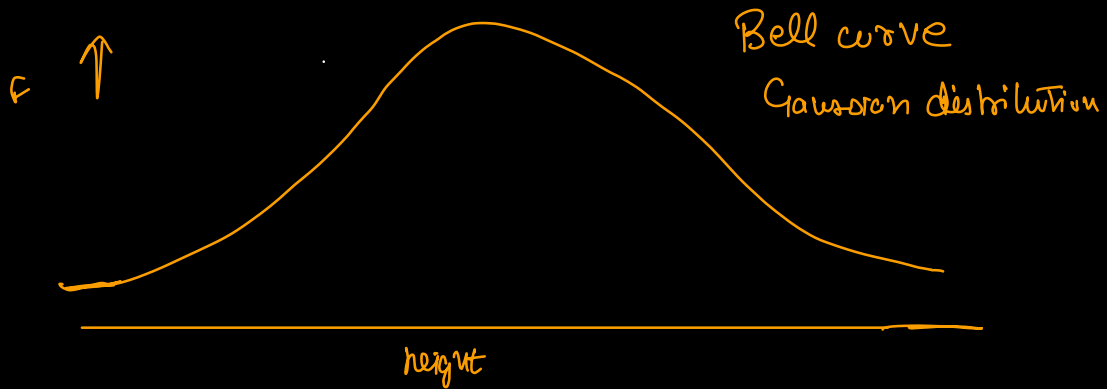## Welcome Back

## Agenda for Today -
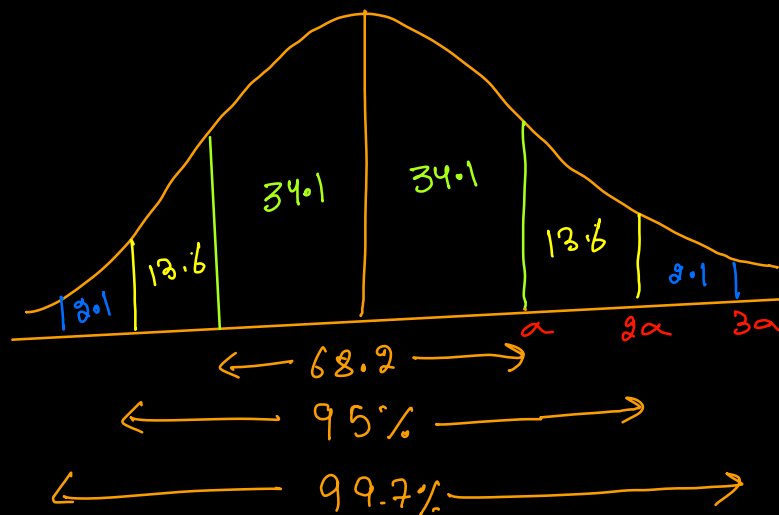
1. Graphs / Distribution in Data Science

2. Inferential Stats / Hyp. testing.

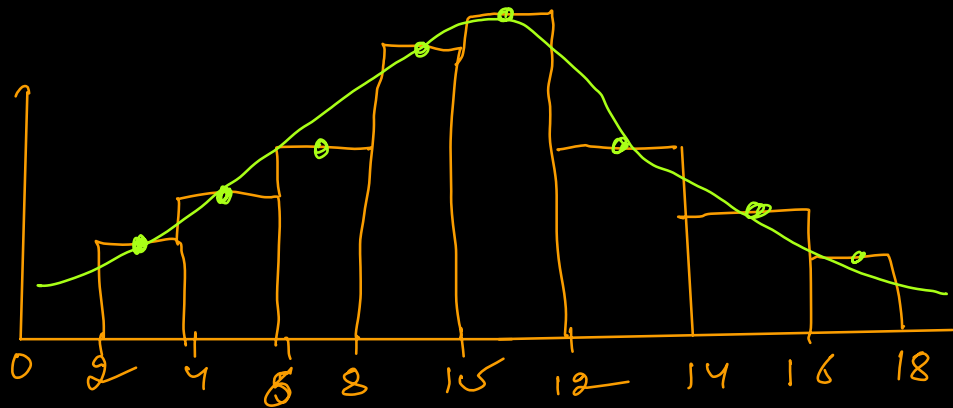# Distribution

1. Normal / Gaussian Distribution



$F \uparrow$

Bell curve
Gaussian distribution

height

⇒ Emphirecal Formula of normal distn.



2.1  13.6  34.1  34.1  13.6  2.1

← 68.2 →

$a$  $2a$  $3a$

← 95% →

← 99.7% →

2 3 5    7  89 10 11 12    17 18



0 2 4 6 8 10 12 14 16 18

## Standard Normal Distribution

$\{1, 2, 3, 4, 5\}$

$\mu = 3$

$\sigma = \sqrt{2} = 1.414$

$\approx 1$ (assume)



1 2 3 4 5

$$z\text{-score} = \frac{x - \mu}{\sigma}$$

# Standordization

| Values (x) | $\mu$ | $x - \mu$ | $\sigma$ | $x - \mu/\sigma$ |
|---|---|---|---|---|
| 1 | 3 | -2 | 1 | -2 |
| 2 | 3 | -1 | 1 | -1 |
| 3 | 3 | 0 | 1 | 0 |
| 4 | 3 | 1 | 1 | 1 |
| 5 | 3 | 2 | 1 | 2 |

-2   -1   0   1$\sigma$   2$\sigma$

100   200   300   400   500

-200 - 100   0   100   200

| Age | Wt | Salary |
|-----|-----|--------|
| 25 | 60 | 3000 |
| 30 | 75 | 40000 |
| 32 | 63 | 500k |
| 34 | 45 | 70R |

diff unit

Age



Age



In order to do this we apply

Z - score

$$x_i = \frac{x - \mu}{\sigma}$$

3.35 σ → 0

# Normalization [Min-max Scale]

we specify the $\min^m$ and $\max^m$ value

$$x_{norm} = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

| f1 | |
|----|----|
| 2 | 0.14 |
| 5 | 0.57 |
| 8 | 1 |
| 6 | 0.71 |
| 1 | 0 |

$$x_2 = \frac{2-1}{8-1} \Rightarrow \frac{1}{7}$$

In most ML cases, standardization

In DL, CNN, we have <u>pixels</u>

We normalize to 0-1

Standard normal distribution

68

$2a$    $1a$    $0$    $1a$    $2a$

95

$x_1$    $x_3$    $x_2$      $\Rightarrow$    $z_1$    $z'_3$    $z_2'$

# Log Normal Distribution

One of the type of skewed data
is log normal
distribution.

Wealth of people

length of
comments
on Youtube

Can we convert this into Gaussian
distribution?

If we take natural log of
each and every value, I
will get normal / Gaussian distri-

$$x_i' = \ln(x_i)$$

$\log_{10}$
normal log

$\log_e$
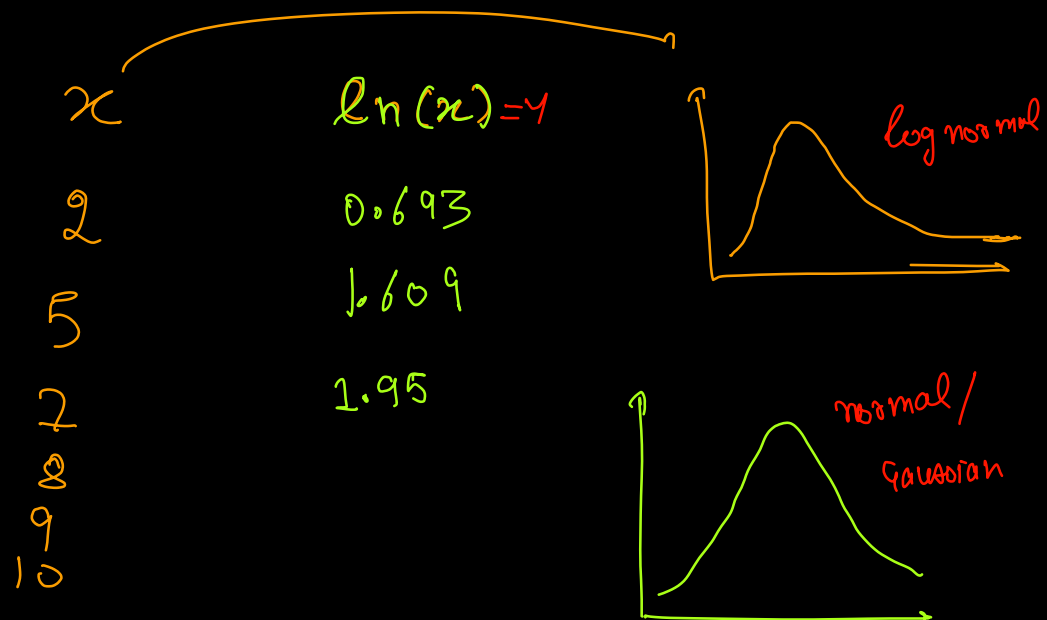natural log

$$2^3 = 8$$
$$3 = \log_{2} 8$$

$$\log_{a} b = x$$
$$a^{x} = b$$

$$e = 2.73 \dots$$

$$\underline{\log_{e}}$$

$$\log_{10} 1000 = 3$$

$$\pi$$

If a variable $x$ is following a log-normal distribution then $y = \ln(x)$

will follow normal distribution

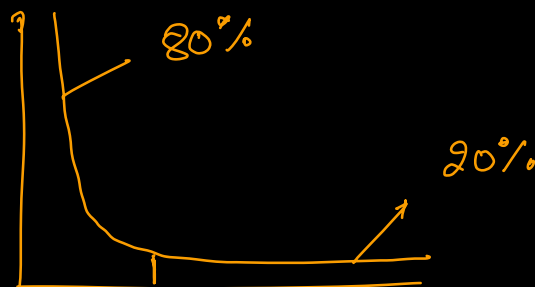| $x$ | $\ln(x) = y$ |
|-----|--------------|
| 2 | 0.693 |
| 5 | 1.609 |
| 7 | 1.95 |
| 8 | |
| 9 | |
| 10 | |

log normal

normal / Gaussian

# Power Law    80:20 rule

pareto principle

80% of your outcome comes from
20% effort


80%

20%

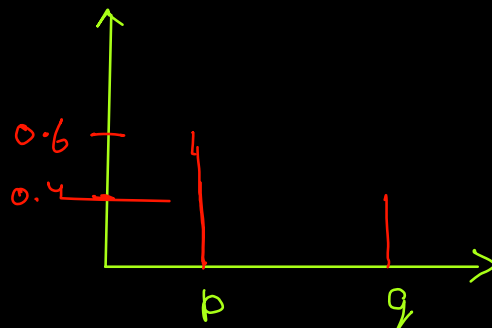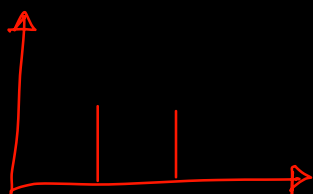## Bernaulis Distribution
→ discrete values

flib a coin

6 heads
10 total

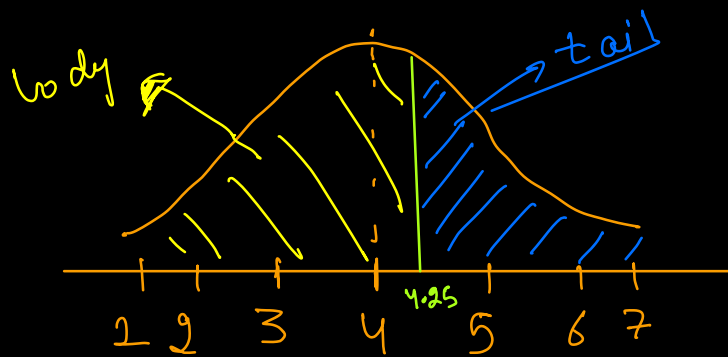$$p = \frac{6}{10}$$

$$q = 1 - p = \frac{4}{10}$$

0.6

0.4

p        q

# Inferential Statistics

$$z = \frac{x_i - \bar{x}}{\sigma}$$

We take a standard normal distn.
with $\mu = 4$ $\sigma = 1$
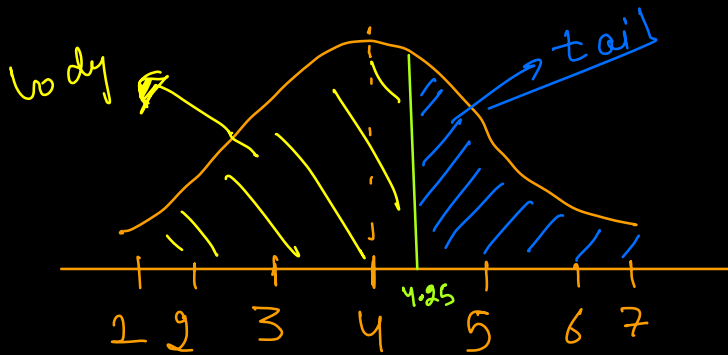
body

tail

4.25

1 2 3 4 5 6 7

Where does
4.25 fall?

⇓

How many
S.D. does
our 4.25
falls from mean?

$$z = \frac{4.25 - 4}{1}$$

$$= .25$$

What % of scores/values falls above
4.25 ?



body

tail

I need to calculate the area of blue

Can I assume full area to be 1 ?

For are, let us calculate &
see z-score.

If body area is .5987

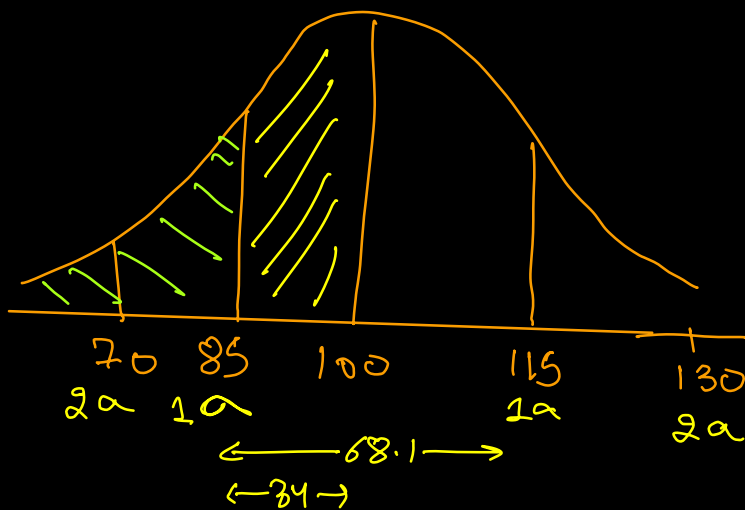area of tail.

$$\Rightarrow 1 - .5987$$

$$\Rightarrow 40.13\%$$

In Mars, avg IQ is 100 with a
S.D. of 15 .

What % of population would
you expect to have IQ less

than 85 ?

$u = 100$
$a = 15$



70    85     100      115      130
$2a$   $1a$           $1a$     $2a$

← 68.1 →
← 34 →

$$Z = \frac{85 - 100}{15}$$

$$= \frac{-15}{15} = -1$$

Value from Z - table
for $z = -1$

$$= 0.1587$$

$$\Rightarrow \quad 15.87\%$$

% of pop$^n$ having IQ b/w
85 −100

0.5 − 0.1587

# Hypothesis testing , Confidence
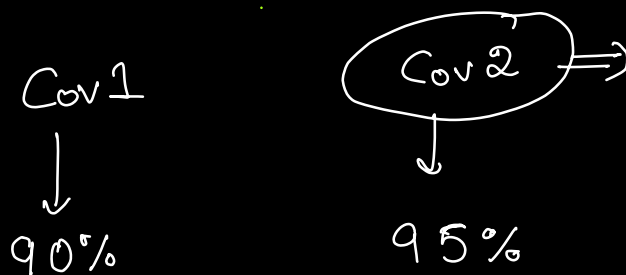Interval ,

Significance
value.

## Coin toss experiment : Check whether
a coin is fair or
not by doing
100 tosses

$p(n) = \underline{0.5}$

$p(t) = 0.5$

Cov 1

$\downarrow$

90%

Cov 2 $\rightarrow$

$\downarrow$

95%

fair

not fair

| 50H | 60H | 80H | 90H |
| 50T | 40T | 20T | 10T |

$\Downarrow$

fair

# Hypothesis testing

$H_0$ = null hypothesis = default assump$^n$

= coin is fair

$H_1$ = alternate hypothesis $\Rightarrow$ opp. of default assump$^n$

= coin is not fair

## Perform experiment

## Reject / Accept our null Hypotheses
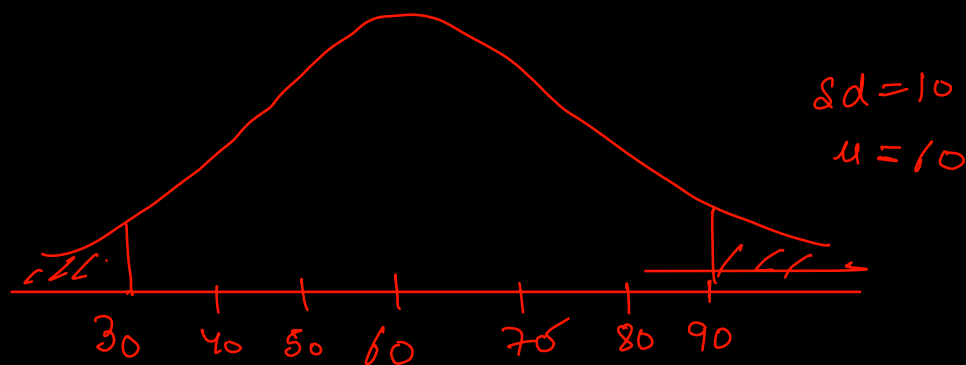
| 50 H | 80 H | 90 H |
|------|------|------|
| 50 T | 20 T | 10 T |

$\Rightarrow$ Whenever we have a condition / situation

like above we define few more param.

2 params.

Confidence Interval , significant value
(α)

Using α, we
calculate confidence
interval

$sd = 10$
$u = 10$



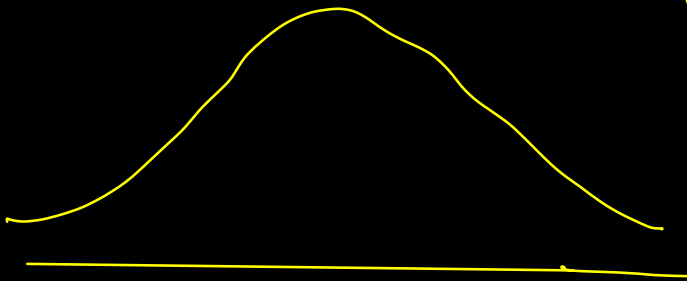30   40  50  60      70   80  90

below 30   , we have a doubt

So, confidence interval basically says that
we will be defining some range &
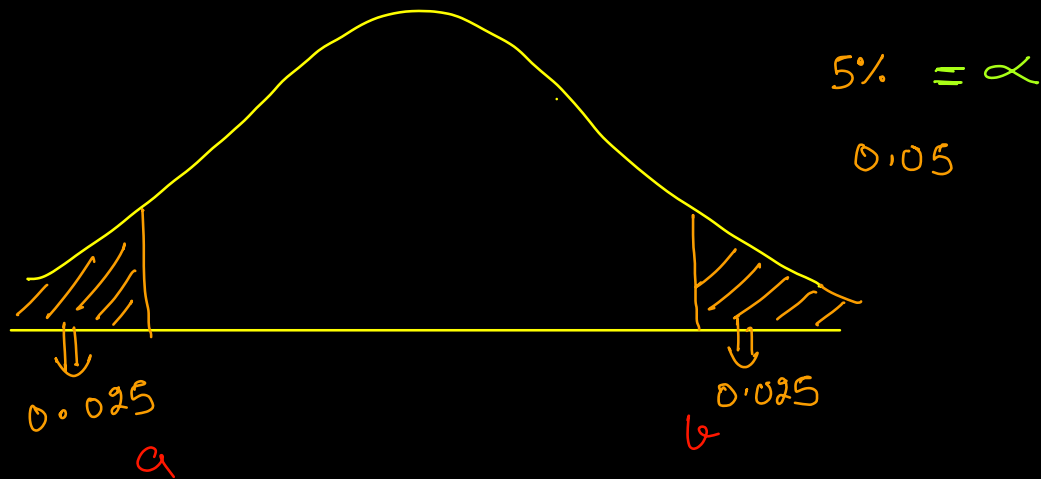if our   value falls b/w that range
then   our coin is fair (Hov)

Who defines/ this confidence Interval ?
      decides ?

   Domain Expert person defines
       this interval,


eg :—   health core

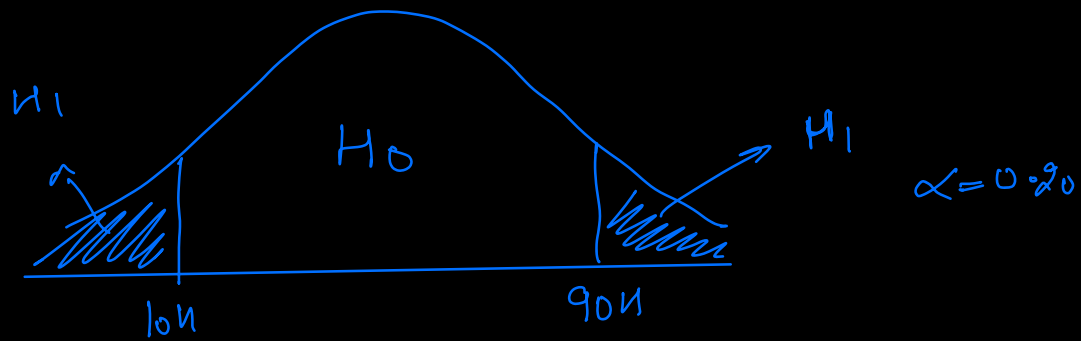   Covid test    ⟶ false -ve   ⟶ symphons.
                      report

5% $= \alpha$

0.05

0.025

0.025

a

b

significance value $= 0.5\%$

How to find confidence Interval (C. Io)
from $\alpha$ ?

$$C.I_o = 1 - 0.025 - 0.025$$

$$= 95\%$$

If my experiment yields that
value fall b/w $a - b$

I will accept Ho

$H_1$
$H_0$
$H_1$
$\alpha = 0.20$

10μ    90μ

Who define this $\alpha$ (significant value)

Vaccine                    100 people take it

                               ⇓

domain                     { 30 don't get covid
expert                     { 60 don't get covid
                           { 75 don't get covid
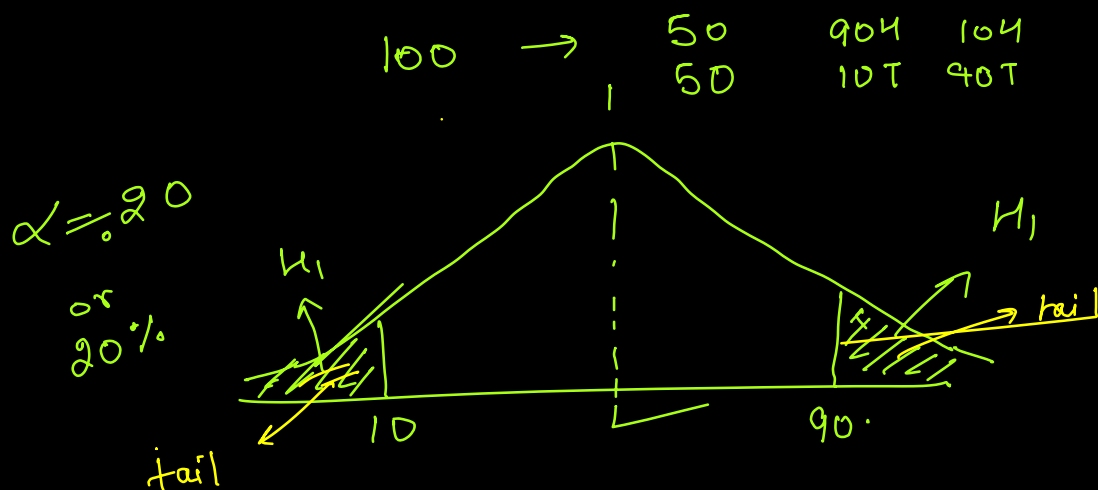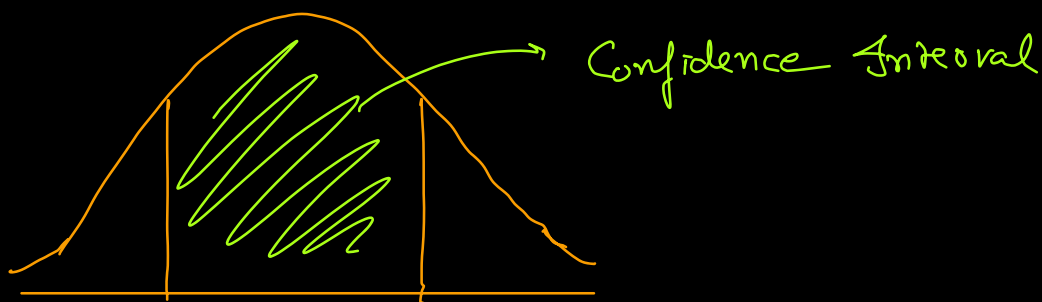
medical cases we normally
have a very high $\alpha$

# Confidence Interval
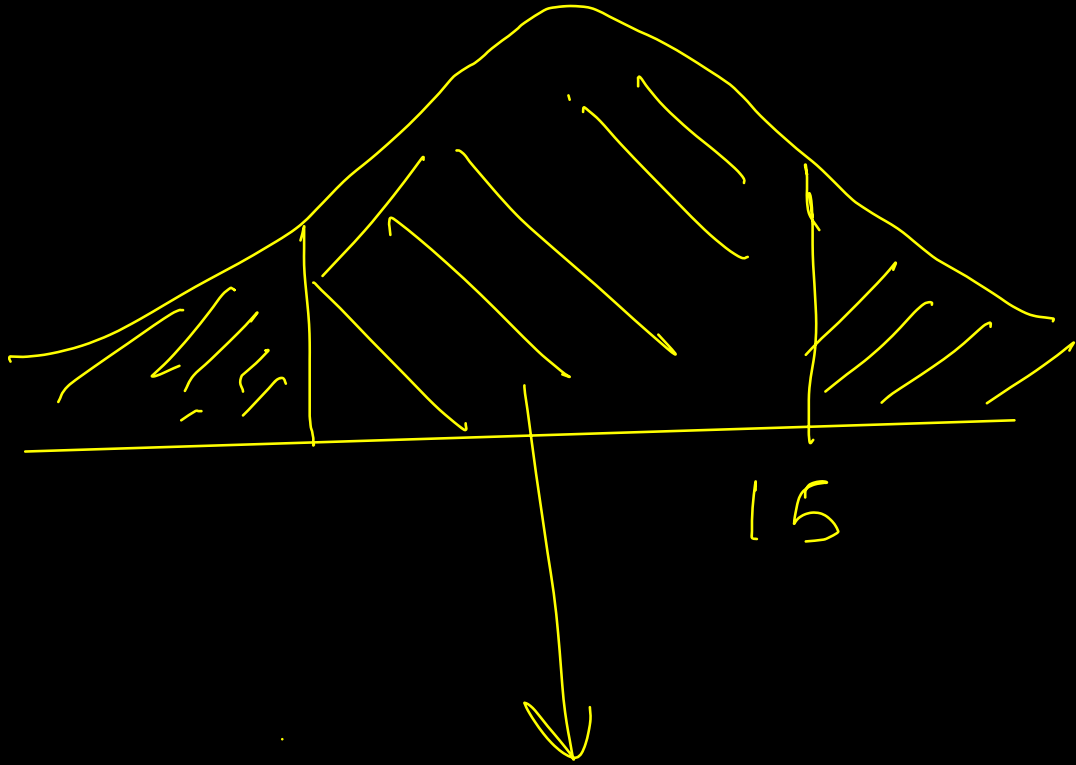
⇓

Some ~~arrange~~ range, within which
if our value falls then
we accept null hypothesis



→ Confidence Interval

$$100 \rightarrow \begin{array}{ccc} 50 & 90H & 10H \\ 50 & 10T & 90T \end{array}$$

$\alpha = .20$

or
20%

$H_1$

$H_1$

tail

tail

10

90.

$$C.I. = 1 - 0.1 - 0.1$$

$$= 80\%$$

15

$85H$
$15T$ } $\Rightarrow$ unfair coin

$\alpha = \cdot 3$

$60H$ $40T$ $\Rightarrow$ unfair