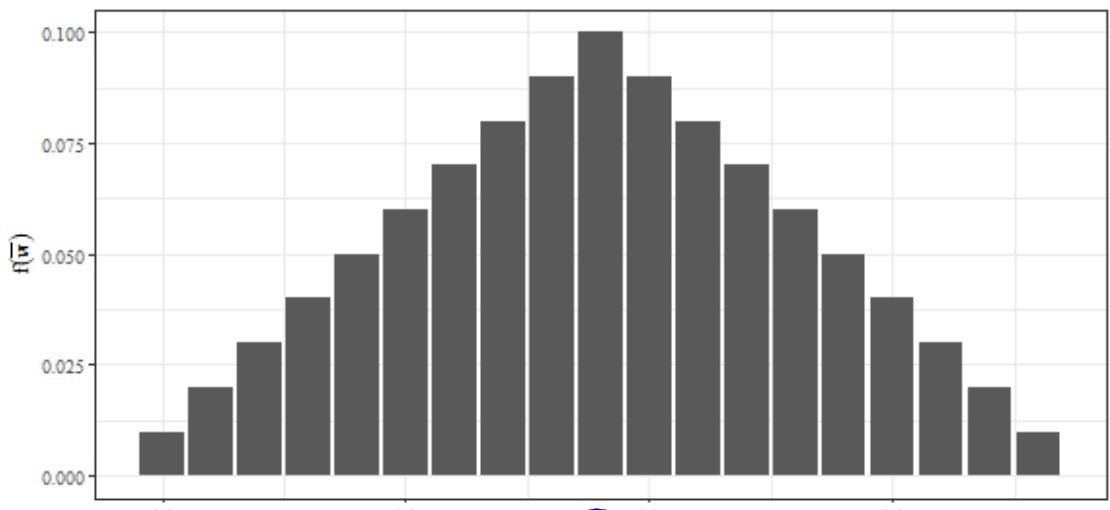
Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Tool serial numbers]



$$E \bar{w} = E w_1 = 4.5$$

$$\text{Var}(\bar{w}) = \frac{\text{Var}(w_1)}{n} = 4.125$$

Comparing the two distributions, it is clear that even for a completely flat/uniform distribution of W and a small sample size of $n = 2$, the probability distribution of \bar{w} looks more bell-shaped than the underlying distribution.

Functions of RVs

Linear Combinations

Sample Mean

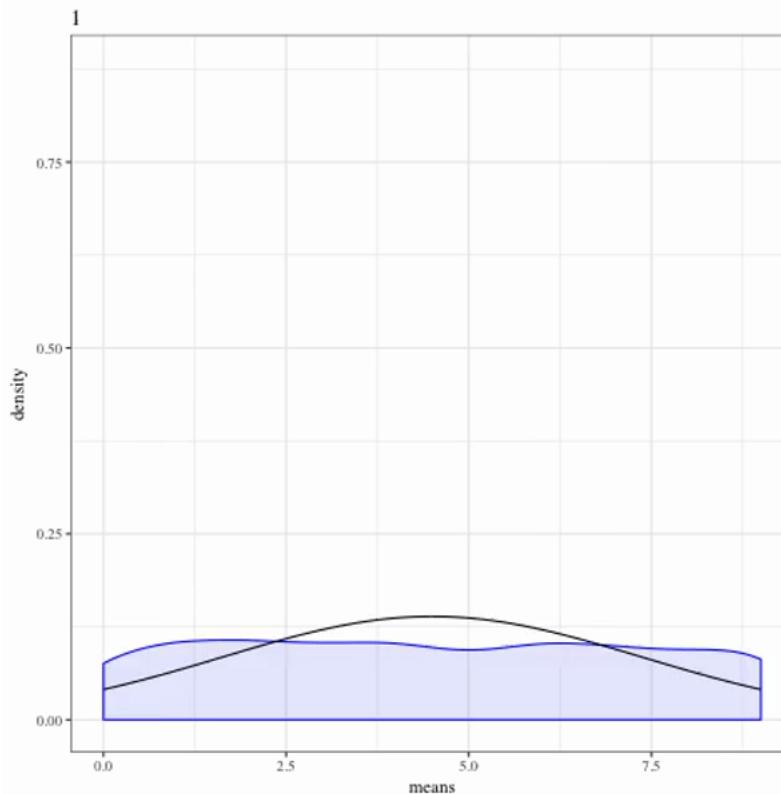
CLT

\bar{w} will always
have $E\bar{w} = Ew_i = 4.5$,
but $\text{var } \bar{w} = \frac{\text{var}(w_i)}{n}$
decreases as $n \rightarrow \infty$

and \bar{w} will approximately have Normal shape
(distribution)

Now consider larger and larger sample sizes,
 $n = 1, \dots, 40$:

Watch how CLT works [here](#)



Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Stamp sale time]

Imagine you are a stamp salesperson (on eBay). Consider the time required to complete a stamp sale as S , and let

assume iid

\bar{S} = the sample mean time required to complete the next 100 sales

$n \geq 25$

Each individual sale time should have an $Exp(\alpha = 16.5s)$ distribution. We want to consider approximating $P[\bar{S} > 17]$.

$$S_i \stackrel{iid}{\sim} Exp(\alpha = 16.5) \rightarrow E S_i = \alpha = 16.5$$

$$\text{var } S_i = \alpha^2 = 16.5^2 = 272.25$$

one of them

↓

$$\text{Now: } E \bar{S} = E S_i = 16.5$$

$$\text{var}(\bar{S}) = \frac{\text{var}(S_i)}{n} = \frac{272.25}{100} = 2.72225$$

Since $n=100 \geq 25$,

$$\xrightarrow{1.65^2}$$

using CLT: $\bar{S} \sim N(\mu=16.5, \sigma^2 = 2.72225)$

$$\rightarrow P(\bar{S} > 17) = P\left(\frac{\bar{S} - 16.5}{\sqrt{2.72225}} > \frac{17 - 16.5}{\sqrt{2.72225}}\right)$$

$$= P(Z > \frac{17 - 16.5}{1.65})$$

0.303

$$= 1 - P(Z \leq 0.303)$$

$$= 1 - \Phi(0.303)$$

$$\text{table } \approx 1 - 0.6217 = 0.3783$$

Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Cars]

Suppose a bunch of cars pass through certain stretch of road. Whenever a car comes, you look at your watch and record the time. Let X_i be the time (in minutes) between when the i^{th} car comes and the $(i + 1)^{th}$ car comes for $i = 1, \dots, 44$. Suppose you know the average time between cars is 1 minute.

Find the probability that the average time gap between cars for the next 44 cars exceeds 1.05 minutes.

x_i : time (minute) between i^{th} car & $(i+1)^{th}$ car \leftarrow

$$x_i \stackrel{\text{iid}}{\sim} \text{Exp}(\alpha=1)$$

$\bar{x} = \frac{1}{44} \sum_{i=1}^{44} x_i$: the average time gap between cars for 44 cars.

$$P(\bar{x} > 1.05) = ?$$

we have iid sample + $n = 44 \geq 25$

use CLT.

$$\bar{X} \sim N \left(E x_i = \alpha = 1, \frac{\text{Var}(x_i)}{n} = \frac{\alpha^2}{n} = \frac{1}{44} \right)$$

$$P(\bar{X} > 1.05) = P\left(\frac{\bar{X} - 1}{\sqrt{\frac{1}{44}}} > \frac{1.05 - 1}{\sqrt{\frac{1}{44}}}\right)$$

$$= P(Z > 0.332) \quad , Z \sim N(0,1)$$

$$= 1 - P(Z \leq 0.332)$$

$$= 1 - \Phi(0.332)$$

$$= 1 - 0.633 = 0.3669.$$

with 36.69% probability, the average time gap between 44

Cars is > 1.05 minute.

Functions of RVs

Linear Combinations

Sample Mean

CLT (assume iid)

Example: [Baby food jars, cont'd]

The process of filling food containers appears to have an inherent standard deviation of measured fill weights on the order of $1.6g$. Suppose we want to calibrate the filling machine by setting an adjustment knob and filling a run of n jars. Their sample mean net contents will serve as an indication of the process mean fill level corresponding to that knob setting.

we're looking for \underline{n} .

You want to choose a sample size, n , large enough that there is an 80% chance the sample mean is within $.3g$ of the actual process mean.



$$* P(\mu - 0.3 < \bar{x} < \mu + 0.3) = 0.8$$

Note that, μ & n are not given!

we're looking for \underline{n}

x_i : the weight of one jar.

\bar{x} : The sample mean weight of n jars.

$$E\bar{x} = E x_1 = E x_i = \mu$$

$$\text{Var}(\bar{x}) = \frac{\text{Var}(x_1)}{n} = \frac{1.6^2}{n}$$

For n large enough, by CLT: $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$
 $\Rightarrow \bar{x} \sim N(\mu, \frac{1.6^2}{n})$

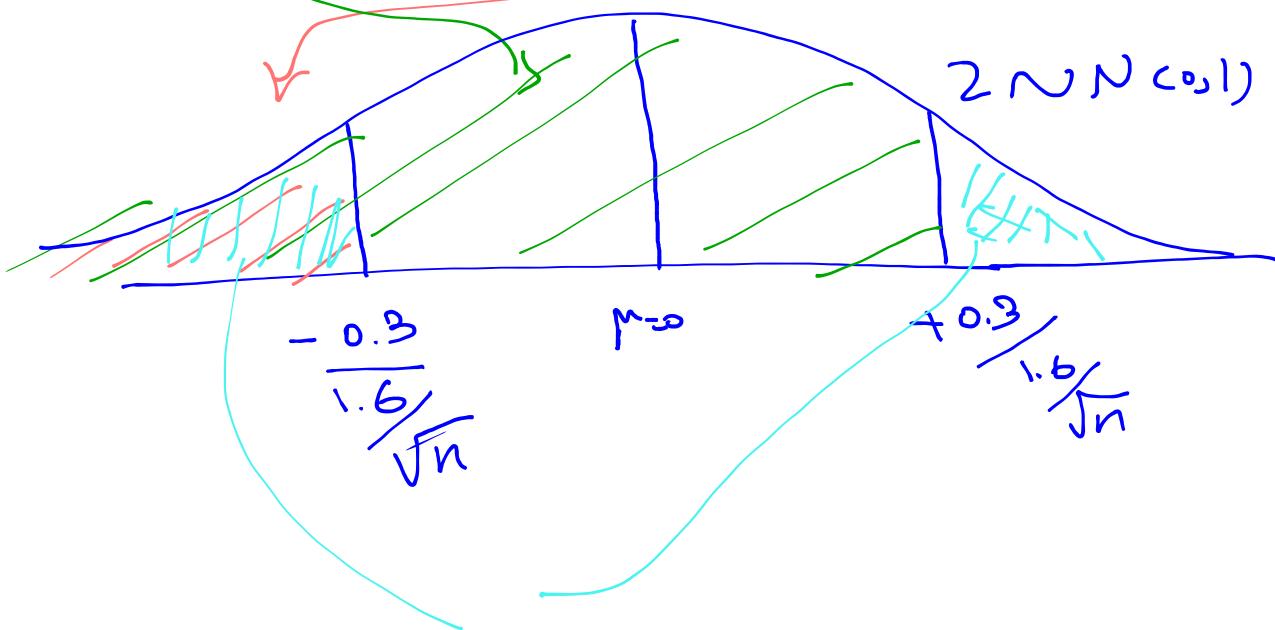
$$P(\mu - 0.3 < \bar{x} < \mu + 0.3) = 0.8$$

$$= P\left(\frac{\mu - 0.3 - \mu}{1.6/\sqrt{n}} < \frac{\bar{x} - \mu}{1.6/\sqrt{n}} < \frac{\mu + 0.3 - \mu}{1.6/\sqrt{n}} \right) = 0.8$$

$$Z \sim N(0, 1)$$

$$= P\left(-\frac{0.3}{1.6/\sqrt{n}} < Z < \frac{+0.3}{1.6/\sqrt{n}}\right) = 0.8$$

$$= \Phi\left(\frac{0.3}{1.6/\sqrt{n}}\right) - \Phi\left(-\frac{0.3}{1.6/\sqrt{n}}\right) = 0.8$$



The same area

$$= \Phi\left(\frac{0.3}{1.6\sqrt{n}}\right) - \left(1 - \Phi\left(\frac{0.3}{1.6\sqrt{n}}\right)\right) = 0.8$$

$$= 2\Phi\left(\frac{0.3}{1.6\sqrt{n}}\right) - 1 = 0.8$$

$$\Rightarrow \Phi\left(\frac{0.3}{1.6\sqrt{n}}\right) = \frac{0.8}{2} = 0.9$$

by the

table $\Rightarrow \frac{0.3}{1.6\sqrt{n}} = 1.29 \Rightarrow \sqrt{n} = \frac{1.29 \times 1.6}{0.3}$

$$\Rightarrow n = \left(\frac{1.29 \times 1.6}{0.3}\right)^2 = 47.3344$$

choose : $\underbrace{n = 48}$

Functions of RVs

Linear Combinations

Sample Mean

CLT

$$\sigma^2 \rightarrow$$

Example: [Printing mistakes]

Suppose the number of printing mistakes on a page follows some unknown distribution with a mean of 4 and a variance of 9. Assume that number of printing mistakes on a printed page are iid.

- What is the approximate probability distribution of the average number of printing mistakes on 50 pages?

$$\bar{x}$$

$$n > 25$$

$$\bar{x} \sim N(4, \frac{9}{50}) \text{ by CLT.}$$

- Can you find the probability that the number of printing mistakes on a single page is less than 3.8?

$$x$$

No, because the probability dist. of # of printing mistakes is unknown on a single page.

Functions of RVs

Linear Combinations

Sample Mean

CLT

Example: [Printing mistakes]

$$\bar{x}$$

- Can you find the probability that the average number of printing mistakes on 10 pages is less than 3.8?

No, because $n=10 < 25$ & we cannot use CLT. Thus, the dist. of \bar{x} is unknown.

- Can you find the probability that the average number of printing mistakes on 50 pages is less than 3.8?

$$\bar{x}$$

Yes. because $n=50 > 25$ & x_i are iid.

So, by CLT $\bar{x} \sim N(\mu=4, \frac{9}{50})$

