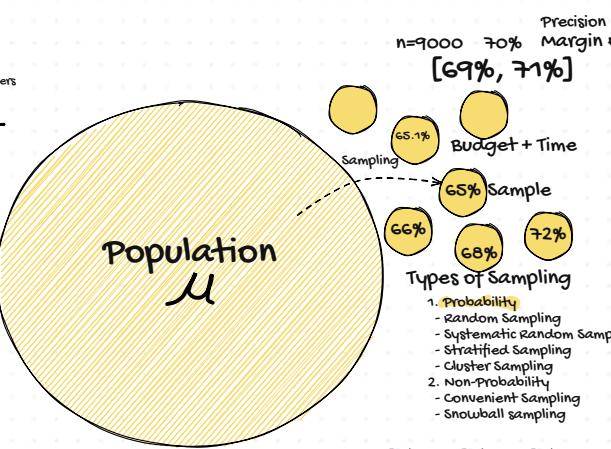
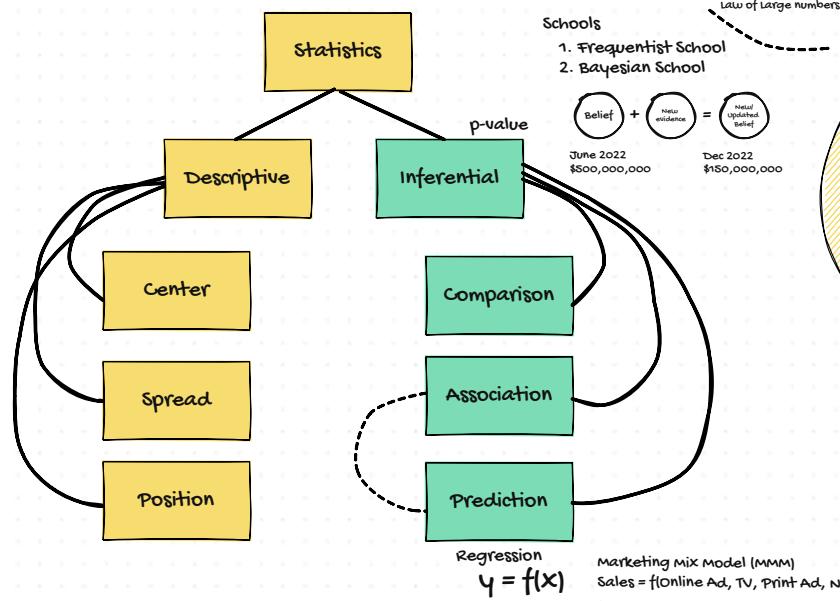
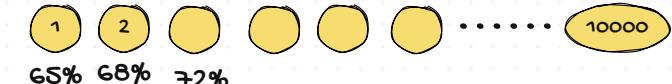


# Intro to Statistics for Data Analyst

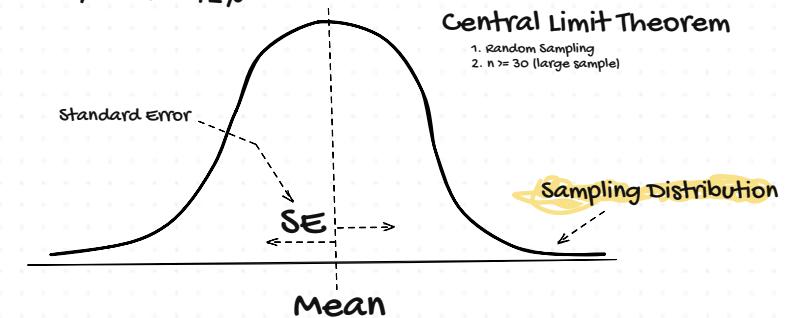
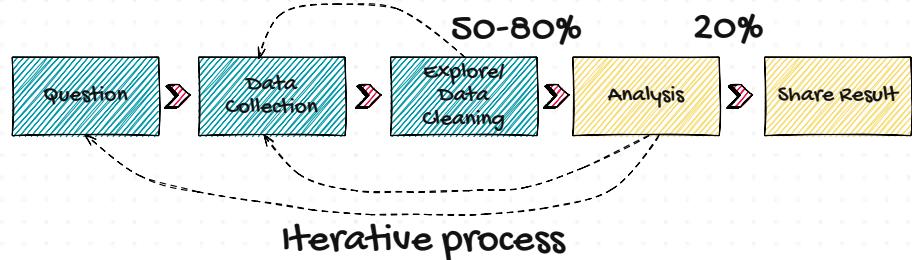


1 → 2 → 3

## Sampling variation



## Statistics Process



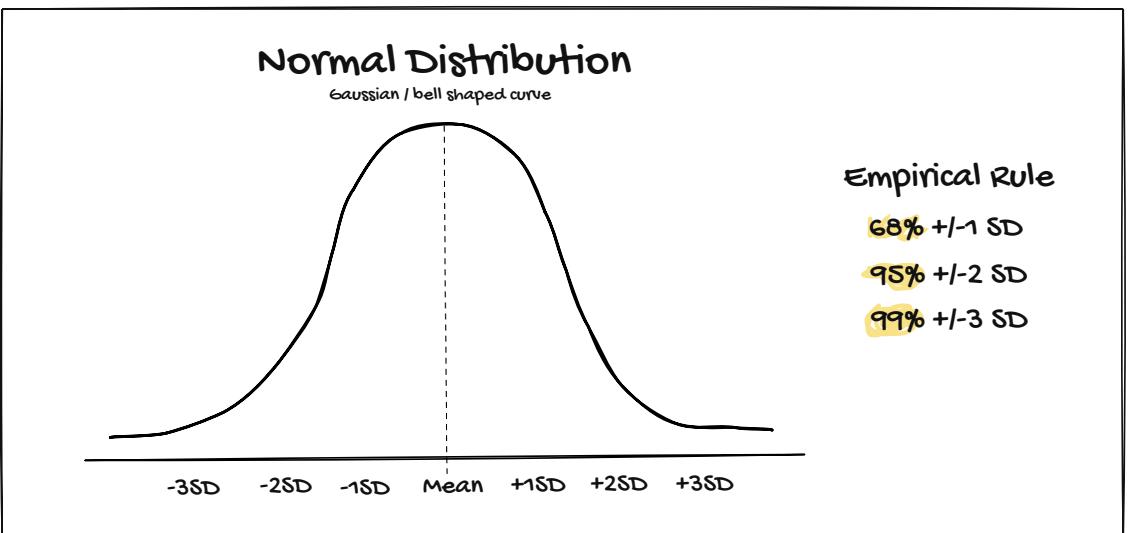
## Center Central Tendency

- Mean
- Median
- Mode

## Spread variation

- Range
  - SD
  - Variance
  - IQR
- Min
  - Max
  - Percentile
  - Quartile

## Position how to find the data



# Agenda

Morning - Concepts

Break 1 hour

Afternoon - Sheets/ R

## Spread

variability/ variation

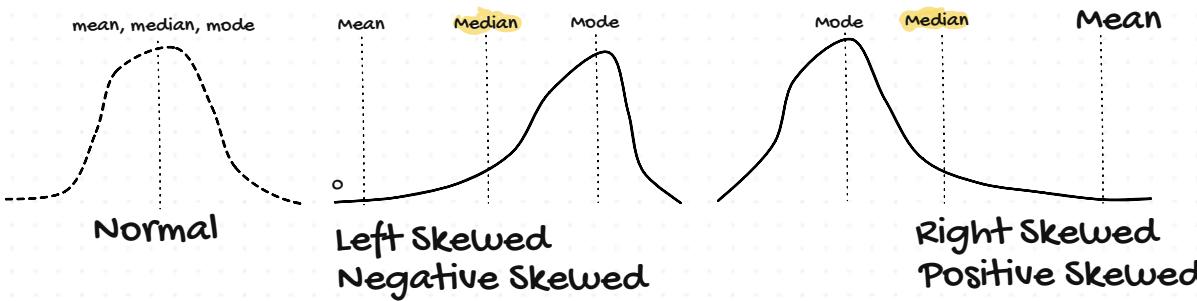
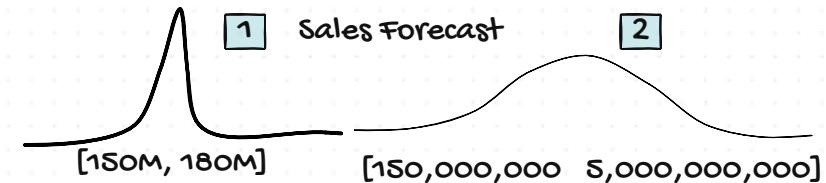
- Range (max- min)
- Variance
- Standard Deviation
- IQR (Interquartile Range)

20, 25, 50, 100, 200

$$\text{Range} = 200-20 = 180$$

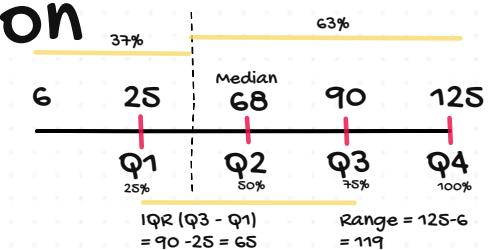
$$\text{var} = \sum ((x - \bar{x})^2) / (n-1)$$

$$\text{sd} = \sqrt{\sum ((x - \bar{x})^2) / (n-1)}$$



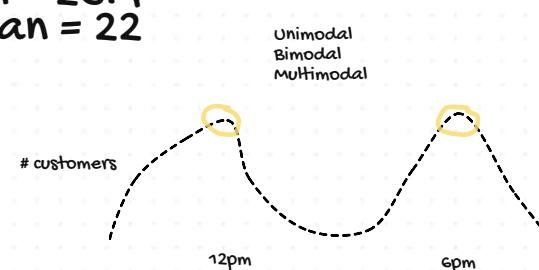
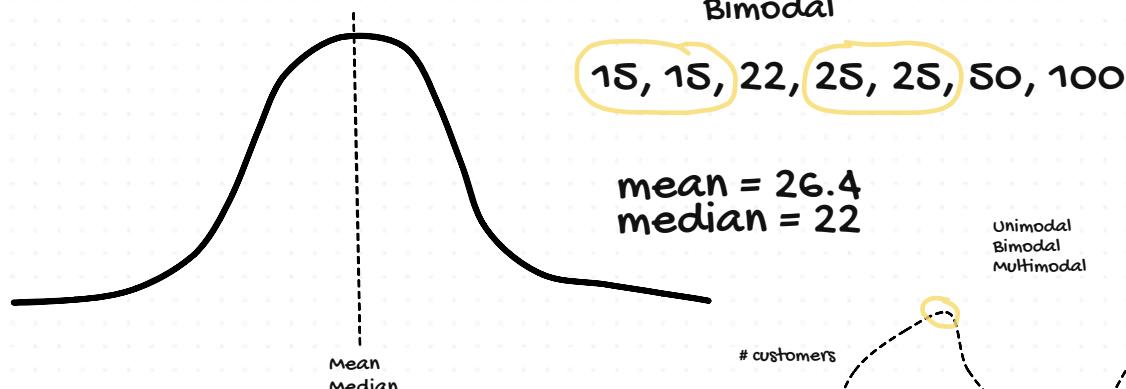
## Position

Min  
Max  
Percentile  
Quartile



## Median < Mean

### Central Tendency



# OUTliers (extreme values)

Boxplot + 5 numbers summary

$$20 - 1.5 * (40)$$

$$Q1 - 1.5 * IQR$$

IQR

$$60 + 1.5 * (40)$$

$$Q3 + 1.5 * IQR$$



Q1  
**20**

Median

Q3  
**60**

X  
Mean

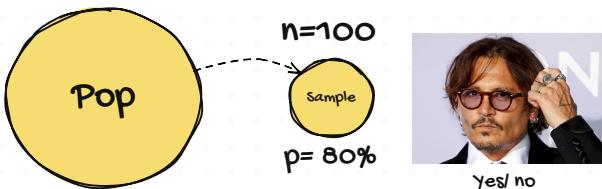


# Confidence Interval (CLT)

Results:  
80% [72%, 88%] CI 95%

Steps to Apply CLT

1. collect data
2. SE
3.  $ME = SE * 1.96$
4. Confidence Interval

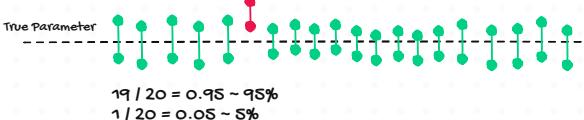


Numbers in each row of the table are values on a t-distribution with  
(df) degrees of freedom for selected right-tail (greater-than) probabilities (p).



df	0.40	0.25	0.10	0.05	0.025	0.01	0.005	0.0005
1	3.24820	3.00000	3.07789	3.13752	3.20025	3.25052	3.30574	3.36192
2	2.88675	0.91497	1.89518	2.91998	4.30285	6.94658	9.32484	31.5991
3	2.78671	0.76492	1.63774	2.55338	3.18425	4.54070	5.84091	12.9240
4	2.70722	0.74069	1.53326	2.13184	2.77495	3.74695	4.60409	8.6103
5	2.67181	0.71756	1.47594	2.01504	2.57058	3.34943	4.02214	6.8888
6	2.64835	0.71158	1.43975	1.94316	2.44691	3.14267	3.70743	5.9588
7	2.63187	0.71117	1.41424	1.89457	2.36462	2.98795	3.46948	5.4079
8	2.61921	0.70837	1.36815	1.85948	2.30600	2.88048	3.35539	5.0413
9	2.60955	0.70272	1.35313	1.83313	2.28144	2.82650	3.28704	4.7009
10	2.59985	0.70118	1.32718	1.79818	2.25210	2.75007	3.19507	4.3869
11	2.59945	0.69145	1.30349	1.78595	2.20093	2.71808	3.19501	4.4270
12	2.59903	0.69448	1.36217	1.73228	2.17811	2.80100	3.08454	4.3178
13	2.59591	0.69329	1.35071	1.70933	2.19037	2.80501	3.01228	4.2298
14	2.59213	0.69217	1.34503	1.69310	2.14473	2.82448	2.97684	4.1405
15	2.57985	0.69197	1.34066	1.75302	2.13145	2.80248	2.94671	4.0728
16	2.57598	0.69012	1.38757	1.74586	2.11991	2.58349	2.92078	4.0150
17	2.57347	0.68919	1.33379	1.73960	2.10982	2.59683	2.88823	3.9651
18	2.57123	0.68834	1.33039	1.73406	2.10052	2.55238	2.87844	3.9216
19	2.56923	0.68749	1.327728	1.729133	2.09302	2.53948	2.86083	3.8834
20	2.56743	0.68654	1.325341	1.724711	2.08596	2.52798	2.84534	3.8495
21	2.56580	0.68552	1.323188	1.720743	2.07961	2.51765	2.81316	3.8193
22	2.56432	0.68546	1.321237	1.717144	2.07381	2.50832	2.81676	3.7921
23	2.56297	0.68506	1.319468	1.713872	2.06823	2.49897	2.80734	3.7676
24	2.56169	0.68458	1.318138	1.710546	2.06326	2.49216	2.80008	3.7454
25	2.56049	0.68449	1.316245	1.708141	2.05954	2.48714	2.79744	3.7251
26	2.55955	0.68404	1.314972	1.706118	2.05553	2.47882	2.77971	3.7066
27	2.55858	0.68395	1.313703	1.703288	2.05182	2.47266	2.77068	3.6996
28	2.55768	0.68393	1.312527	1.701131	2.04841	2.46714	2.76326	3.6739
29	2.55684	0.68394	1.311434	1.699117	2.04523	2.46202	2.76583	3.6594
30	2.55605	0.682756	1.310415	1.697281	2.04203	2.45728	2.76000	3.6460
31	2.55537	0.674490	1.281552	1.644881	2.03263	2.57583	3.2905	
32	0.347	0.674490	88%	90%	95%	98%	99%	99.9%

True Definition CI  
Confidence 95%



SE for mean vs. SE for proportion

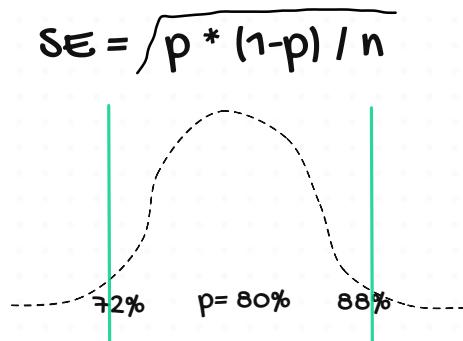
$$SE_m = SD / \sqrt{N}$$

$$SE_p = \sqrt{p * (1-p) / n}$$



$$z$$

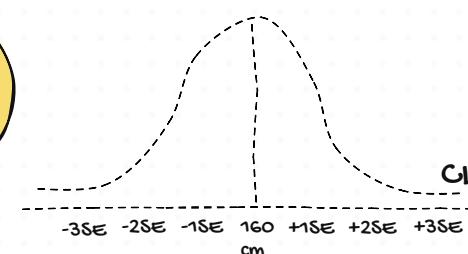
90% 1.645  
95% 1.96  
99% 2.58



n = 100  
mean = 160 cm  
sd = 5 cm



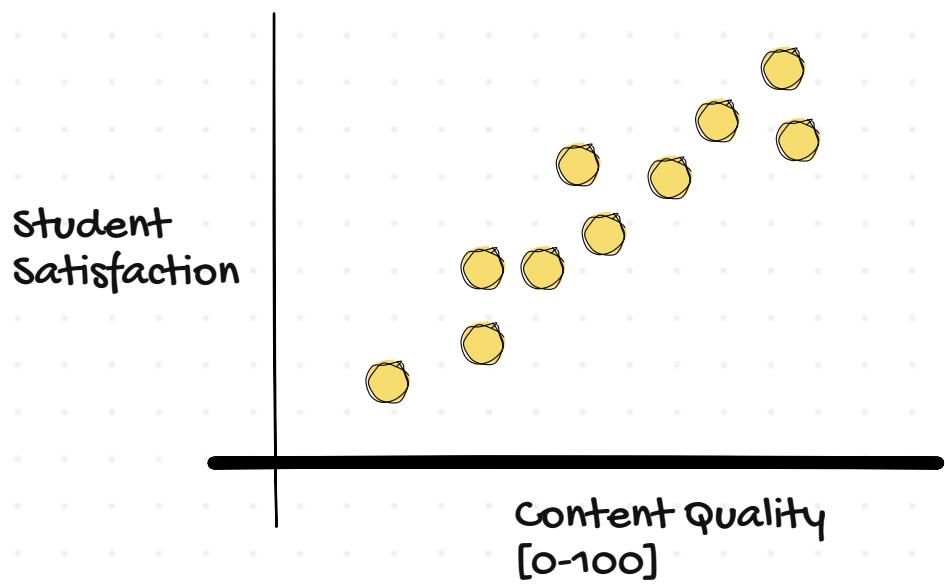
height Thai Pop



- 1  $SE = SD / \sqrt{N}$   
 $SE = 5 / \sqrt{100} = 5/10 = 0.5$
- 2 Margin Error =  $SE * 2 = 0.5 * 2 = 1$
- 3 Confidence Interval = [160-1, 160 + 1]

# Correlation

Linear association



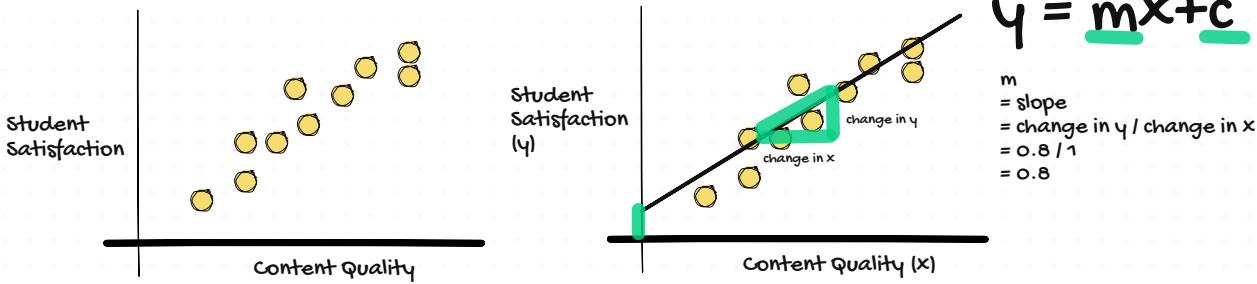
correlation ( $r$ )

- $[-1, +1]$
- + positive
- negative
- no relationship

=CORREL(x, y)

strong  $r > 0.7$   
moderate  $r \geq 0.4$   
weak  $r < 0.4$

# Linear Regression



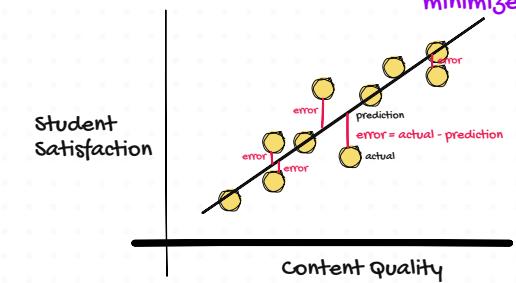
$$y = f(x)$$

$$y = b_0 + b_1 \cdot x$$

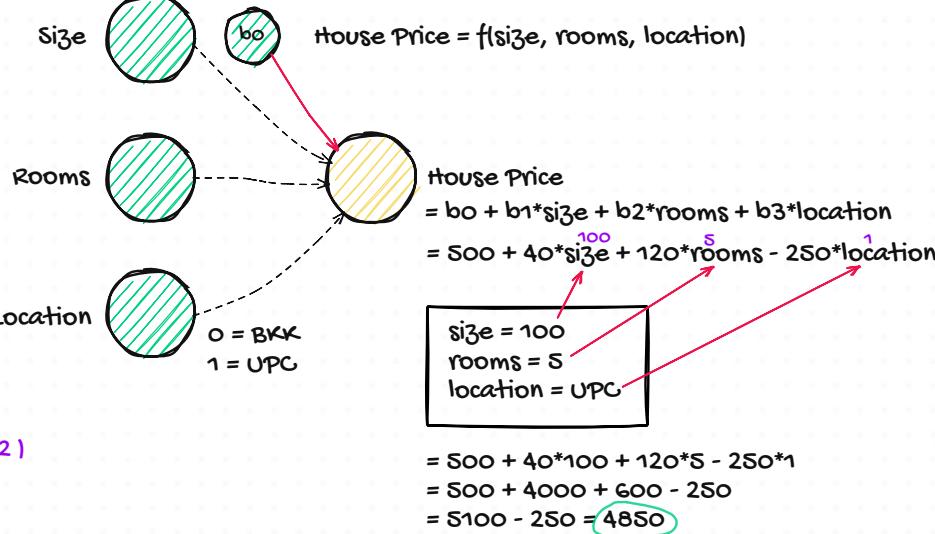
$b_0$  = y intercept  
 $b_1$  = slope  
 Coefficients best fitted line

## Least Squared Error

minimize  $\Rightarrow \text{sum}(\text{error}^2)$



Find  $b_0, b_1$   
to minimize  $\text{sum}(\text{error}^2)$

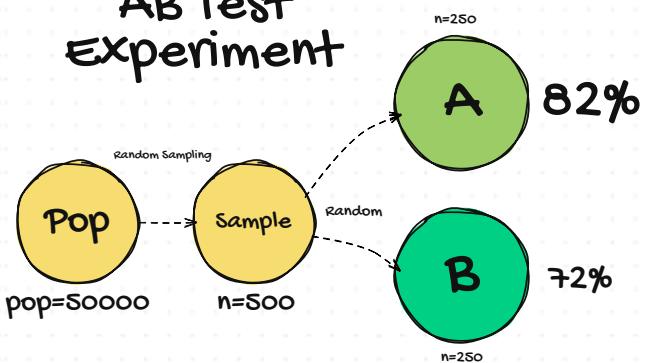


## Significance Test

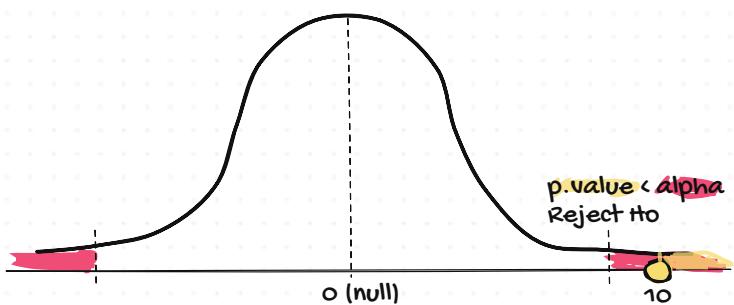
x is significance if  $p.\text{value} < 0.05$

Break 10 min (2.51pm)

## AB Test Experiment



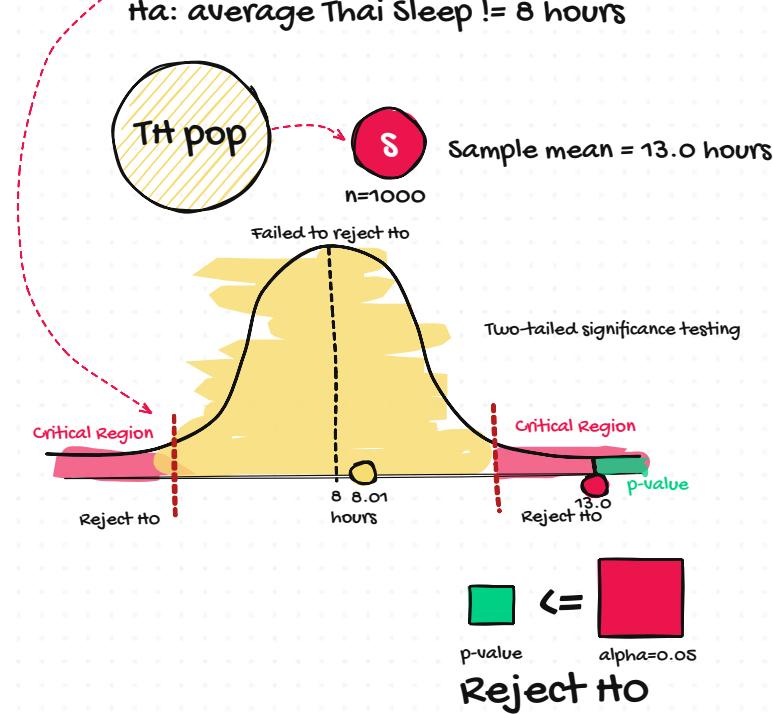
$H_0: \% \text{ like A} - \% \text{ like B} = 0$   
 $H_a: \% \text{ like A} - \% \text{ like B} \neq 0$



## Foundations of Hypothesis Testing

# Frequentist Approach

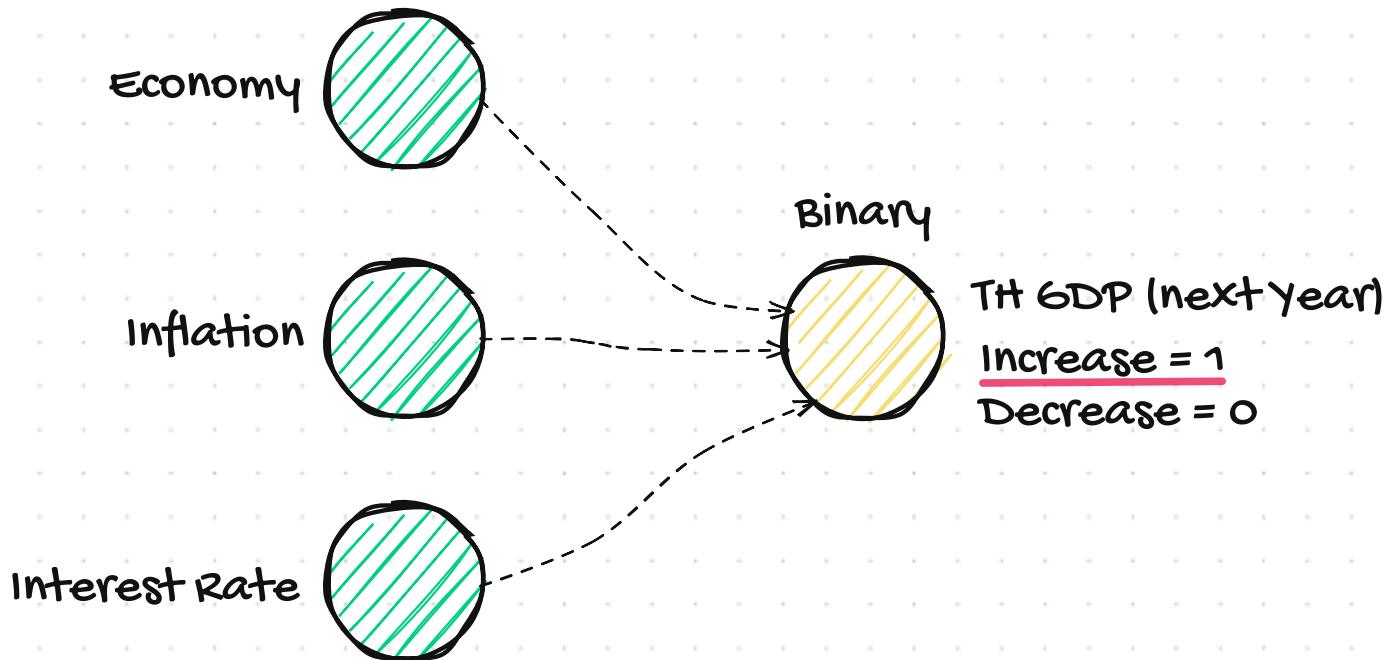
$H_0: \text{average Thai sleep} = 8 \text{ hours}$   
 $H_a: \text{average Thai Sleep} \neq 8 \text{ hours}$



$$\begin{matrix} \text{p-value} & \leq \\ & \alpha = 0.05 \end{matrix} \quad \text{Reject H}_0$$

$p\text{-value} = p(\text{observed data or more extreme} \mid H_0 \text{ is true})$

# Logistic Regression



GDP = sigmoid( [economy, inflation, interest rate] )

$p(Y=1) = \text{sigmoid}( b_0 + b_1 * \text{economy} + b_2 * \text{inflation} + b_3 * \text{ir} )$

$$\frac{1}{1+e^{-z}}$$

[0, 1]

=IF( p(Y=1) > 0.5, 1, 0 )

# Intro Neural Network

