

Tossing a coin in the long run

$$\# \text{Heads} = \frac{1}{2} \# \text{Tosses} + \text{Chance Error}$$

(observed) value (^{expected} value)

$$\text{OV} = \text{EV} + \text{CE}$$

① Common mistake = Gambler's Fallacy

As # of tosses ↑, the size of CE ↓

② Correct thinking = Law of Large Numbers (LLN)

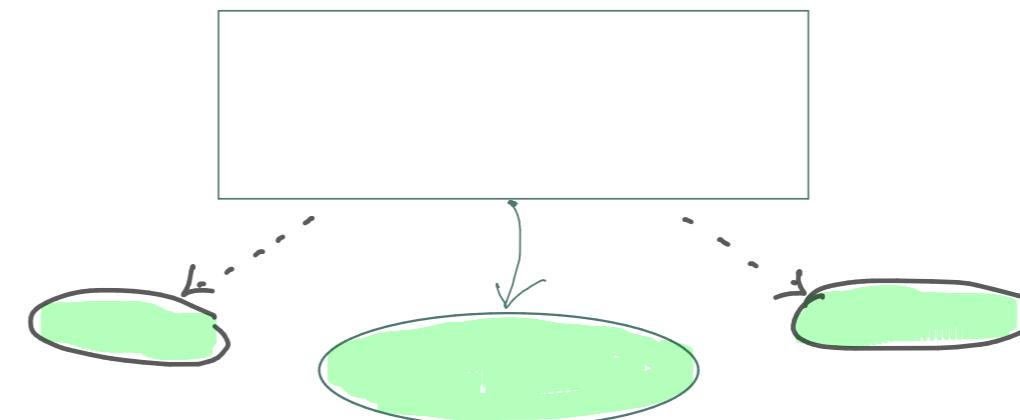
As # of tosses ↑, % Heads $\xrightarrow{\text{stabilises}}$ P(Head)

The Central Limit Theorem (CLT)

When drawing at random from any box, if the sample size is sufficiently large, then

the shape of the Sum or Mean of the Sample is approximately Normal.

Module 3: Sampling Data



T6: Understanding Chance

T7: Chance Variability
(The Box Model)

T8: Sample Surveys

LO6: Use the box model to describe chance & chance variability, including sample surveys & the CLT.

The Box Model

population



with replacement

sample

represented by "tickets"
need to know
- what values
- how many of each

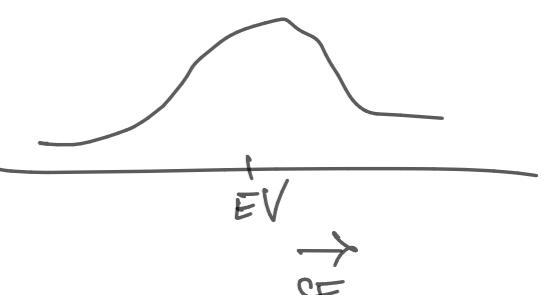
① We want to model the Sum or Mean of an (unknown) Sample.

② $\text{OV} = \text{EV} + \text{CE}$

calculate from a known sample OR simulate samples from the population

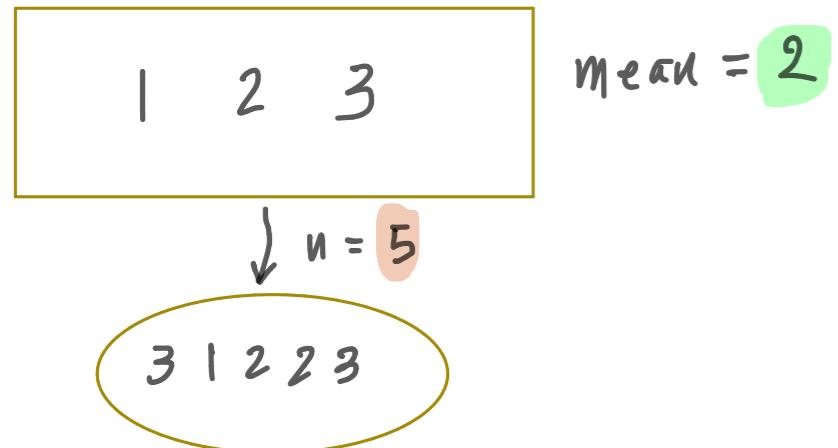
③ Use the CLT.

Eg) Sample Sum $\sim N(\text{EV}, \text{SE}^2)$



Examples of the Box Model

① DV22 p23



② DV22 p29

It costs \$1 to play a game

- if you roll a 6
- you get \$1 back
- plus an extra \$1
- if you don't roll a 6
- you lose your \$1

You play 25 times. What is your expected win/loss?

Consider the Sum of the Sample:

$$\cdot EV = 5 * 2 = 10$$

$$\cdot OV = 3 + 1 + 2 + 2 + 3 = 11$$

$$\text{So } CE = OV - EV = 11 - 10 = 1$$

1 game



Consider the Sum of the Sample

$$\cdot EV = 25 * -\frac{4}{6} = -\$16.67$$

$$\cdot OV = -13$$

$$\text{So } CE = OV - EV = \$3.67$$

simulation

set.seed(1)
die.tosses = sample(c(1, -1), 25, rep = T)
sum(die.tosses)

; ; ;
output is -13
< - - -

📌 Summary of the Box Model

$E_0 / \text{[0,1]}$

$\text{mean} = 0.5$

$SD = 0.5$

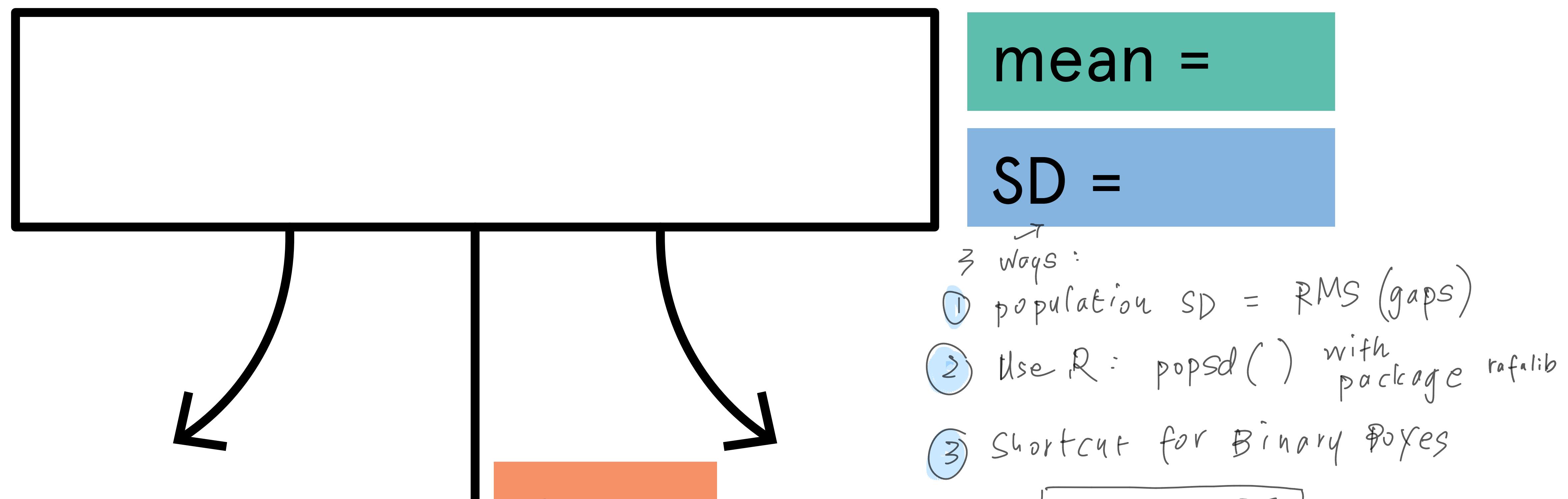
$SD = \sqrt{\frac{(0-0.5)^2 + (1-0.5)^2}{2}}$

② $\text{box} = c(D_{\text{rf}})$
library(rffq.lib)

③ $(1-\theta) \sqrt{0.5 \times 0.5} = 0.5$

The Box Model

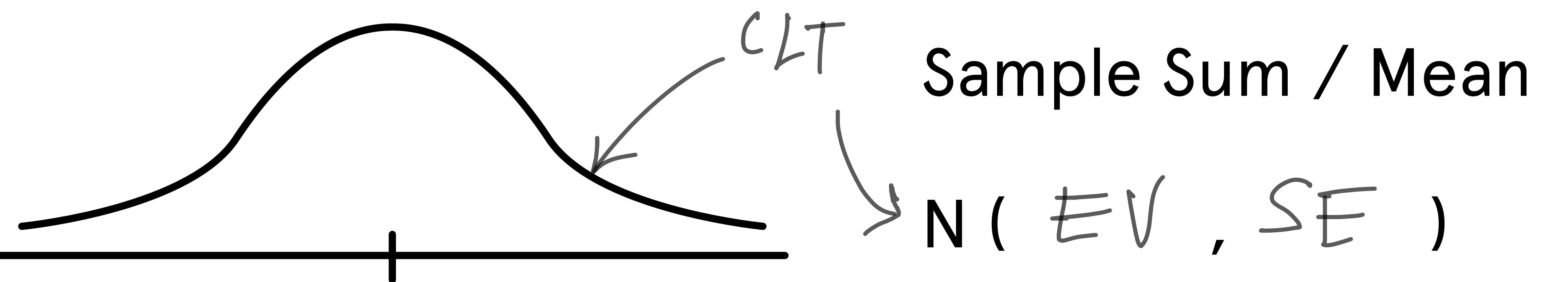
population = "box" = ...



observed sample

observed mean =
observed SD =

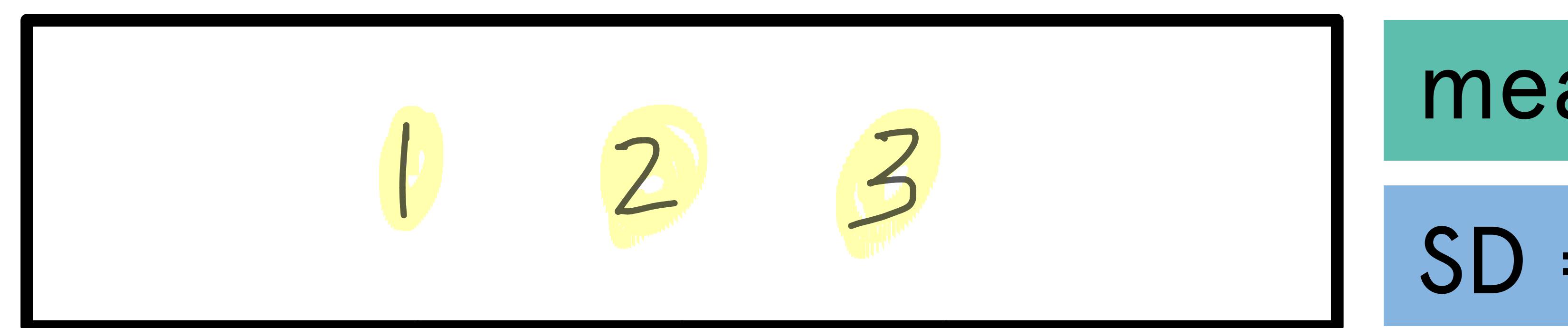
Sample Sum	$EV = n \text{ mean}$	
Sample Mean	$SE = \sqrt{n} SD$	
Sample Sum	$EV = \text{mean}$	
Sample Mean	$SE = \frac{SD}{\sqrt{n}}$	



📌 Summary of the Box Model

The Box Model

population = "box" = ... Box 1



$$\text{mean} = 2$$

$$\text{SD} \doteq 0.82$$

$$① \text{SD} = \text{RMS}(\text{gap}) = \sqrt{\frac{(1-2)^2 + (2-2)^2 + (3-2)^2}{3}} = \sqrt{\frac{3}{3}} = 0.816$$

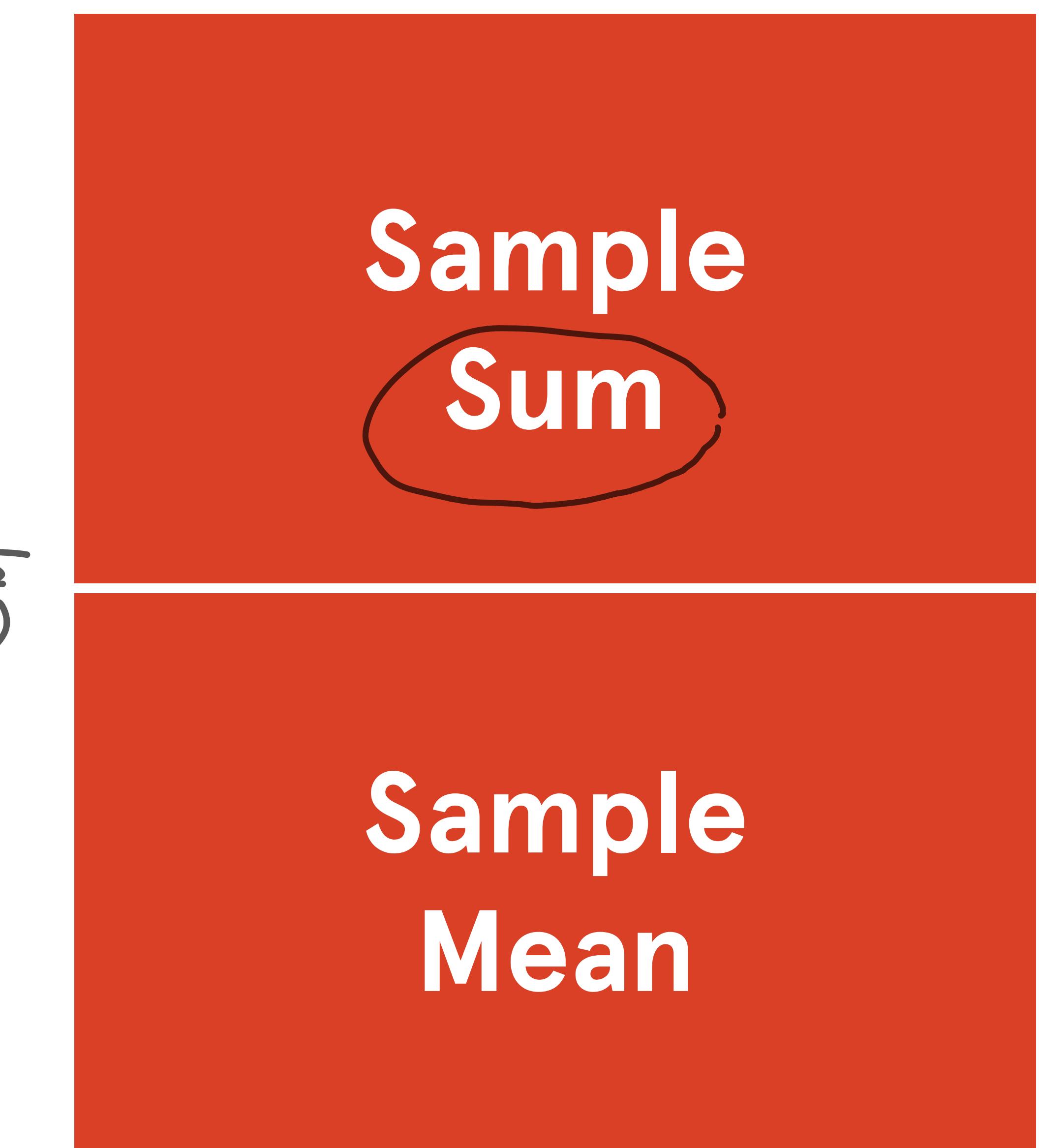
② $\text{box1} = c(1, 2, 3)$
library (`rafal`)
`popsd`(box1)

③ not Binary Box

observed sample

observed mean =
observed SD =

Model of Sample



$$\text{EV} = n \text{ mean}$$

$$25 \times 2 = 50$$

$$\text{SE} = \sqrt{n} \text{ SD}$$

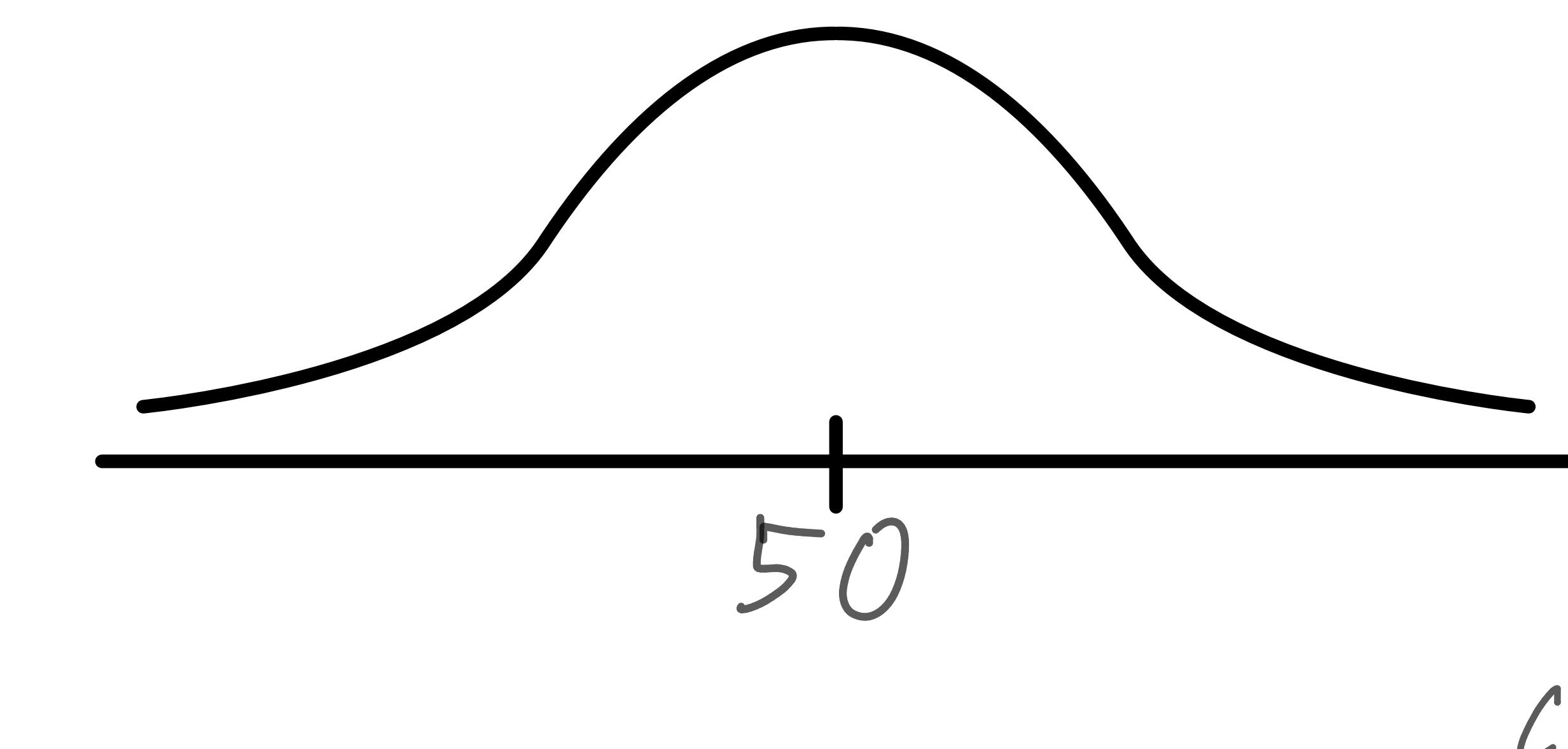
$$\sqrt{25} \times \sqrt{\frac{2}{3}} = 4.1$$

$$\text{EV} = \text{mean}$$

$$\text{SE} = \frac{\text{SD}}{\sqrt{n}}$$

Sample Sum / Mean

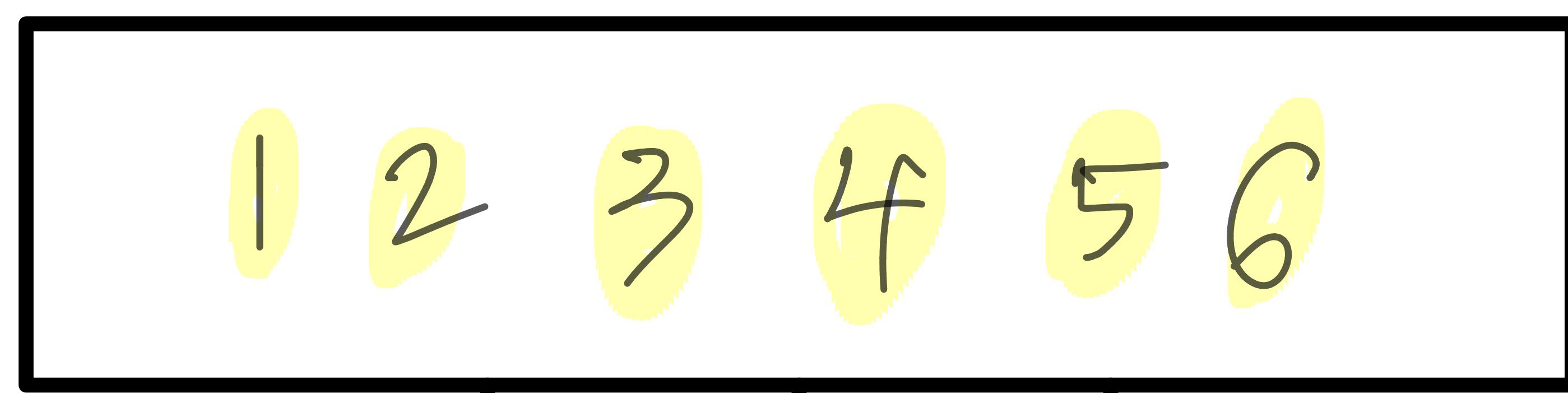
$$N(50, 4.1^2)$$



📌 Summary of the Box Model

The Box Model

population = "box" = ... Box 2



$$\text{mean} = 3.5$$

$$\text{SD} \doteq 1.71$$

$$1 \quad \text{SD} = \text{RMS}(\text{gap}) = \sqrt{\frac{(1-3.5)^2 + \dots + (6-3.5)^2}{6}} \div 1.71$$

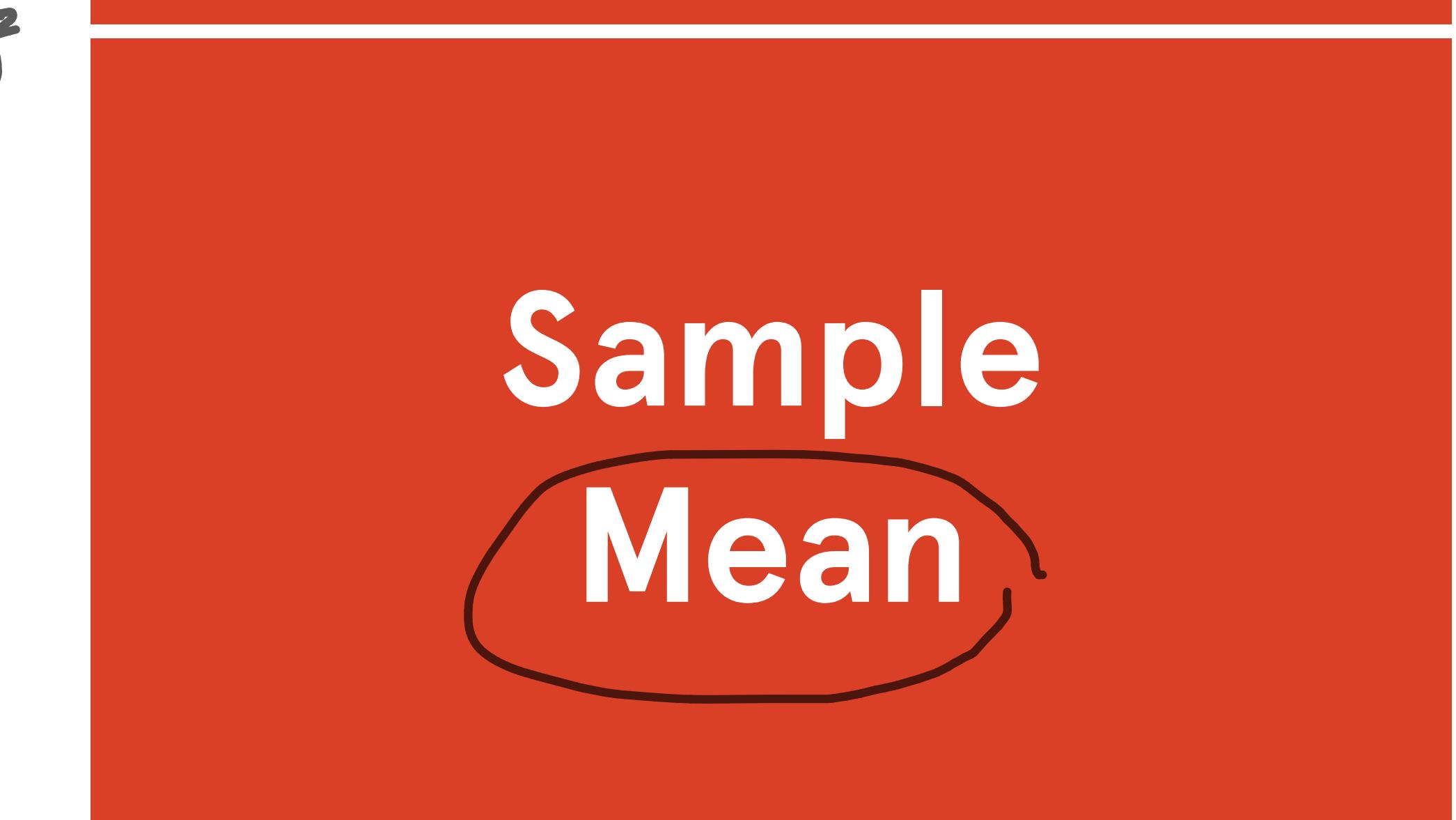
$$2 \quad \text{Box 2} = 1:6 \\ \text{library(rafalib)} \\ \text{popsd(Box2)}$$

$$n = 16$$

observed sample

observed mean =
observed SD =

Model of Sample



$$\text{EV} = n \text{ mean}$$

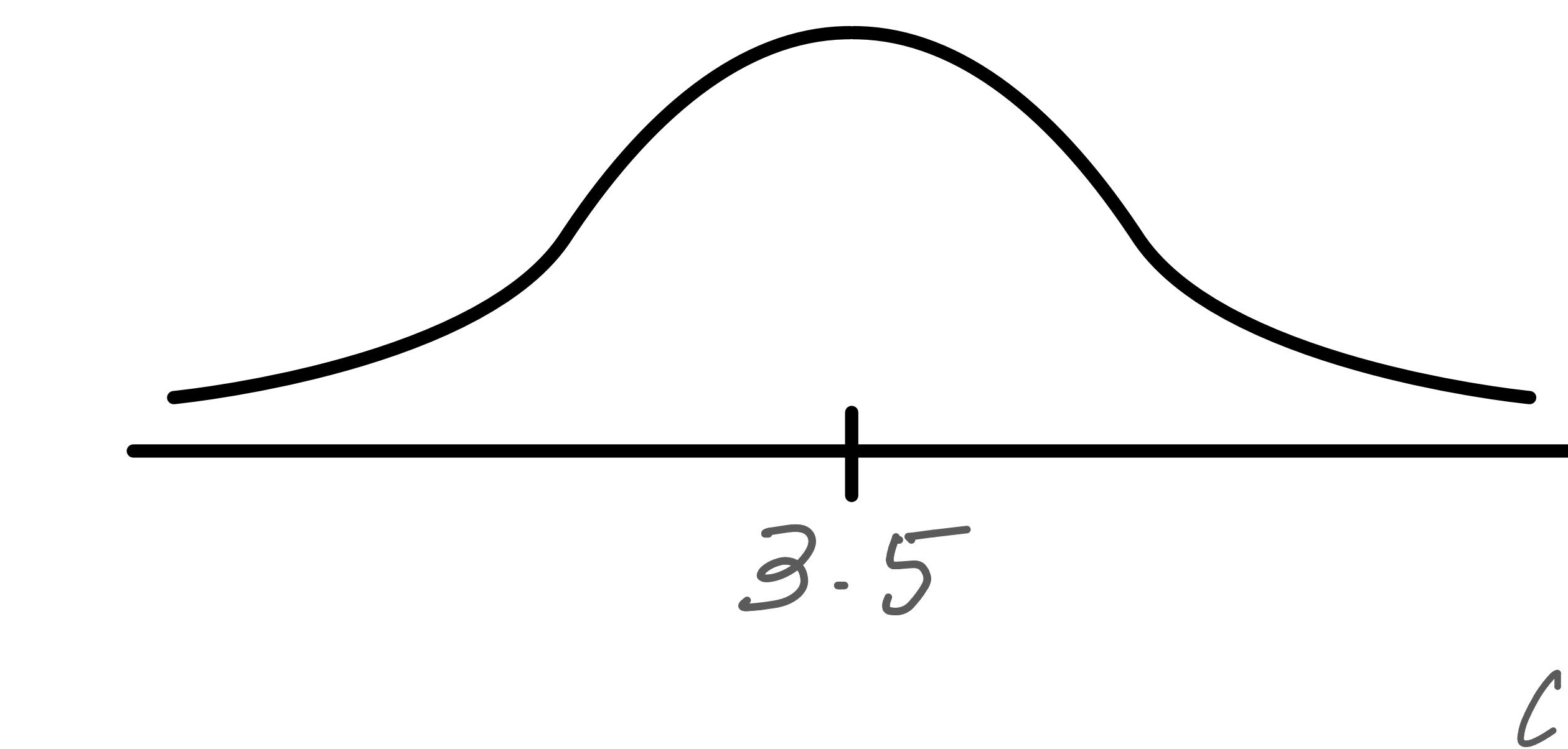
$$\text{SE} = \sqrt{n} \text{ SD}$$

$$\text{EV} = \text{mean}$$

$$\text{SE} = \frac{\text{SD}}{\sqrt{n}}$$

$$3.5$$

$$\frac{1.71}{\sqrt{16}} \doteq 0.43$$



Sample Sum / Mean

$$N(3.5, 0.43^2)$$

Continuity Correction

When we apply the Box Model:

As the box is discrete,

then we can (often) improve

the use of the Normal approximation

by slightly adjusting the thresholds.

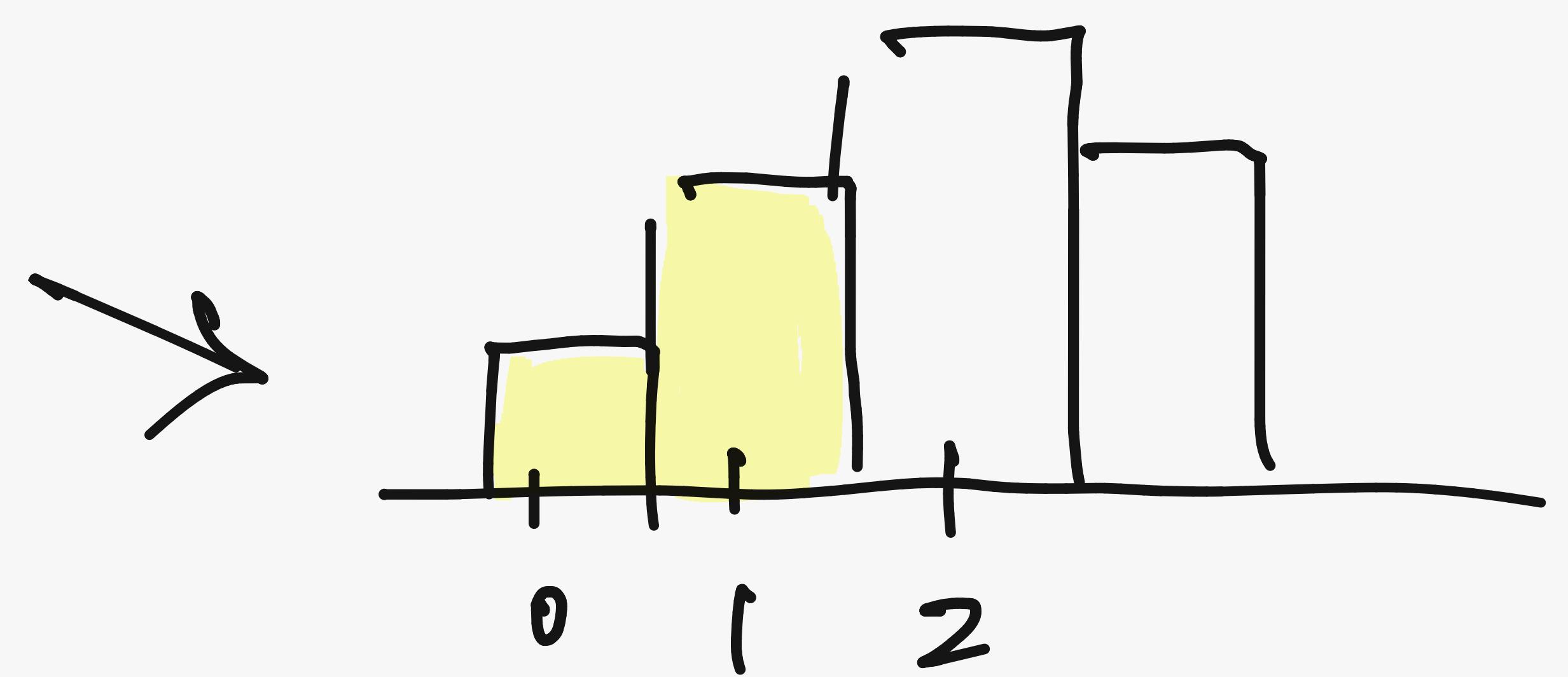
Example:

Suppose the Sample Sum $\sim N(2, 3^2)$
(discrete) \nwarrow continuous

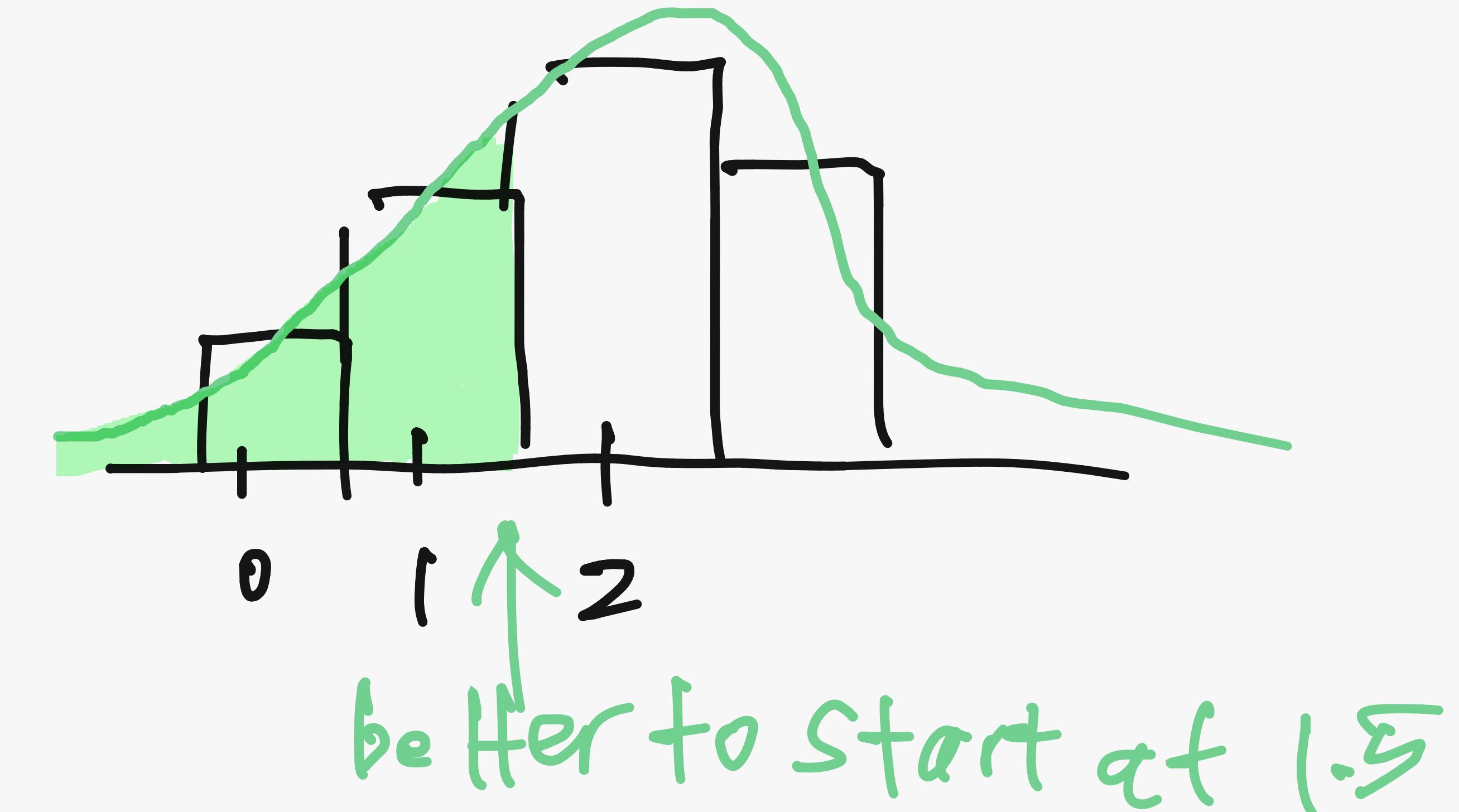
Find $P(X \leq 1)$

① Exact

$$P(X \leq r)$$



② Approximation



Hence we work out

$P(X \leq 1.5)$ for the Normal curve.

`> pnorm(1.5, 2, 3)`