

Probability and Statistics

~

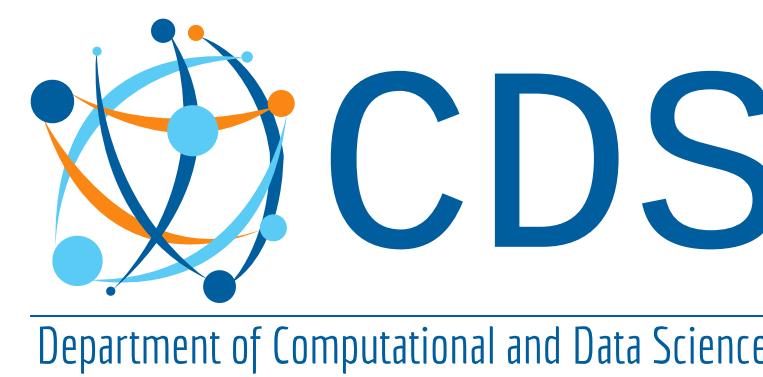
Module 1: revision

- Information.
- testing sheet.
- exp, β, γ .

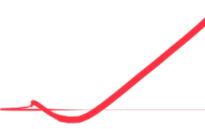
Konduri Aditya

Department of Computational and Data Sciences
Indian Institute of Science, Bengaluru

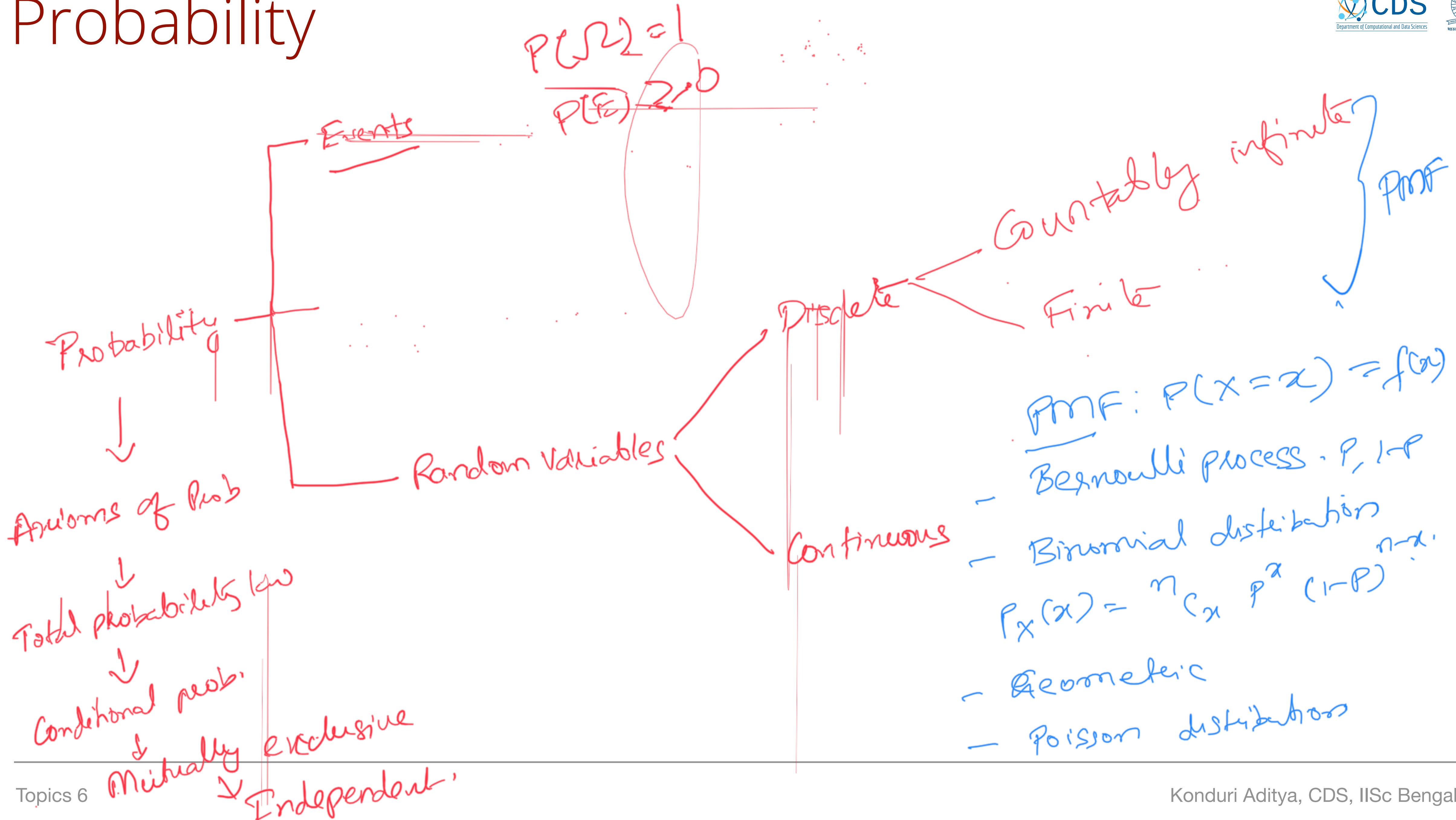
Konduriadi @ iisc.ac.in



Outline

- ▶ Probability 
- ▶ Bayesian Inference
- ▶ Hypothesis testing 
- ▶ Overall summary

Probability



Probability

Continuous R.V: $X \rightarrow$ Any real value. $(-\infty, \infty)$

PDF $f(x=x) = f_X(x) = g(x)$

- We obtain the probability associated with different interval.

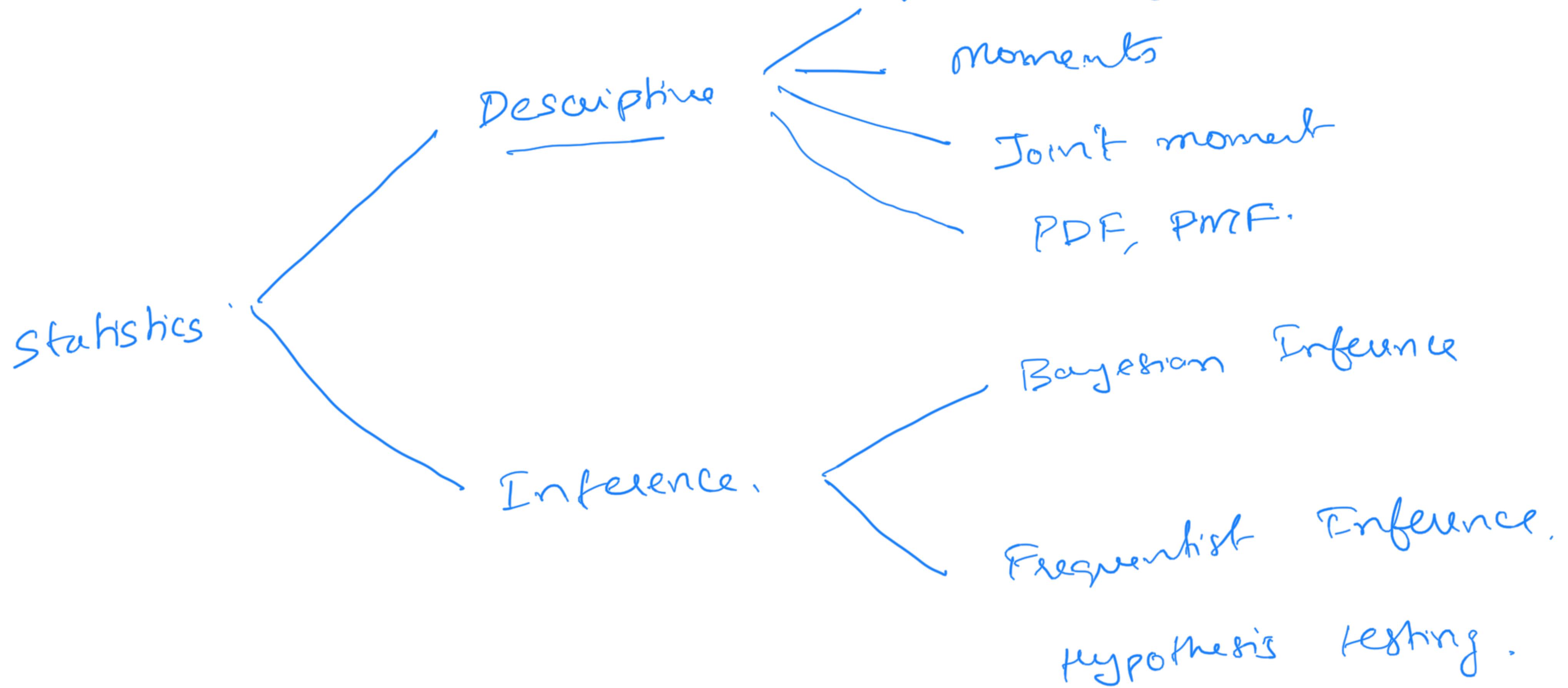
$$P_X(a \leq x \leq b) = \int_a^b f_X(x) dx.$$

- cumulative distribution function

$$CDF(x) = P_X(X \leq x) = \int_{-\infty}^x f_X(t) dt.$$

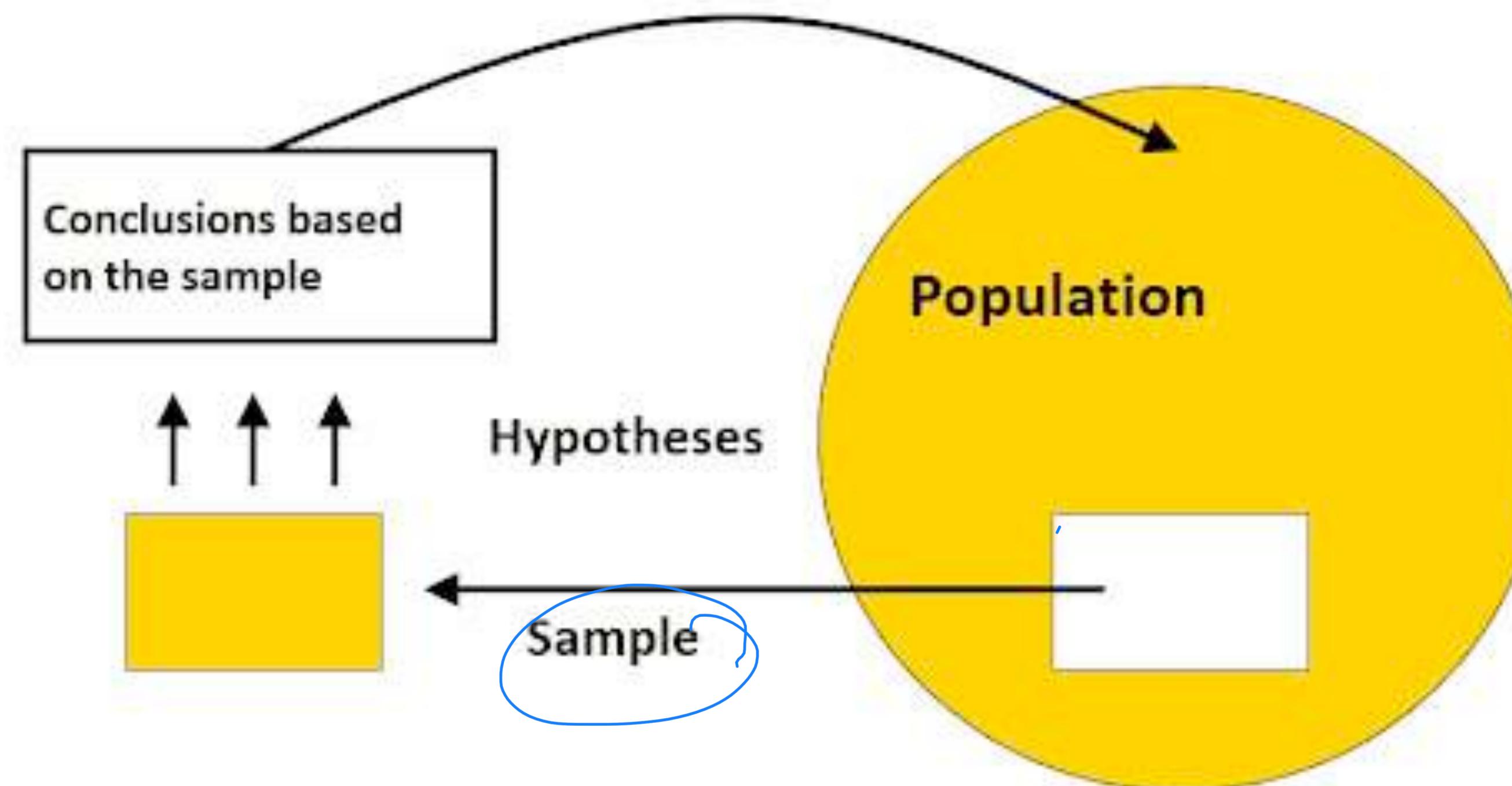
- Normal distribution, exponential distribution, Beta distribution

Statistics



The idea of statistical inference

Generalisation to the population



Source: howmed.net

Bayesian statistical inference

- ▶ Bayes' rule is the key:

$$P(\text{hypothesis is true} | \text{data}) = \frac{P(\text{data} | \text{hypothesis is true}) P(\text{hypothesis is true})}{P(\text{data})}$$

↓ Posterior
 Likelihood
 Prior

$$P(H | D) = \frac{P(D | H) P(H)}{P(D)}$$

↓
~~P(D)~~

Discrete priors

- ▶ A certain disease has a prevalence of 0.005
Prior
- ▶ A screening test has 2% false positives and 1% false negatives
- ▶ Suppose a patient is screened and has a positive test
 - 1. Represent this information with a tree and use Bayes' theorem to compute the probabilities the patient does and doesn't have the disease
 - 2. Identify the data, hypotheses, likelihoods, prior probabilities and posterior probabilities
 - 3. Make a full likelihood table containing all hypotheses and possible test data

Discrete priors

1. Represent this information with a tree and use Bayes' theorem to compute the probabilities the patient does and doesn't have the disease

H_+ : patient has the disease, H_- : patient does not have the disease

T_+ : patient tests positive, T_- : patient tests negative

$$P(H_+ | T_+) = \frac{P(H_+) P(T_+ | H_+)}{P(H_+) P(T_+ | H_+) + P(H_-) P(T_+ | H_-)} =$$

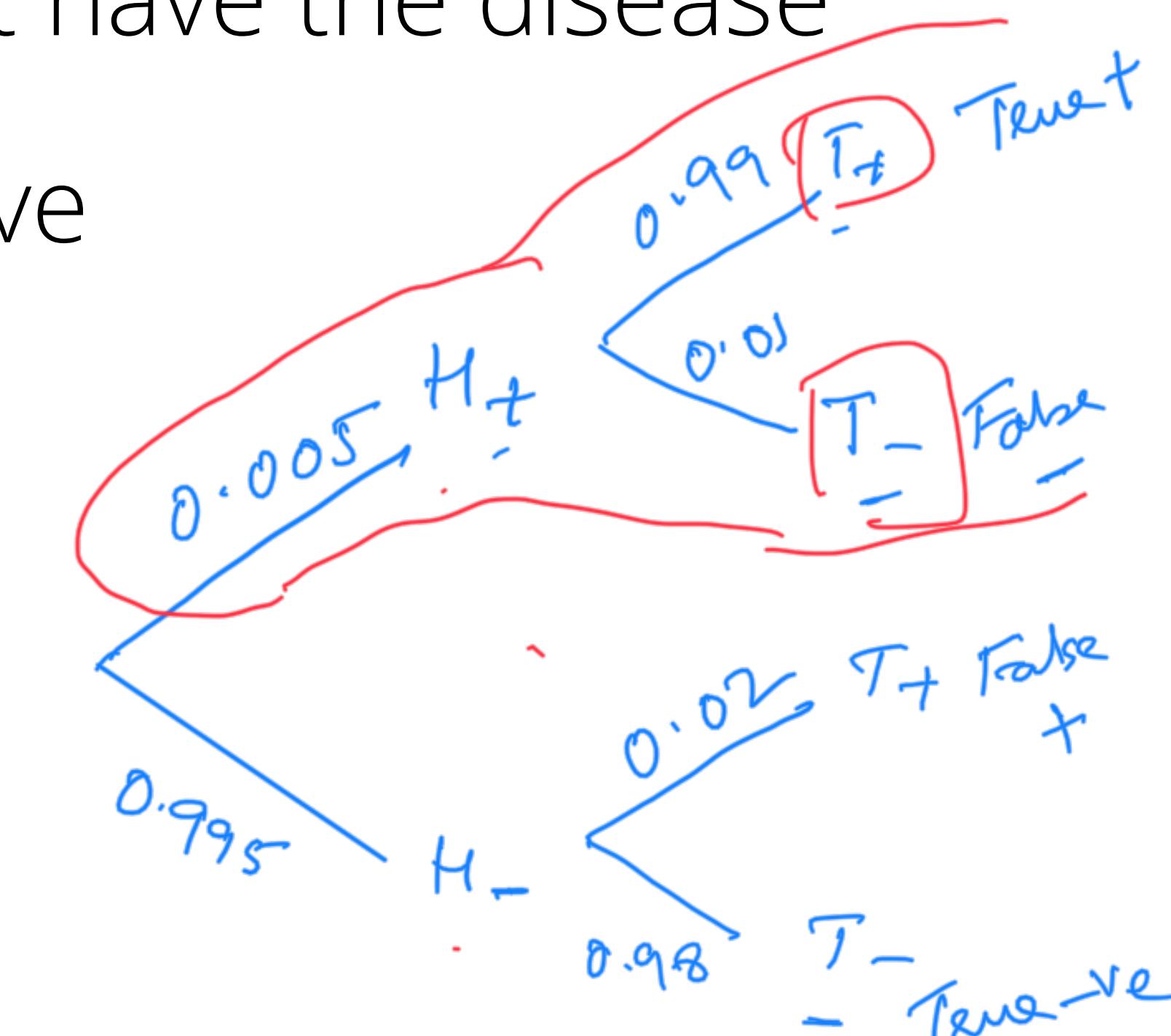
$$P(H_- | T_+) =$$

Posterior
 ↓
 Prior.
 ↓
 likelihood
 ↗

$P(H_+) P(T_+ | H_+)$
 ~~$P(H_+) P(T_+ | H_+)$~~
 ~~H_+~~
 ~~T_+~~
~~True~~
~~+~~

$P(H_-) P(T_+ | H_-)$
 ~~$P(H_-) P(T_+ | H_-)$~~
 ~~H_-~~
 ~~T_+~~
~~False~~
~~-~~

$TP \quad TN \quad FP \quad FN$



Discrete priors

Data: The data are the results of the experiment. In this case, the positive test

Hypotheses: The hypotheses are the possible answers to the question being

asked. In this case they are H_+ the patient has the disease; H_- they don't

Likelihoods: The likelihood given a hypothesis is the probability of the data given that hypothesis. In this case there are two likelihoods, one for each hypothesis

$$P(T_+ | H_+) = 0.99 \quad \text{and} \quad P(T_+ | H_-) = 0.02$$

=

Likelihood is the probability given the hypothesis

Discrete priors

Prior probabilities of the hypotheses: The priors are the probabilities of the hypotheses prior to collecting data. In this case,

$$P(H_+) = 0.005 \quad \text{and} \quad P(H_-) = 0.995$$

Posterior probabilities of the hypotheses: The posteriors are the probabilities of the hypotheses given the data. In this case

$$P(H_+ | T_+) = 0.199 \quad \text{and} \quad P(H_- | T_+) = 0.801$$

$$P(H_+ | T_+) = \frac{P(H_+)P(T_+ | H_+)}{P(T_+)}$$

Discrete priors

3. Full likelihood table

The table holds likelihoods $P(D | H)$ for every possible hypothesis and data combination

hypothesis \mathcal{H}	likelihood $P(\mathcal{D} \mathcal{H})$	
disease state	$P(\mathcal{T}_+ \mathcal{H})$	$P(\mathcal{T}_- \mathcal{H})$
\mathcal{H}_+	0.99	0.01
\mathcal{H}_-	0.02	0.98

Beta prior: problem

- Treatment has prior $f(\theta) = \underline{\underline{\text{Beta}(5,5)}}$

1. Suppose you test it on 10 patients and have 6 successes. Find the posterior distribution on $\underline{\underline{\theta}}$. Identify the type of the posterior distribution
2. Suppose you recorded the order of the results and got SSSFFSSSFF. Find the posterior based on this data

Beta prior: solution

- ▶ Prior PDF $f(\theta) = Beta(5,5)$ = $\frac{9!}{4!4!} \theta^4(1-\theta)^4 = c_1 \theta^4(1-\theta)^4$ $0 < \theta \leq 1$ ✓
 - ▶ Hypothesis: θ Continuous R.V.
 - ▶ Likelihood: $P(Data | \theta) = {}^{10}C_6 \theta^6(1-\theta)^4$ → Binomial distribution
 - ▶ Posterior PDF:
- $$f(\theta | Data) = \frac{c_1 \theta^4(1-\theta)^4 \cdot {}^{10}C_6 \theta^6(1-\theta)^4}{\int_0^1 c_1 t^4(1-t)^4 \cdot {}^{10}C_6 t^6(1-t)^4 dt} = Beta(11,9)$$
- ~~c_1~~ ~~\int_0^1~~ ~~$t^4(1-t)^4$~~ ~~${}^{10}C_6$~~ ~~$t^6(1-t)^4 dt$~~
- $\downarrow f(\theta)$ $\downarrow P(Data | \theta)$ Beta - Binomial \Rightarrow conjugate pair.
- $y = ax + b$
-

Hypothesis testing

- ▶ H_0 : the null hypothesis. This is the default assumption for the model generating the data
- ▶ H_A : the alternative hypothesis. If we reject the null hypothesis we accept this alternative as the best explanation for the data
- ▶ X : the test statistic. We compute this from the data

Ex:

H_0 : On an avg children watch 3 shows of TV per day.

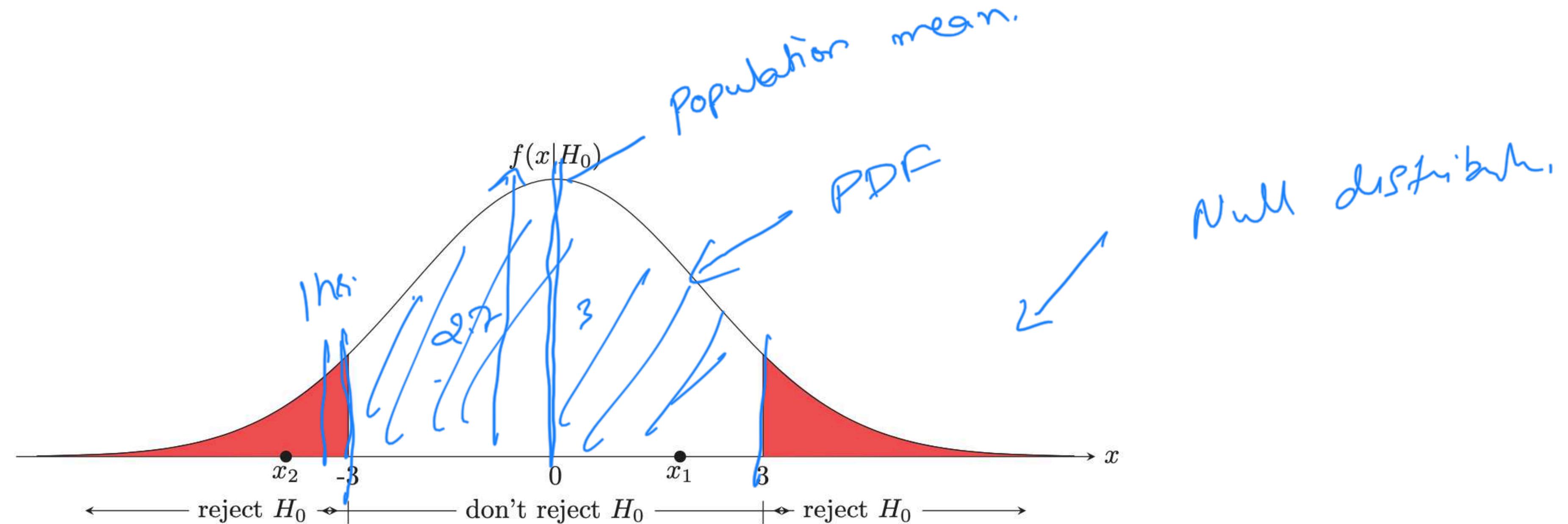
H_0 : $\mu = 3$

H_A : $\mu \neq 3$

$\bar{X} \rightarrow$ computed from the sample : 2.9
2.5

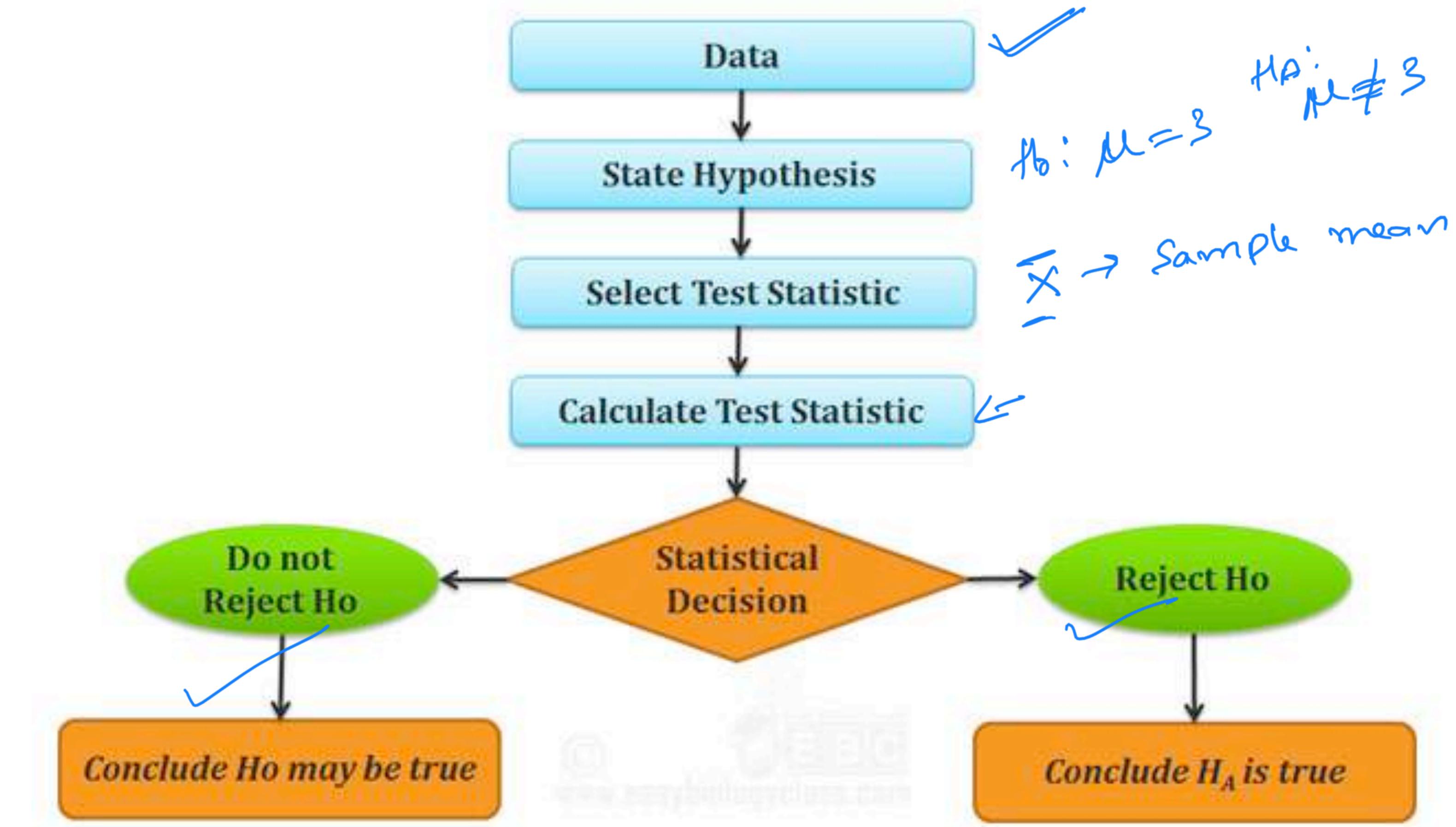
Hypothesis testing

- ▶ Null distribution: the probability distribution of X assuming $\underline{H_0}$
- ▶ Rejection region: if X is in the rejection region we reject H_0 in favor of H_A
- ▶ Non-rejection region: the complement to the rejection region



Hypothesis testing

STEPS IN HYPOTHESIS TESTING



Source: <https://medium.com/analytics-vidhya/hypothesis-testing-steps-235d2670cad4>

Examples of hypothesis

Significance level $\alpha: 0.05$

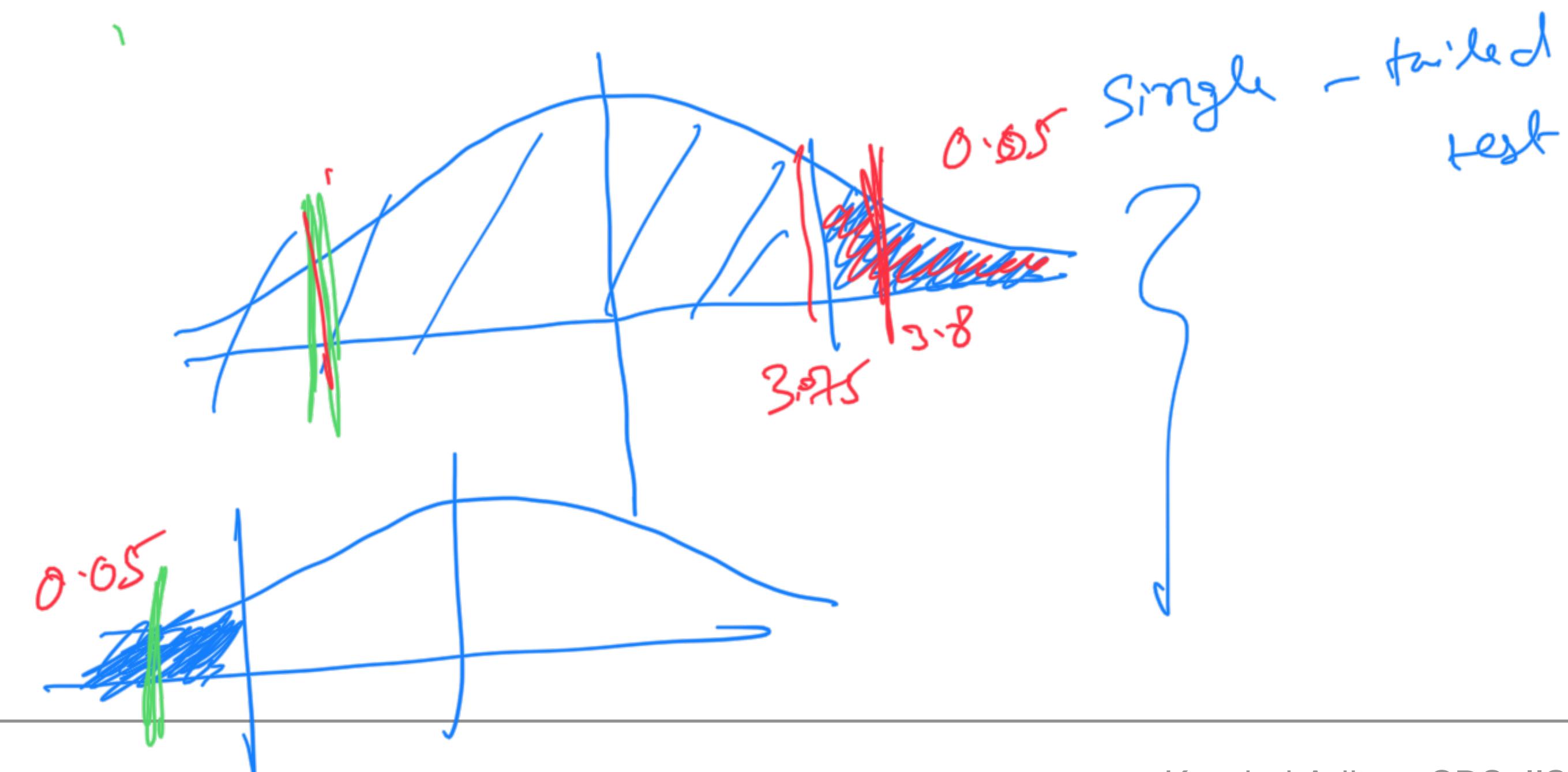
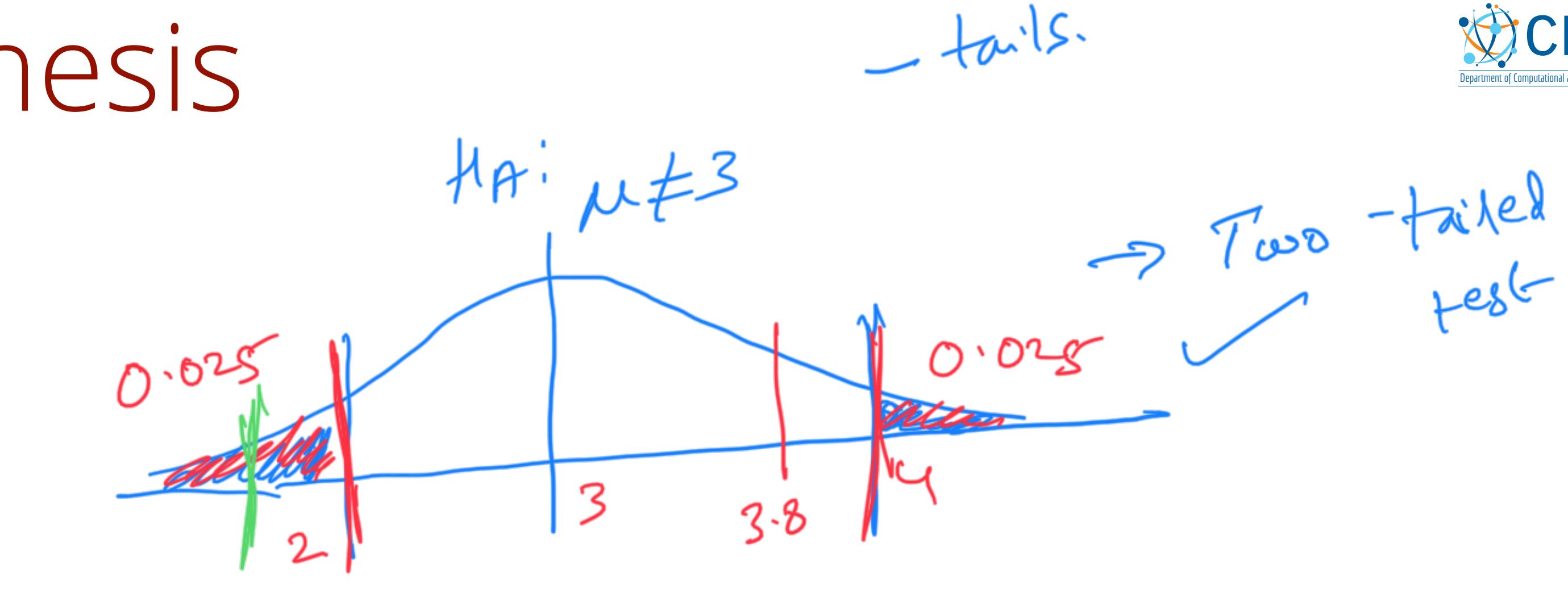
$$H_0: \mu = 3$$

$$H_A: \mu \neq 3$$

$$H_0: \mu = 3$$

$$H_A: \mu > 3$$

$$H_A: \mu < 3$$



Types of errors, significance level, power

- ▶ Type I: false rejection of H_0
- ▶ Type II: false non-rejection ('acceptance') of H_0

		True state of nature	
		H_0 ✓	H_A ✓
Our decision	Reject H_0 ✓	Type I error ↗	correct decision ↘
	'Don't reject' H_0 ✓	correct decision ↘	Type II error ↗

Significance level = $P(\text{reject } H_0 | H_0)$ *ideal case 0*
 = probability we incorrectly reject H_0
 = $P(\text{type I error})$.

Power = probability we correctly reject H_0
 = $P(\text{reject } H_0 | H_A)$ *ideal case 1*
 = $1 - P(\text{type II error})$. = ⊕

Ideal values??

P-value

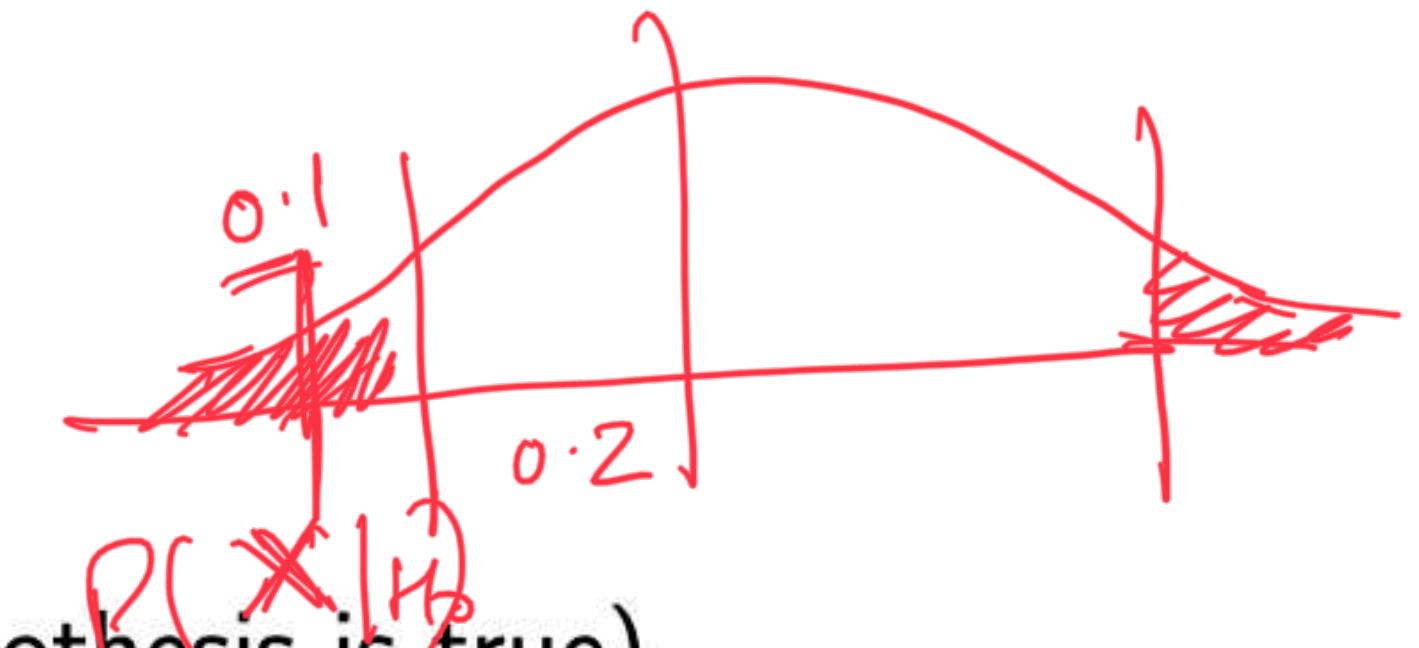
Hypotheses: H_0, H_A .

Test statistic: value: x , random variable X .

Null distribution: $f(x|H_0)$ (assumes the null hypothesis is true)

Sides: H_A determines if the rejection region is one or two-sided.

Rejection region/Significance: $P(x \text{ in rejection region} | H_0) = \alpha$.



The p -value is a computational tool to check if the test statistic is in the rejection region. It is also a measure of the evidence for rejecting H_0 .

p-value: $P(\text{data at least as extreme as } x | H_0)$

Data at least as extreme: Determined by the sided-ness of the rejection region.

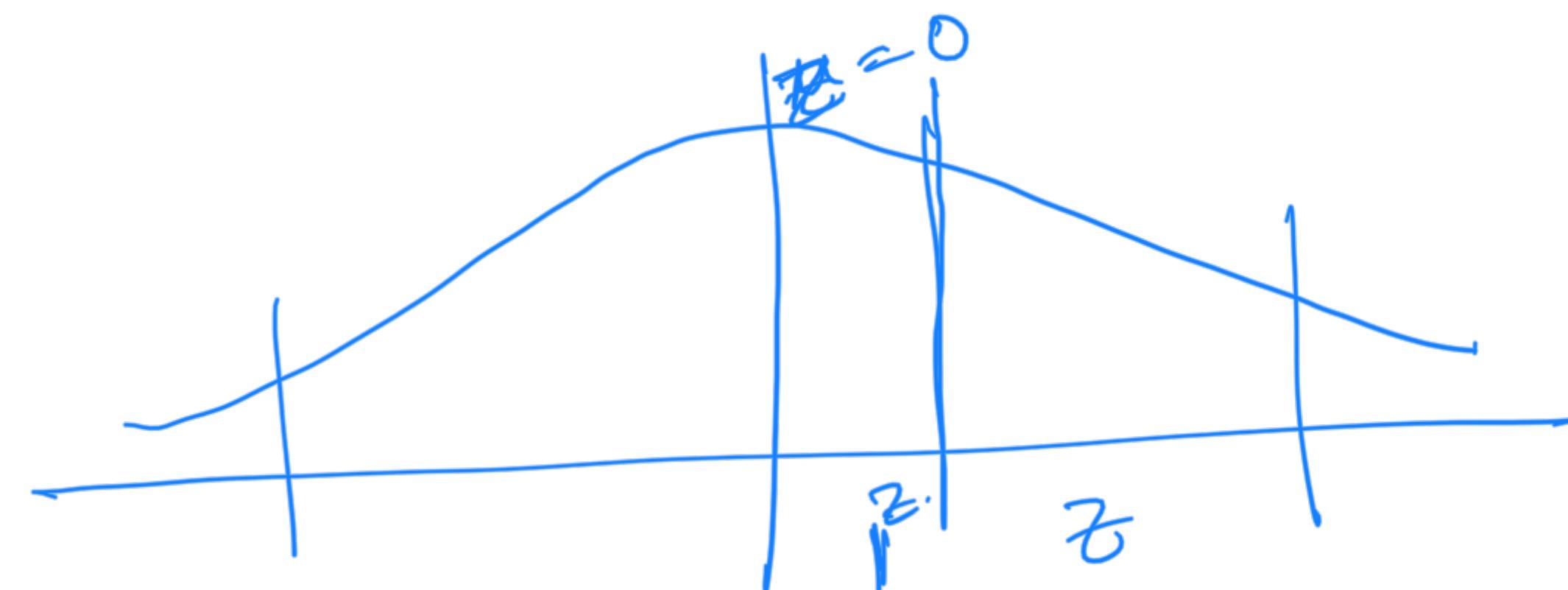
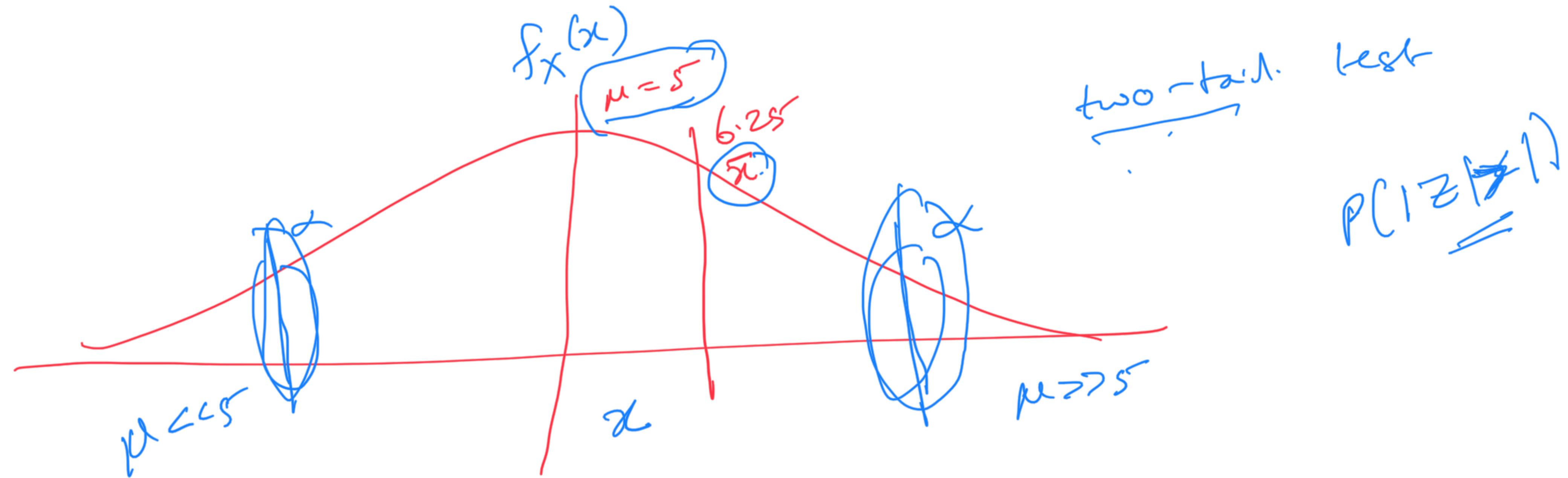
Screenshot from MIT OCW 18.05 slides

Example

- H_0 : data follows a $\underline{N(5, 10^2)}$
- H_A : data follows a $\underline{N(\mu, 10^2)}$ where $\mu \neq 5$.
- Test statistic: $z = \text{standardized } \bar{x}$ → *sample mean*
- Data: 64 data points with $\bar{x} = 6.25$.
- Significance level set to $\alpha = 0.05$.

- (i) Find the rejection region; draw a picture.
- (ii) Find the z -value; add it to your picture.
- (iii) Decide whether or not to reject H_0 in favor of H_A .
- (iv) Find the p -value for this data; add to your picture.
- (v) What's the connection between the answers to (ii), (iii) and (iv).

Screenshot from MIT OCW 18.05 slides



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

$$= \frac{6.25 - 5}{5 / \sqrt{4}} = 0$$

$$\bar{z} = \frac{6.25 - 5}{5 / \sqrt{4}}$$

Example: solution

The null distribution $f(z | H_0) \sim N(0, 1)$

$$\alpha = 0.05$$

- (i) The rejection region is $|z| > 1.96$, i.e. 1.96 or more standard deviations from the mean.

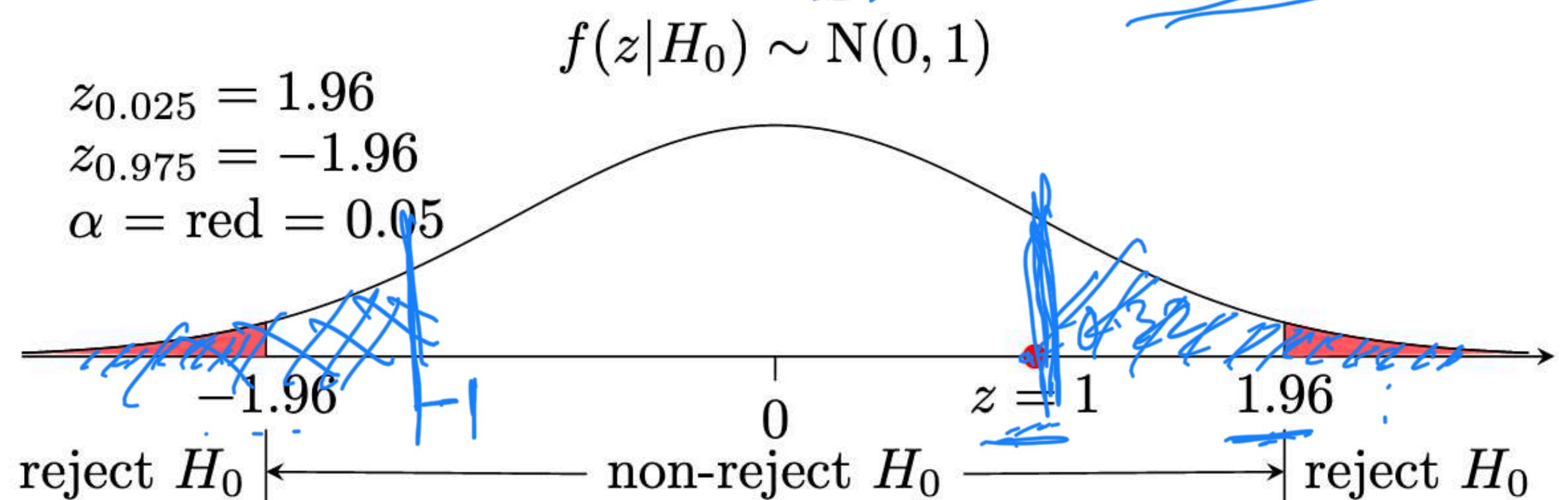
(ii) Standardizing $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{1.25}{1.25} = 1$.

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

- (iii) Do not reject since z is not in the rejection region

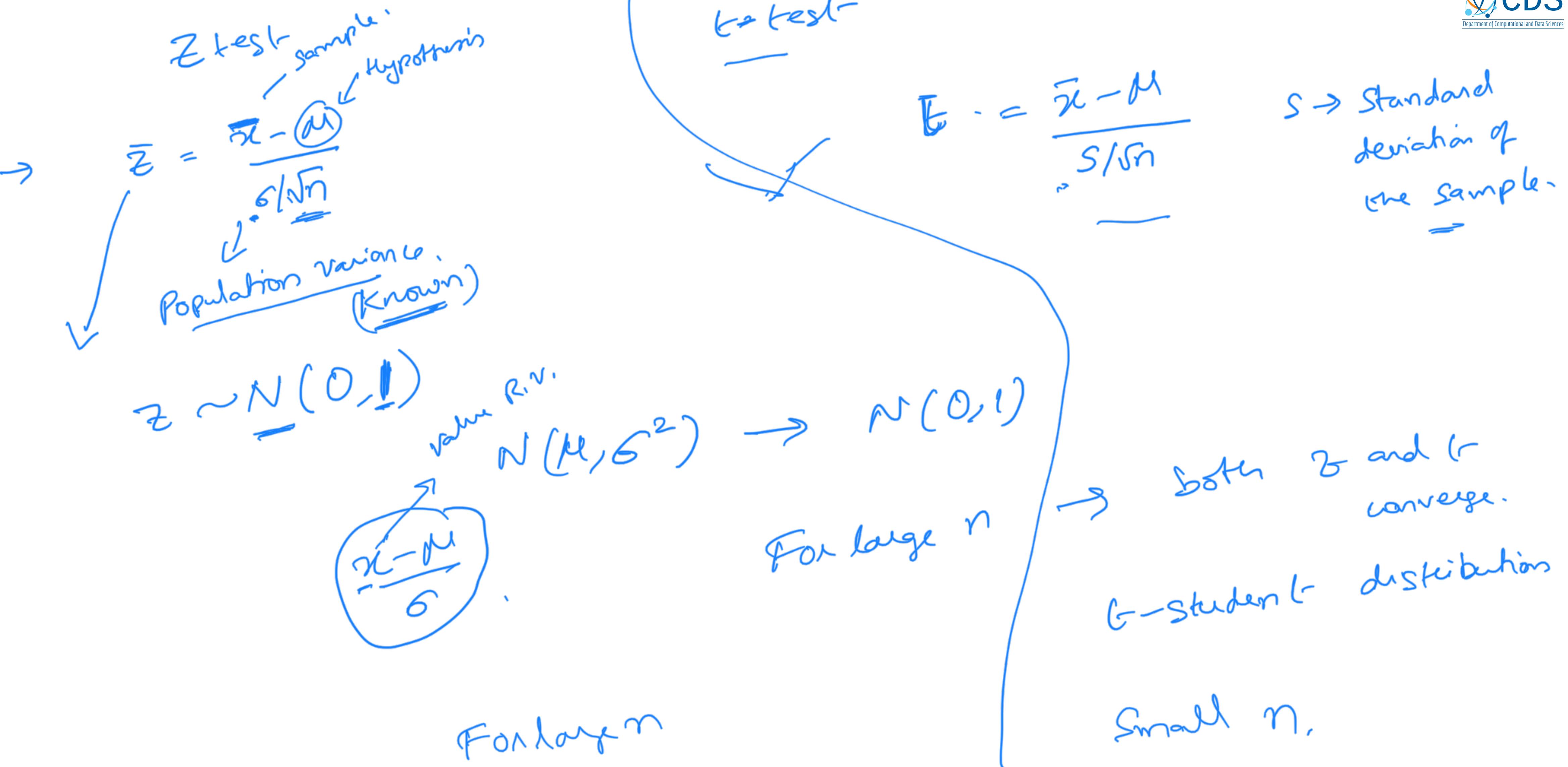
- (iv) Use a two-sided p -value $p = P(|Z| > 1) = .32 > 0.05$

$\cancel{z \text{-table}}$

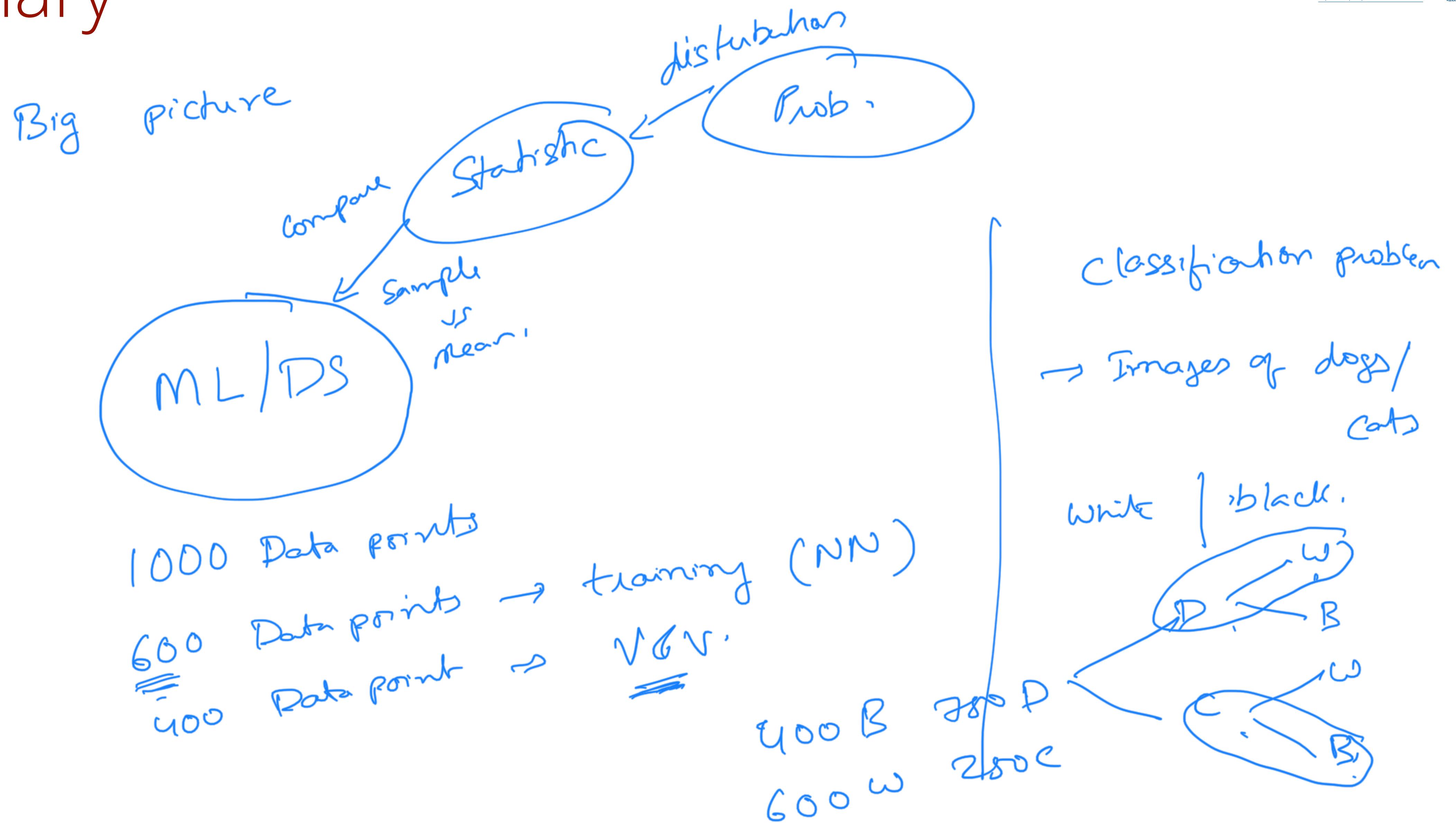


- (v) The z -value is not in the rejection region tells us exactly the same thing as the p -value being greater than the significance, i.e. don't reject the null hypothesis H_0 .

Screenshot from MIT OCW 18.05 slides



Summary



Training : 600

WD + BC

BD / WC

Testing :

BD
 \downarrow
 Input image

Falsely \rightarrow BC

Compare the diversity between the sample and the population

→ Use hypothesis test.
 \downarrow
 based on a particular parameter.
 \downarrow
Distribution \rightarrow Null distribution

$$y = \hat{a}x + \hat{b}$$

→ Ex: Bayesian linear regression.

→ Bayesian NN.
⋮

Question Paper

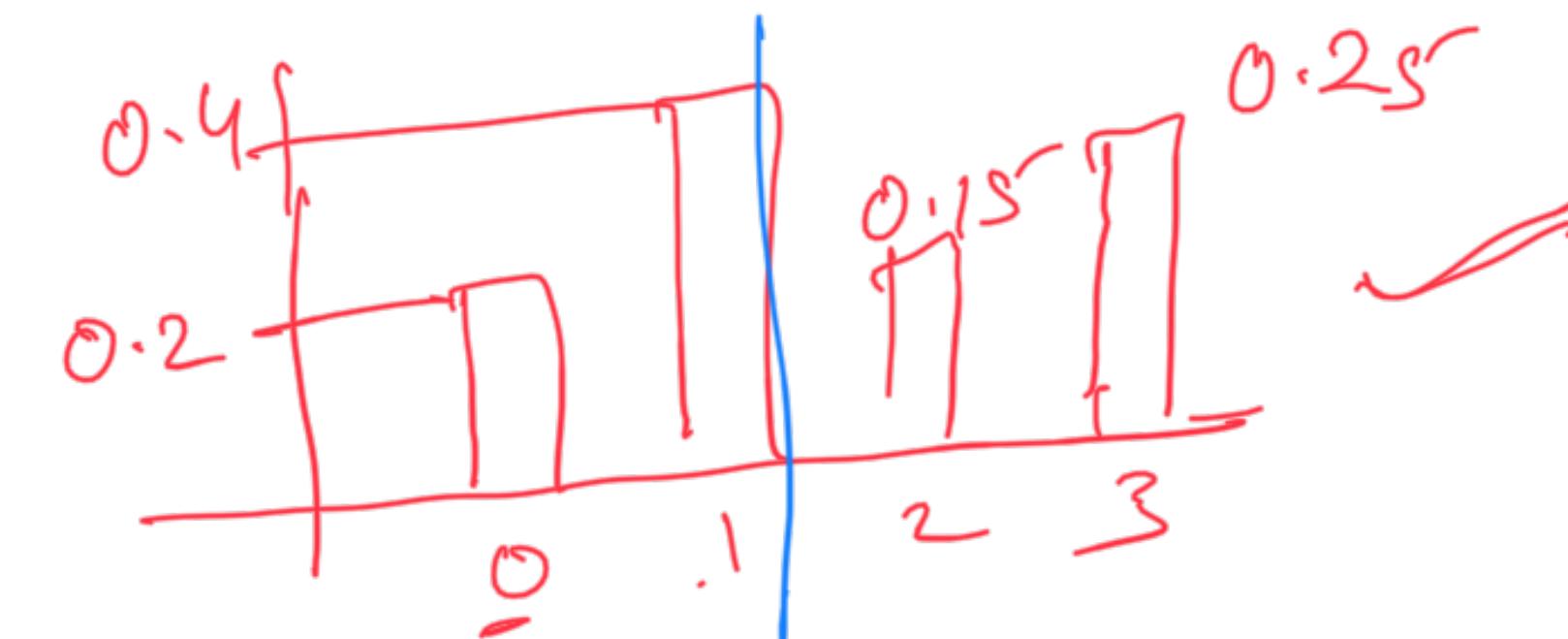
① T/F:

② MCQ

③ Descriptive Questions

Moment Generating function

Discrete R.V. $X \rightarrow$ PMF: $P_X(x) = g(x)$



~~Mean~~ + Mean = $\sum_x x \cdot P_X(x)$

Second moment = $\sum_x x^2 P_X(x)$ MGF

n^{th} moment = $\sum_x x^n P_X(x)$