

January 7, 2023

Hypothesis

Recap:

- * Null Hypothesis : H_0
- * Alternative Hypothesis : H_a
- * p-value.
- * Test statistic
- * Significance level: α .
- * Z-test
- * t-test (1-sample).

DSML: CC Fundamentals

$$\mu_w \approx 80$$

$$\mu_e \approx 53.$$

Testing - 3

→ Q] Do fewer people attend class on weekend?



Class begins @ 9:05 p.m.

TERMINOLOGY
EVERYWHERE!

Hypothesis testing Framework.

- 1] Set up the Null and Alternate hypothesis.
- 2] Choose the correct test statistic —
- 3] Select a type of test: left-tailed, right-tailed
or two-tailed.
- 4] Compute the p-value —
- 5] If p-value is less than α , the significance level, then reject the null hypothesis.

Recap



$$\mu = 1800$$
$$\sigma = 100$$

A retailer has 2000 stores in the country

Historical data tells us that weekly sales of shampoo bottles has an average of 1800, with a standard deviation of 100

Sales team wants to improve sales by hiring a marketing team

Hiring a marketing team can be expensive, so we need to be very sure that they will improve sales

Before deploying their strategy for all 2000 stores, they are tested in 50 stores

On the 50 stores, their average sales for that week was 1850

You are the data scientist who should tell your sales team whether this is statistically significant

Sales team has said that we will hire only if we are 99 % confident $\alpha = 0.01$

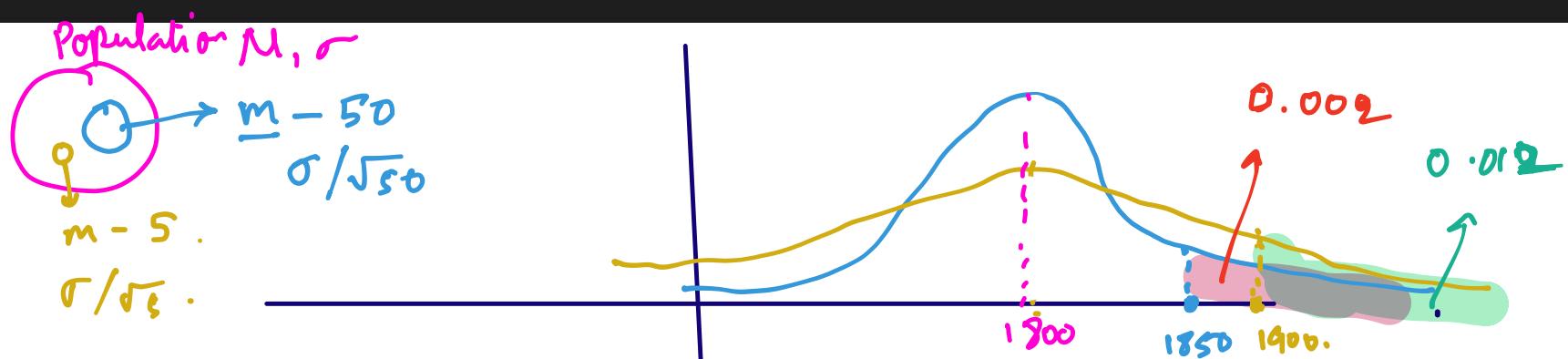
Another marketing team is also being considered

They are tested on 5 stores

On the 5 stores, their average sales for that week was 1900

Would you say this team is better than the first one?

Between the “blue team” and the “yellow team”, whom will you choose?



Recap

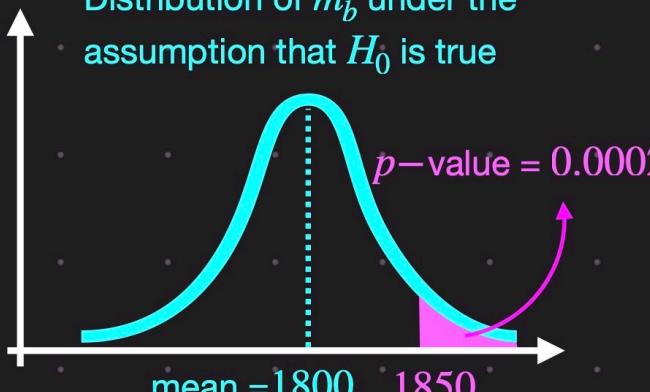
$$\alpha = 0.01$$

50 stores with average of 1850

$$H_0 : \mu_b = 1800$$

$$H_a : \mu_b > 1800$$

Distribution of m_b under the assumption that H_0 is true



$$\text{std dev} = \frac{100}{\sqrt{50}}$$

$$z = \frac{1850 - 1800}{100/\sqrt{50}} = 3.53$$

Reject H_0

$$\mu = 1800$$

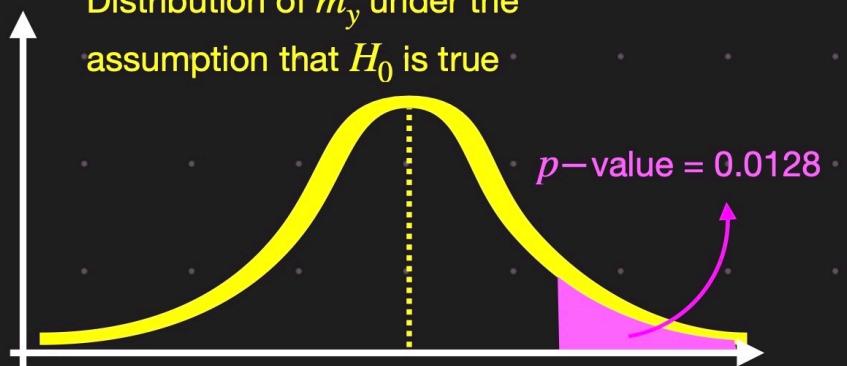
$$\sigma = 100$$

5 stores with average of 1900

$$H_0 : \mu_y = 1800$$

$$H_a : \mu_y > 1800$$

Distribution of m_y under the assumption that H_0 is true



$$\text{std dev} = \frac{100}{\sqrt{5}}$$

$$z = \frac{1900 - 1800}{100/\sqrt{5}} = 2.23$$

Fail to reject H_0

Recap: Sachin Tendulkar's batting.

$$\rightarrow \text{Average} \approx 43 \quad \rightarrow \sigma \approx 42$$

Q1] Does Sachin bat better in the first innings vs. the second?

$$\rightarrow \begin{matrix} \text{first} \approx 46 \\ \text{second} \approx 40 \end{matrix} \rightarrow \text{Not significant}$$

Q2] Does Sachin score more when India wins vs. loses?

$$\begin{matrix} \text{win} \approx 51 \\ \text{loss} \approx 35 \end{matrix} \rightarrow \text{significant}$$

Recap : t - test vs. z - test .
population mean.

$$z = \frac{m - \mu}{\sigma / \sqrt{n}}$$

$$t = \frac{m - \mu}{s / \sqrt{n}}$$

$\rightarrow \sigma \rightarrow$ Population std

$\rightarrow s \rightarrow$ Sample std .

Drug Recovery

suppose 2 companies develop drugs to cure a disease.

Drug 1 was tested on 100 people, and the recovery days look like :

$$\rightarrow [8, 5, 9, 10, \dots, 16], \underline{\mu = 7.1} \text{ days.}$$

Drug 2 was tested on 120 people, and the recovery days look like :

$$- [12, 4, 7, 13, \dots, 8], \underline{\mu = 8.07} \text{ days.}$$

$H_0: \underline{\mu = \bar{\mu}}$. Can we say one drug is better/worse than the other?

$H_a:$

$$\underline{\mu < \bar{\mu}}$$

Drug Recovery : Hypothesis testing . $M_b \rightarrow 1^{\text{st}}$ drug
 $M_y \rightarrow 2^{\text{nd}}$ drug.

1] $H_0 : M_b = M_y$
 $H_a : M_b > M_y$.

" Right - tailed
2 - sample test . "

2] $H_0 : M_b = M_y$
 $H_a : M_b < M_y$

" Left - tailed
2 - sample test "

3] $H_0 : M_b = M_y$
 $H_a : M_b \neq M_y$.

" 2 - Sided
2 - sample test ":

→ 2 - sampled t - test .
test statistic

$$\frac{M_1 - M_2}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$\circ \rightarrow M_1, n_1$ $\circ \rightarrow M_2, n_2$
 $s \rightarrow$ sample std.

IQ across two schools.

Suppose there are two schools competing against each other. Each school says that their students have a higher IQ. So, we conduct a test.

Say the first school had numbers like :

$$\rightarrow [101, 115, 90, \dots, 112, 97], \underline{\mu_1} = \underline{103.7}$$

Say the second school had numbers like :

$$\rightarrow [108, 105, 109, \dots, 111, 98], \underline{\mu_2} = \underline{102.9}$$

- H_a
① $\mu_1 > \mu_2$ → right
② $\mu_1 < \mu_2$ - left
③ $\mu_1 \neq \mu_2 \rightarrow 2$

$$\mu_1 > \mu_2 \rightarrow \text{right}$$

$$\mu_1 < \mu_2 \rightarrow \text{left}$$

$$\mu_1 \neq \mu_2 \rightarrow 2$$

$$H_0: \underline{\mu_1} = \underline{\mu_2}$$

$$\frac{103.7 > 102.9}{\underline{\mu_1} \gg \underline{\mu_2}}$$

1 or 3 are correct.

Summary: 1-sample test vs. 2-sample test.

Supply Chain

1-Sample

Compare avg. sales
to a fixed value.

$$H_0: \mu_b = 1800$$

$$H_a: \mu_b > 1800$$

2-Sample

Directly compare
the two strategies.

$$H_0: \mu_1 = \mu_2$$

$$H_a: \mu_1 \neq \mu_2$$

Drug Recovery

observed $\rightarrow \bar{m}$

$M_p = 8$ days. $H_0: m = 8$.
 \leftarrow compare recovery time to pop. mean.

$$H_a: m < 8.$$

\leftarrow Directly compare the two drugs.
 μ_1, μ_2 . $H_0: \mu_1 = \mu_2$
 $H_a: \mu_1 \neq \mu_2$.

IQ test.

observed $\rightarrow m$

$$\mu_p = 100$$

\leftarrow Compare IQ test results with pop. mean.

$$H_0: m = 100$$

$$H_a: m < 100,
- m > 100$$

\leftarrow Directly compare IQ test results of schools.

$$\mu_1, \mu_2$$

$$H_0: \mu_1 = \mu_2$$

$$H_a: \mu_1 \neq \mu_2,
m_1 > m_2.$$

Sachin Tendulkar's Batting : K S - test .

Youtube ads:

Youtube wants to increase its ad revenue. They decide to put 2 ads instead of 1.

Is this a good move? → For them, yes & for us, no!
can we test this effect?

↓
I sample which
watches 2 ads

* treatment
group.

↓
I sample which
watches 1 ad

* . Watch-time in hours.
↓
of the video.

control

1 sample test



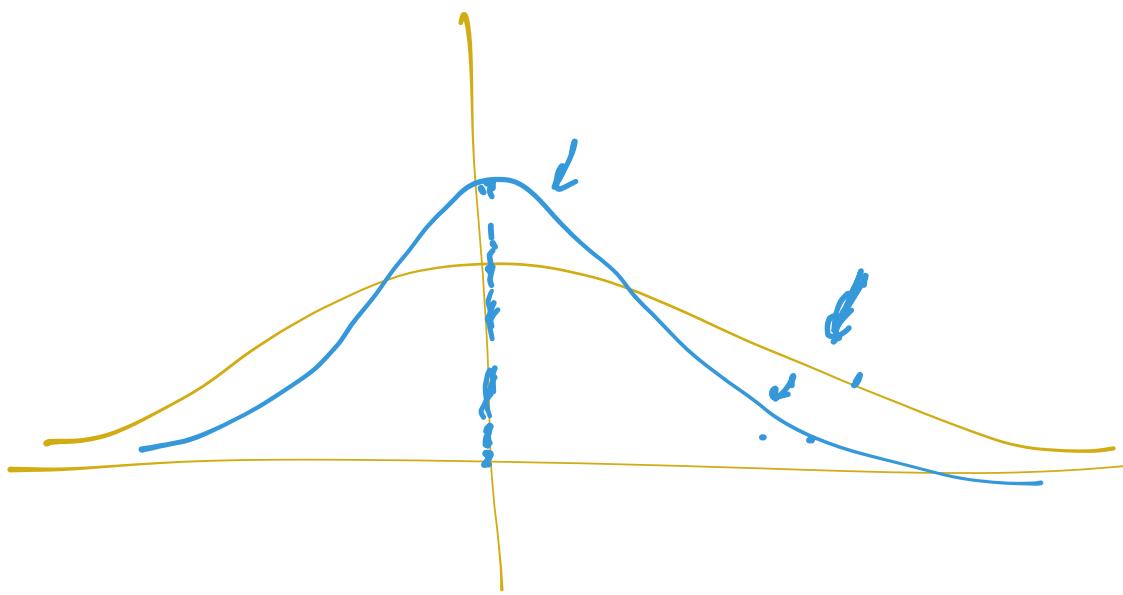
- ① Do this when you have
 - /population mean,
 - /population std.
- ② Compare sample mean to population mean.

2-sample test .



- ① Do this when you have 2 samples and no data about pop. M, σ .

- ② We compare 2 sample means .

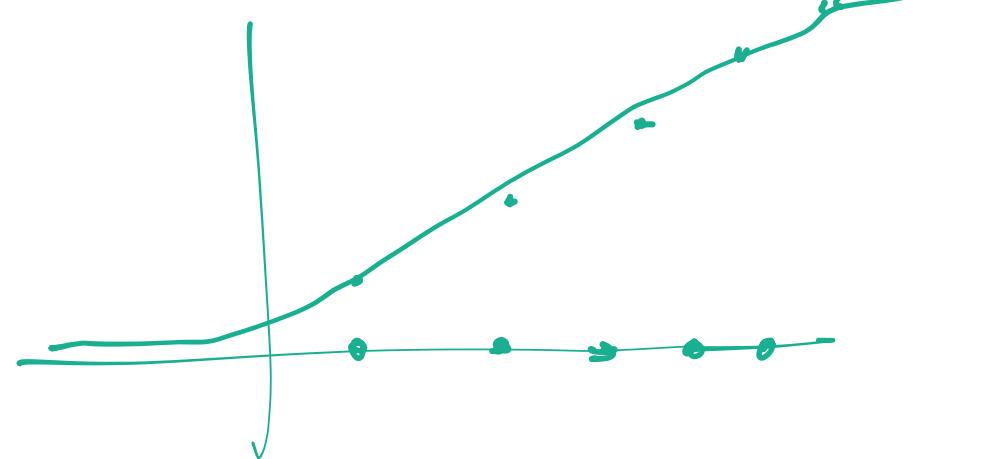
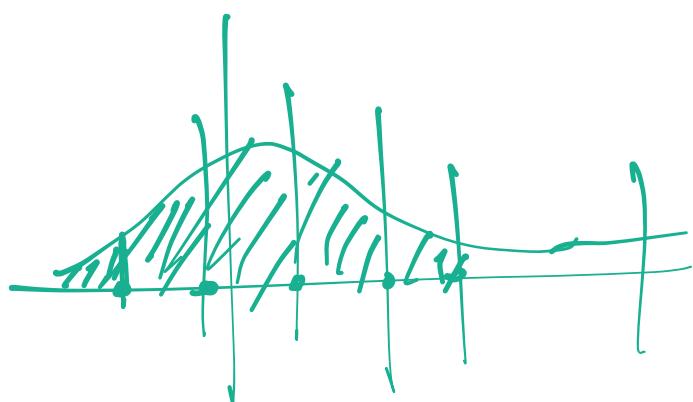


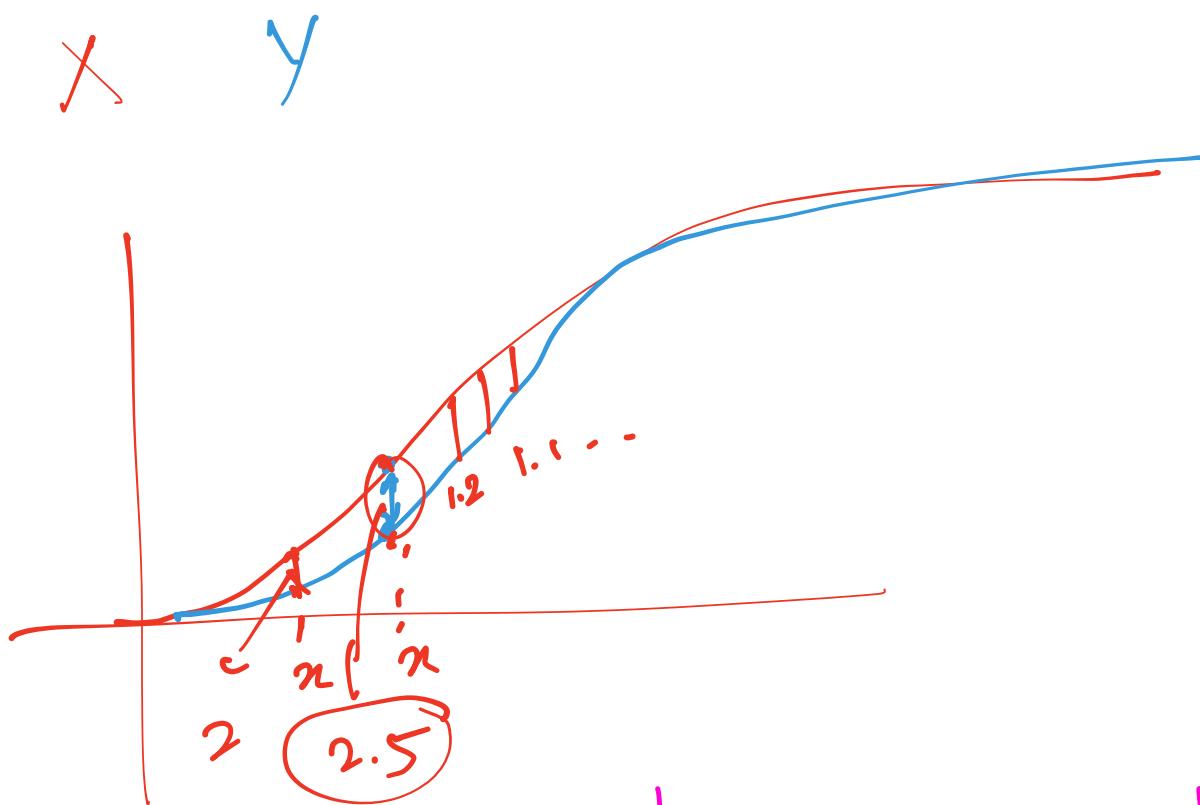
$H_0 : \mu_1 = \mu_2$.
→ 2-sample test.

K-S test → test which
helps us differentiate
between distributions.

K-S test : To check whether 2 Random variables have different distributions or not.

① Calculate the Empirical CDF.





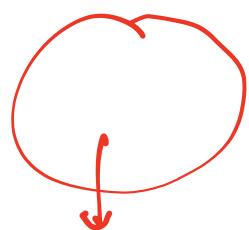
K-S test : $\max \left| \text{cdf}_Y(y) - \text{cdf}_X(x) \right|$

K.S. Test: $X, Y \rightarrow R.V.s$

Do they have the same distribution?

H_0 : Yes, $X \& Y$ have the same dist.

H_a : No, $X \& Y$ do not have the same dist.

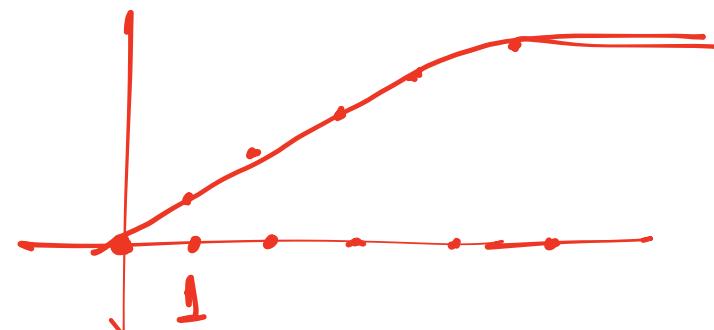


$$a = [1, 2, 3, 4, 5] \leftarrow$$

$$\text{Cum}(a) = 15$$

$$\frac{1}{15}$$

$$, \frac{2}{15}, \frac{3}{15}, \frac{4}{15}, \frac{5}{15}.$$



Recap: Two concepts!

I] 1 sample test vs. 2 sample test :

(a) compares the sample mean to pop. mean

(b) Done when we have pop. μ

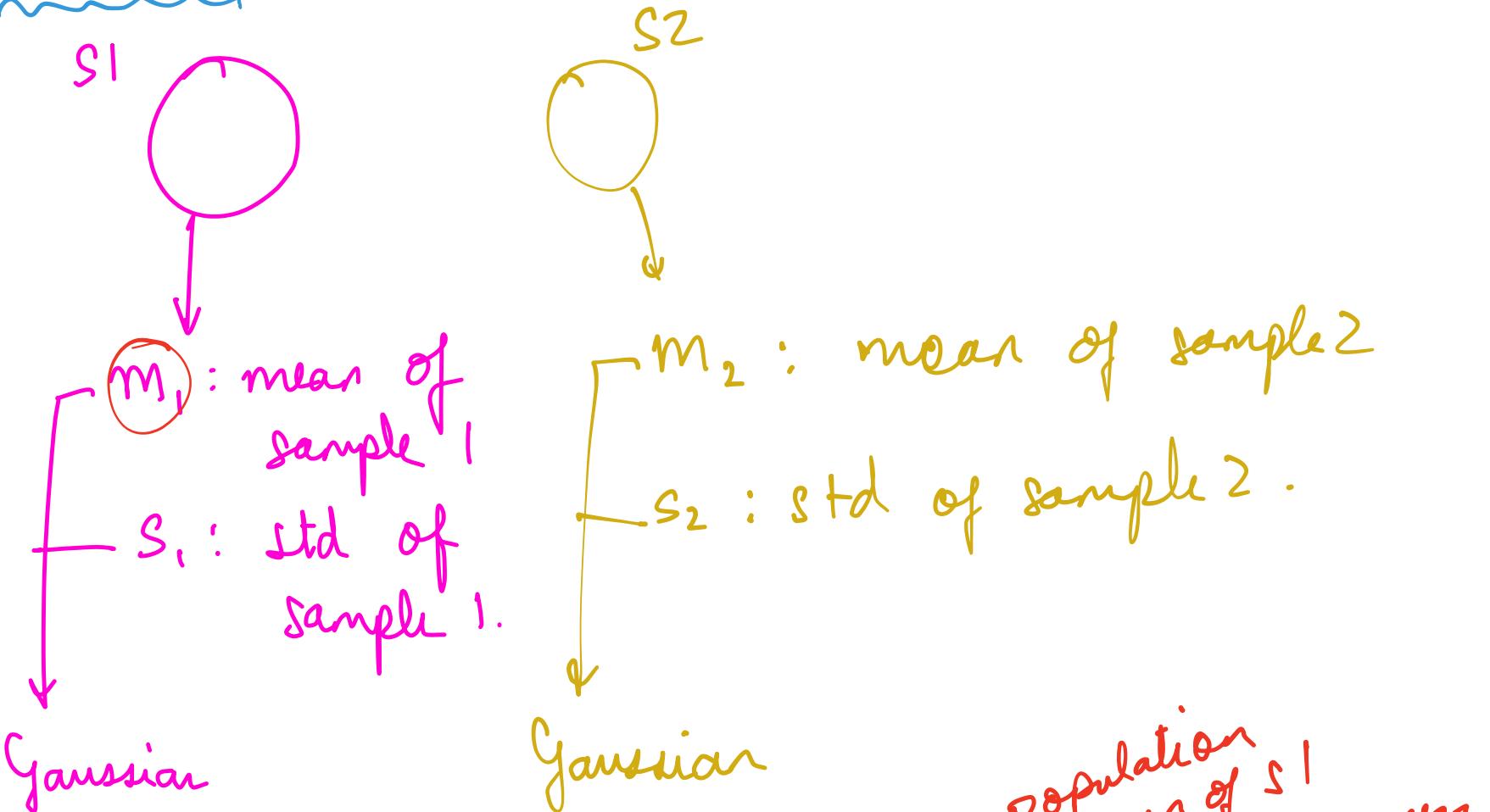
(a) compares 2 sample means.

(b) Done when we only have samples.

II] K-S test : (a) Tests whether the ECDFs of two samples is similar or not.

(b) Null : The distributions are the same.

2 - Sample test :



$$Y = \bar{m}_1 - \bar{m}_2$$

$$E[Y] = E[\bar{m}_1] - E[\bar{m}_2] = \mu_1 - \mu_2$$

Population mean of S_1 Population mean of S_2 .

$$\text{std}(Y) = \sigma_1 / \sqrt{n_1} + \sigma_2 / \sqrt{n_2}$$

$$\text{Var}(x) = \frac{E[x^2]}{n} - (E[x])^2 = \sigma^2.$$

$$\text{Var}(m_1) = E[m_1^2] - (E[m_1])^2 = \sigma_1^2$$

$$\text{Var}(m_1) = \frac{\sigma_1^2}{n} \quad \text{Var}(m_2) = \frac{\sigma_2^2}{n}$$

$$\text{Std}(m_1) = \frac{\sigma_1}{\sqrt{n_1}} \quad \text{Std}(m_2) = \frac{\sigma_2}{\sqrt{n_2}}$$

$$Y = m_1 - m_2$$

$$\begin{aligned}\text{Var}(Y) &= \text{Var}(m_1 - m_2) \\ &= \text{Var}(m_1) + \text{Var}(m_2) \\ &= \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}.\end{aligned}$$

$$\text{std } (\bar{Y}) = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Assume: $\sigma_1^2 = \sigma_2^2 = s^2$.

$$\rightarrow \sqrt{s^2 \left(\frac{1}{n_1^2} + \frac{1}{n_2^2} \right)}$$

test-statistic: $\frac{\bar{m}_1 - \bar{m}_2}{\sqrt{s^2 \left(\frac{1}{n_1^2} + \frac{1}{n_2^2} \right)}}$

