

# Distribution

pdf, cdf

$EX$  (or  $\mu$ )

$$\sum k P(X=k)$$

$\text{Var } X$

$$E\{X - E(X)\}^2$$

# Samples

$X_1, X_2, X_3, \dots$

$$\left\{ \begin{array}{l} \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n} \\ S^2 = \frac{(X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n-1} \end{array} \right.$$

$$E(S^2) = \text{Var } X$$



Sum: I have samples  $x_1, x_2, \dots, x_n \rightarrow (\mu, \sigma^2)$   
 $\underbrace{\hspace{1cm}}_{\text{independent}}$

Question: How does  $Y = x_1 + x_2 + x_3 + \dots + x_n$  look like?

$\rightarrow$  Each time this number will be diff

$$E[Y] = n\mu$$

$$\text{Var}[Y] = n\sigma^2$$

$$\begin{aligned} & \text{Var}(x_1 + x_2) \\ &= \text{Var } x_1 + \text{Var } x_2 \\ & n = \text{num\_samples} \end{aligned}$$

Dice:

Goal is to get trial values =  $[x_1, x_2, x_3, \dots, x_{1000}] \rightarrow \text{hist}$

$$\overline{Y} = \frac{x_1 + x_2 + \dots + x_{1000}}{1000} = (350.5) \rightarrow 350 = 100(3.5)$$

$$Y = X_1 + X_2, \quad X_1, X_2 \text{ are indep}$$

$$\begin{aligned} \text{Var } Y &= \text{Var } X_1 + \text{Var } X_2 \\ &= \sigma^2 + \sigma^2 \\ &= 2\sigma^2 \end{aligned}$$



Central Limit Theorem:  $\{X_i\}$  have mean  $\mu$  & variance  $\sigma^2$

$Y = X_1 + X_2 + \dots + X_n$  is approximately distributed

as Normal / Gaussian with mean =  $n\mu$

and variance =  $n\sigma^2$

$$\frac{Y - n\mu}{\sqrt{n}\sigma} \sim N(0, 1)$$

Standardisation :  $Y \sim N(n\mu, n\sigma^2)$

$$Z = \frac{Y - n\mu}{\sqrt{n}\sigma} \sim N(0, 1)$$

"Standard Normal"

$$E[Y] = n\mu$$
$$\text{Var}[Y] = n\sigma^2$$

$$Z = \frac{Y - E[Y]}{\sqrt{\text{Var } Y}} \text{ is } N(0, 1)$$

$$Y = X_1 + X_2 + \dots + X_n$$

$$Y - n\mu = (X_1 - \mu) + (X_2 - \mu) + \dots + (X_n - \mu)$$



Person 1:  $[6, 3, 1, 3, 4, 5, 6, 2, 1] \xrightarrow{(100 \text{ times})} (27)$   
 $6 + 3 + 1 + 3 + 4$

Person 2:  $[1, 5, 1, 3, 4, 4, 5, 2, 1] \rightarrow$   
 $1 + 5 \xrightarrow{(35)}$

$[27, 35, \dots] \rightarrow$  what distribution

Ex: Insurance company  $\rightarrow$  25,000 policy holders  
Yearly claim  $\rightarrow$  mean 320, S.D = 540

Approximately, Prob. claim is greater than 8.3 million

$Y$ : Total Yearly claim

$$Y = X_1 + X_2 + \dots + X_{25,000}$$
$$E(Y) = 25,000 \times 320 = 8 \text{ mil}$$

$$P[Y > 8.3 \text{ m}] = P\left[\frac{Y - 8 \text{ mil}}{540 \sqrt{25000}} > \frac{8.3 \text{ mil} - 8 \text{ mil}}{540 \sqrt{25000}}\right]$$
$$\left\{ \begin{array}{l} EY = 8 \text{ mil} \\ \text{Var } Y = n 540^2 \\ \sqrt{\text{Var } Y} = \sqrt{n} 540 \end{array} \right.$$

$$= P[Z > 3.51]$$

$$= 1 - P(Z \leq 3.51) = 1 - \text{CDF}(3.51)$$
$$= 1 - \text{norm.cdf}(3.51)$$



Binomial:  $Y$ :  $n$  tosses,  $k$  heads  $k \leq n$   $P\{H\} = p$   
 $P[Y=k] \rightarrow$  "Binomial Random variable"

$$Y = X_1 + X_2 + X_3 + \dots + X_n$$

$X_1 \rightarrow \begin{cases} 0 & \text{if first toss is tail} \\ 1 & \text{if first toss is head} \end{cases}$

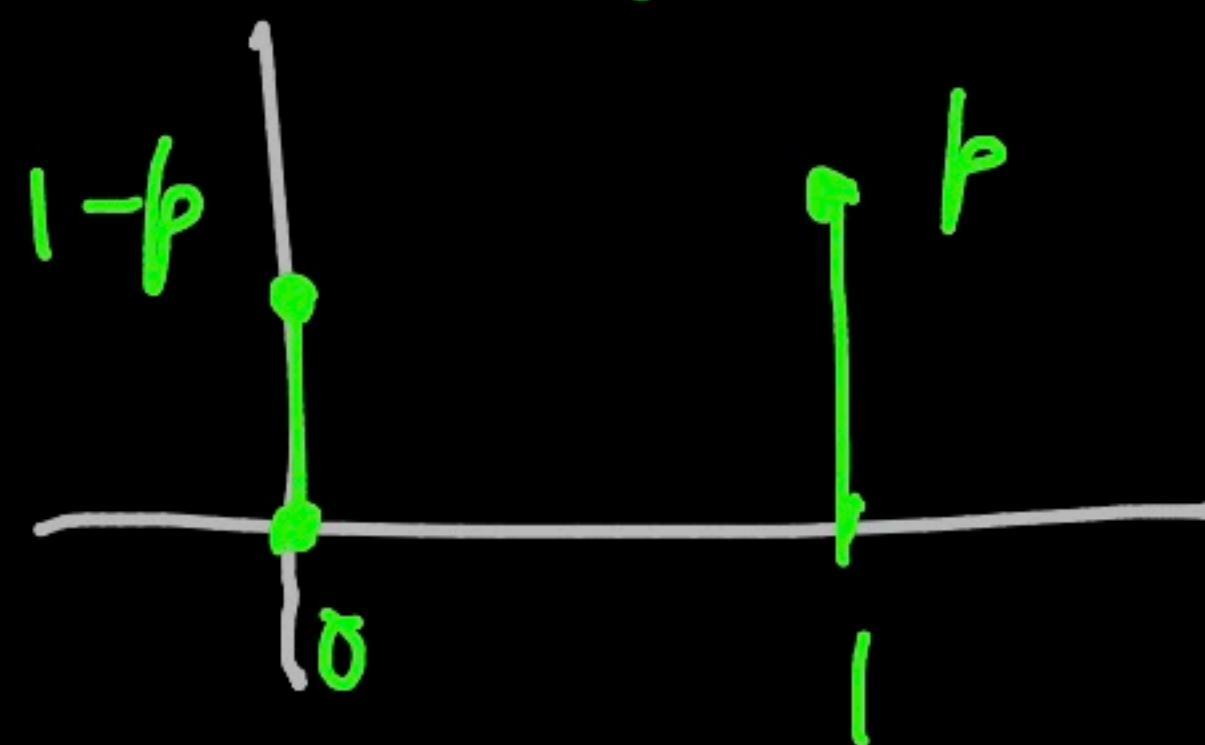
$X_i \rightarrow \begin{cases} 0 & i^{\text{th}} \text{ toss is tail} \\ 1 & \text{heads} \end{cases}$

What is  $\text{Var } Y = n p (1-p)$

$$E Y = n p$$

Bernoulli

$X_i$   $E X_i = p$   
 $\text{Var } X_i = p(1-p)$



# Gaussian Approximation to Binomial :

$$Y = X_1 + X_2 + \dots + X_n$$

$$\frac{Y - np}{\sqrt{np(1-p)}} \sim N(0,1)$$

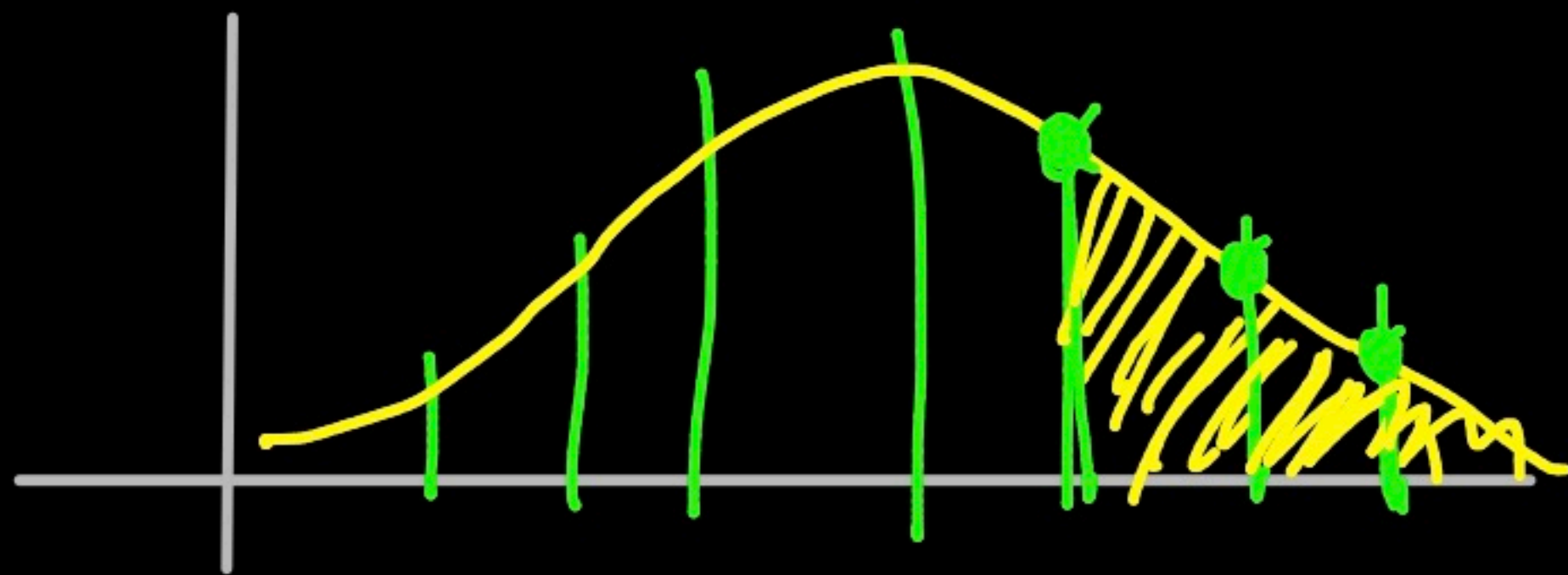
$X_i$  : Bernoulli  $p$

$$EY = np$$

$$\text{Var } Y = np(1-p)$$

$$np(1-p) \geq 40$$

$np$  : random binomial  
 $\uparrow$   
 $(n, p)$





Eg: 150-capacity class.

Among those accepted, only 30% attend  $\rightarrow$  others drop out

College accepts 450 student.

What is the prob of more than 150 people attending

$$\rightarrow n = 450 \quad p = 0.3 \quad E[Y] = np = (450)(0.3) = 135$$

$$\text{Var}[Y] = 94.5 \rightarrow np(1-p) = (450)(0.3)(0.7)$$

$$P[Y > 150] = P\left[\frac{Y - 135}{\sqrt{94.5}} > \frac{150 - 135}{\sqrt{94.5}}\right]$$

$$= P[Z > 1.5] \approx 0.06$$

Sample mean:  $X_1, X_2, X_3, \dots$  mean  $\mu$ , Variance  $\sigma^2$

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

$$E[\bar{X}] = \mu, \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

$$Z = \frac{\bar{X} - E[\bar{X}]}{\sqrt{\text{Var}[\bar{X}]}} = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}$$

$\downarrow$   
 $N(0,1)$



Eg: Weights "population"  $\rightarrow$  everyone  
 $\mu = 167$  pounds,  $\sigma = 27$   $\rightarrow$  std dev.

1) A sample of 36 people are taken.  $\bar{X} \rightarrow$  specific sample mean for those 36 people [163, 170]

$$E[\bar{X}] = 167 \quad \text{Var}[\bar{X}] = \frac{27^2}{36} \quad \sqrt{\text{Var}(\bar{X})} = 4.5$$

$$P[163 < \bar{X} < 170] = P\left[\frac{163-167}{4.5} < \frac{\bar{X}-167}{4.5} < \frac{170-167}{4.5}\right]$$

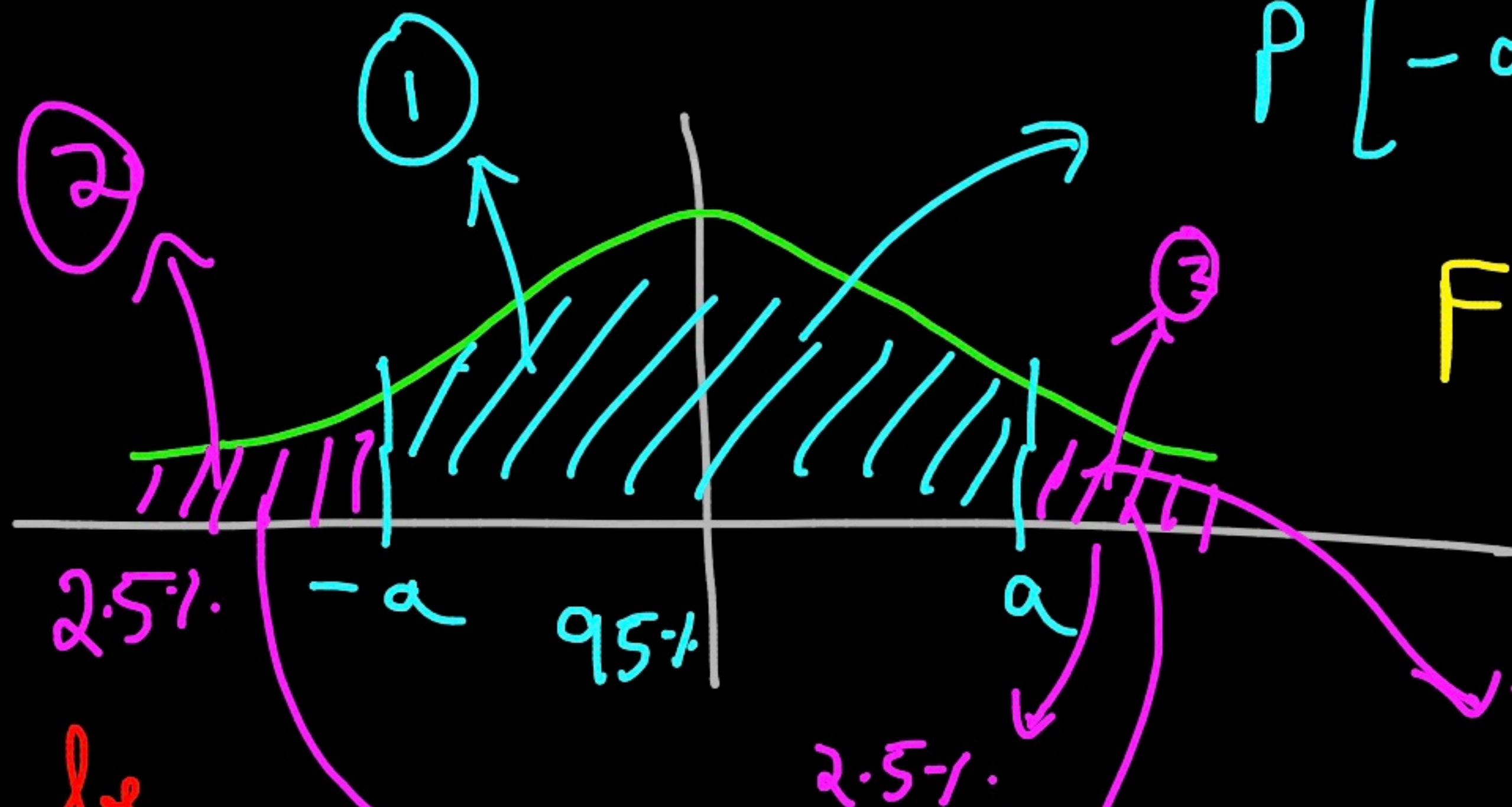
$$= P[-0.889 < Z < \underline{\underline{0.889}}]$$

$\rightarrow$  not exact  
later  
discretized cont



$$P[-a < z < a]$$

$$F(a) = \text{blue} + \text{one margin} \\ \textcircled{1} + \textcircled{2}$$



What will be  
the value of "a"  
if blue area is  
95%.

$$F(a) = 0.975$$

These two areas  
are equal

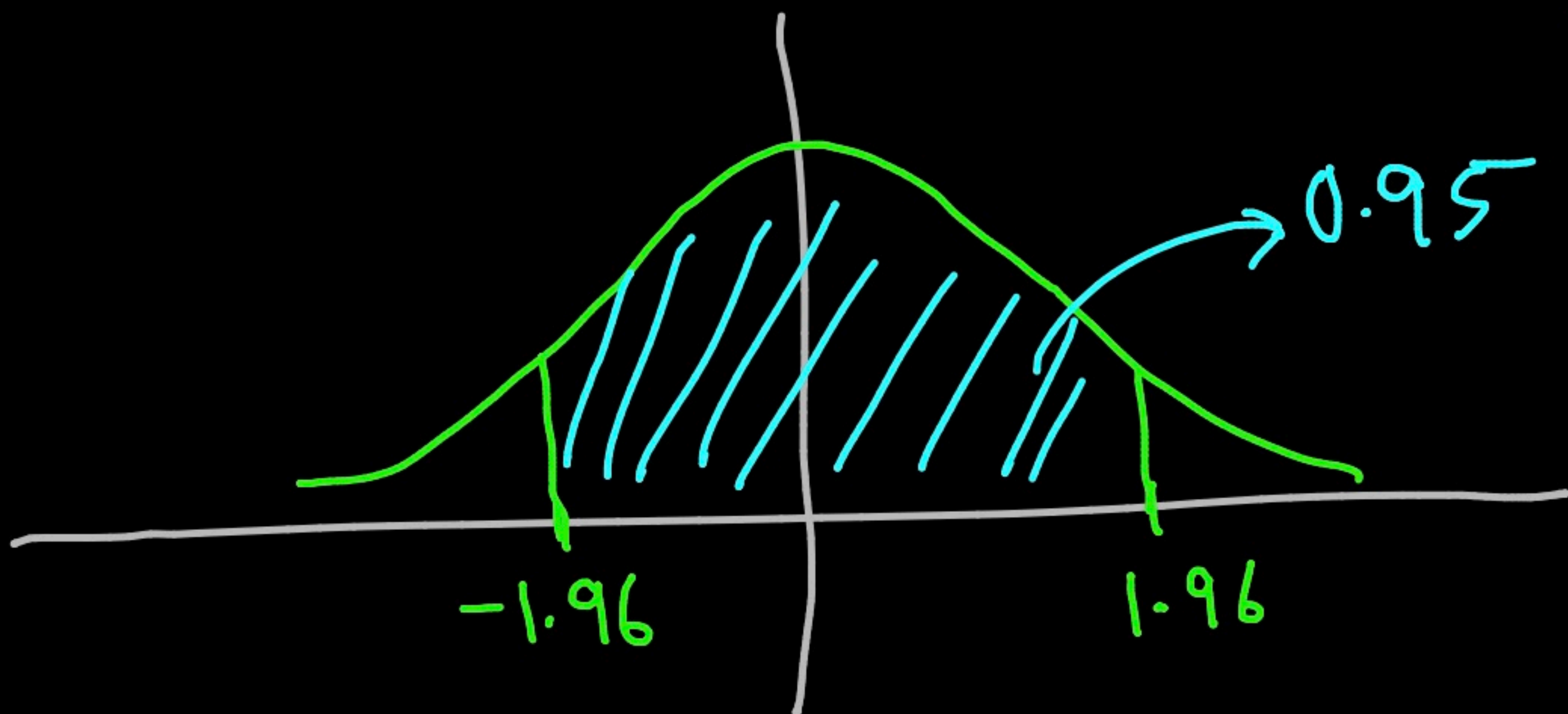
this piece → area?  
 $1 - F(a)$

$$\textcircled{2} = \textcircled{3} = 1 - F(a)$$

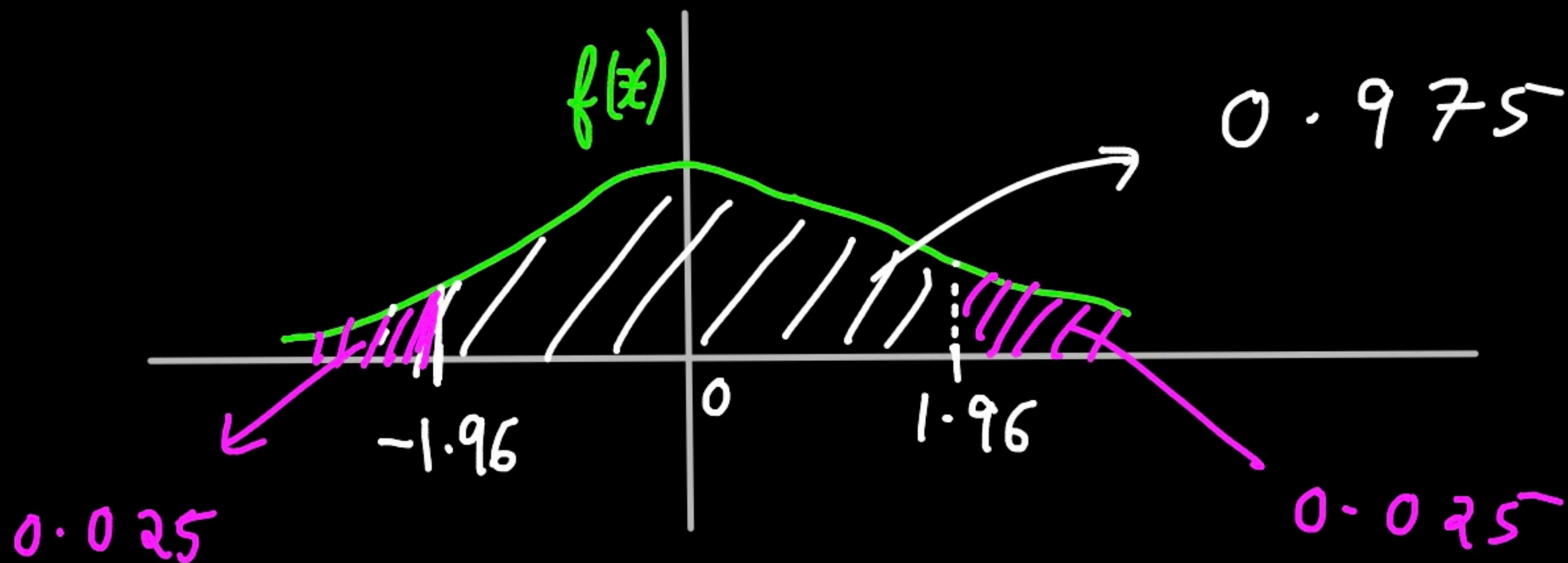
$$\text{blue} = F(a) - (1 - F(a)) \\ = 2F(a) - 1$$



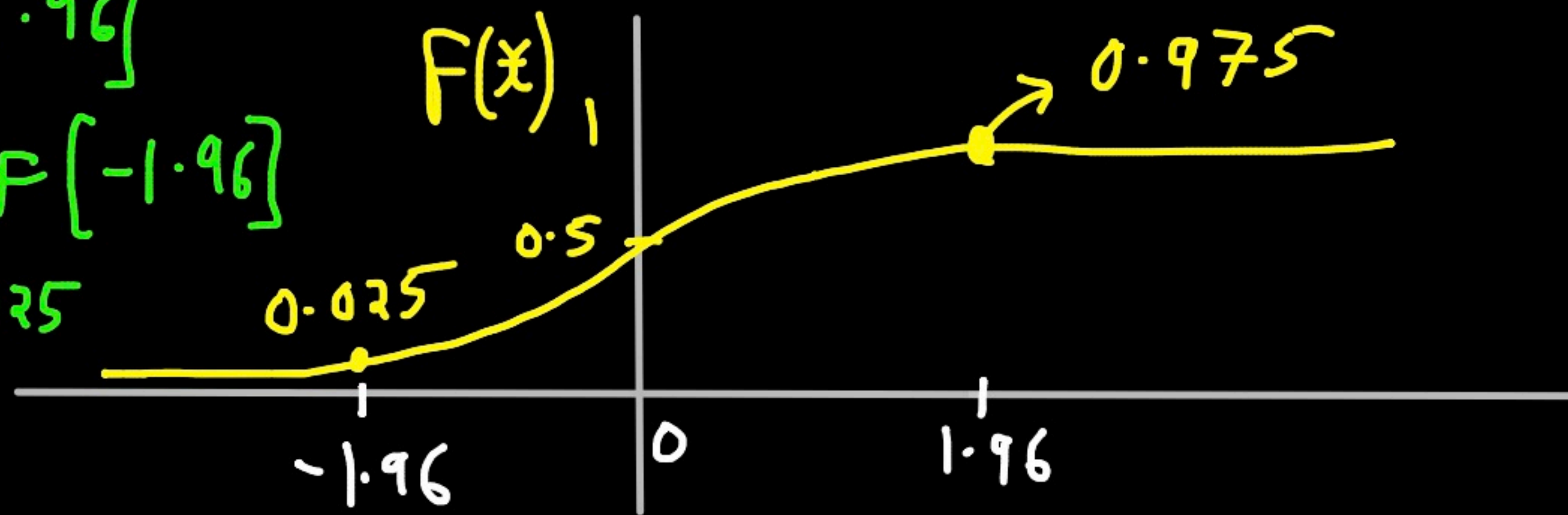
$$\mu = 0$$
$$\sigma = 1$$



$$P[-1.96 < Z < 1.96] = 0.95$$



$$\begin{aligned}
 P(-1.96 < Z < 1.96) \\
 &= F(1.96) - F(-1.96) \\
 &= 0.975 - 0.025 \\
 &= 0.95
 \end{aligned}$$



$$\begin{aligned}
 P(F(0.975)) \\
 &= 1.96 \\
 \text{cdf}(1.96) \\
 &= 0.975
 \end{aligned}$$



Samples  $X_1, X_2, X_3, \dots \rightarrow N(\mu, \sigma)$  (for now)

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

$$Z = \frac{\sqrt{n} \bar{X} - \mu}{\sigma} \rightarrow N(0, 1)$$

$$P\left[-1.96 < \frac{\sqrt{n} \bar{X} - \mu}{\sigma} < 1.96\right] = 0.95$$

$$P\left[\frac{-1.96 \sigma}{\sqrt{n}} < \bar{X} - \mu < 1.96 \frac{\sigma}{\sqrt{n}}\right] = 0.95 \quad - (1)$$

$$P\left[\frac{-1.96 \sigma}{\sqrt{n}} < \mu - \bar{X} < 1.96 \frac{\sigma}{\sqrt{n}}\right] = 0.95$$

$$P\left[\bar{x} - \frac{1.96\sigma}{\sqrt{n}} < \underline{\mu} < \bar{x} + \frac{1.96\sigma}{\sqrt{n}}\right] = 0.95$$

Summary: Population mean lies in an interval

$$\left[\bar{x} - \frac{1.96\sigma}{\sqrt{n}}, \bar{x} + \frac{1.96\sigma}{\sqrt{n}}\right] \text{ with}$$

95% Confidence Interval

95% confidence



Schwag

400  $\rightarrow \mu_s$

"population mean"



50 matches randomly

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_{50}}{50}$$

How far from  $\mu_s$

Dravid

400  $\rightarrow \mu_d$

"population mean"



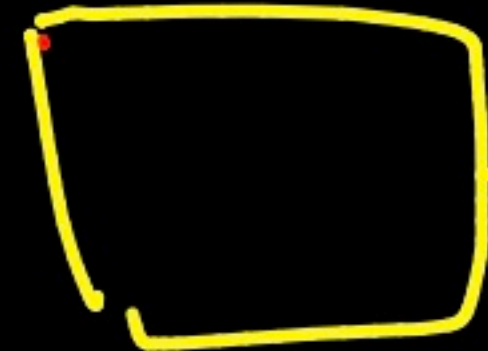
50 matches randomly

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_{50}}{50}$$

How far from  $\mu_d$



Source



signal: number



" $\mu$ "  $\rightarrow$  actual number

noise is added



Receiver

$\mu + N$



$\rightarrow \mathcal{N}(\mu, \sigma^2=4)$

$\mathcal{N}(0, \sigma^2=4)$

5, 8.5, 12, 15, 7, 9, 7.5, 6.5, 10.5  $\rightarrow \frac{81}{9} = 9$   
 $x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_9$

$\bar{x} = 9$ . Can you comment on  $\mu$ ?

$$\left[ 9 - \frac{(1.96)(2)}{\sqrt{9}}, 9 + \frac{(1.96)(2)}{\sqrt{9}} \right] = \underbrace{[7.69, 10.31]}_{95\% \text{ CI}}$$