Agenda:
Framework of problem statement in Data Science, Machine Learning, AI
Different types of problems
Different types of roles


About me:
I love teaching
Current:
Lead AI Scientist and Instructor at Scaler
Past:
Lead AI Scientist at Target
AI Scientist at AlphaICs
PhD from Indian Institute of Science,
B Tech from NITK, Surathkal

Yes or no
↓

| Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|
| 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| 5 | 116 | 74 | 0 | 0 | 25.6 | 0.201 | 30 | 0 |
| 3 | 78 | 50 | 32 | 88 | 31.0 | 0.248 | 26 | 1 |

Which attributes are most relevant
→ Statistics / Correlation / hypothesis test
Can you predict whether a new person will have diabetes?
→ Machine Learning "SKLEARN" → Python

*"Logistic regression"*
*"Decision tree"*
*"Random forest" "XGBoost"*
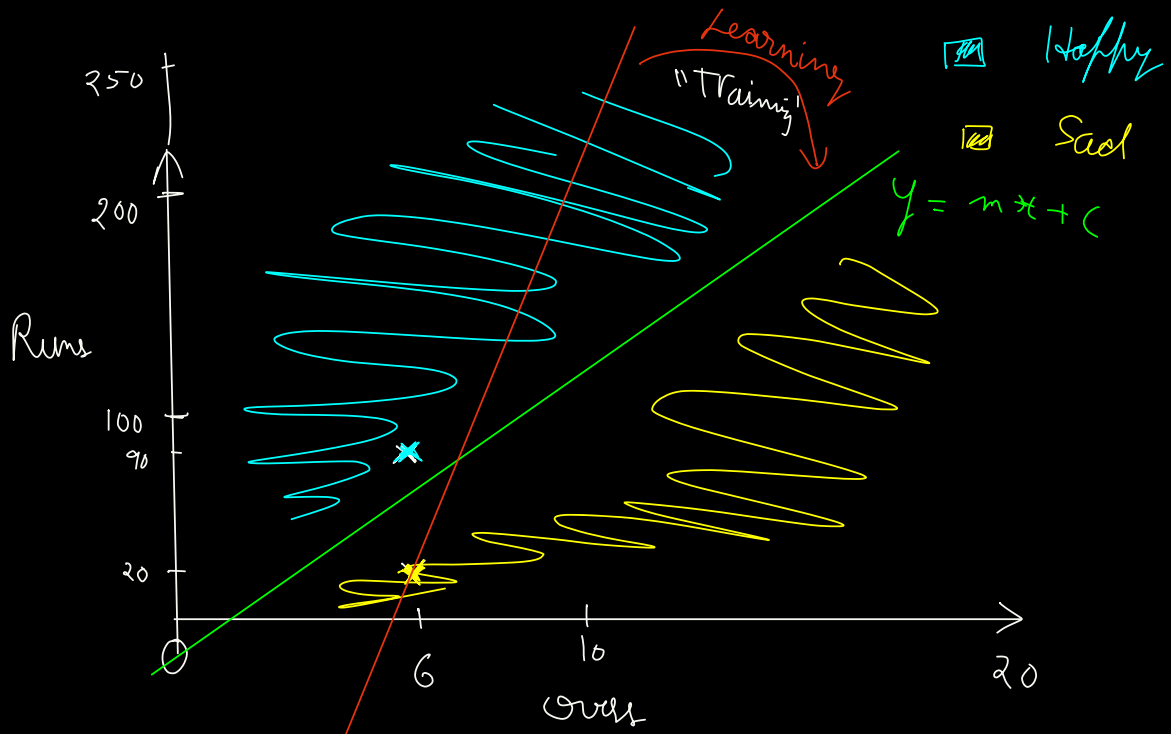*"Support vector Machine"*
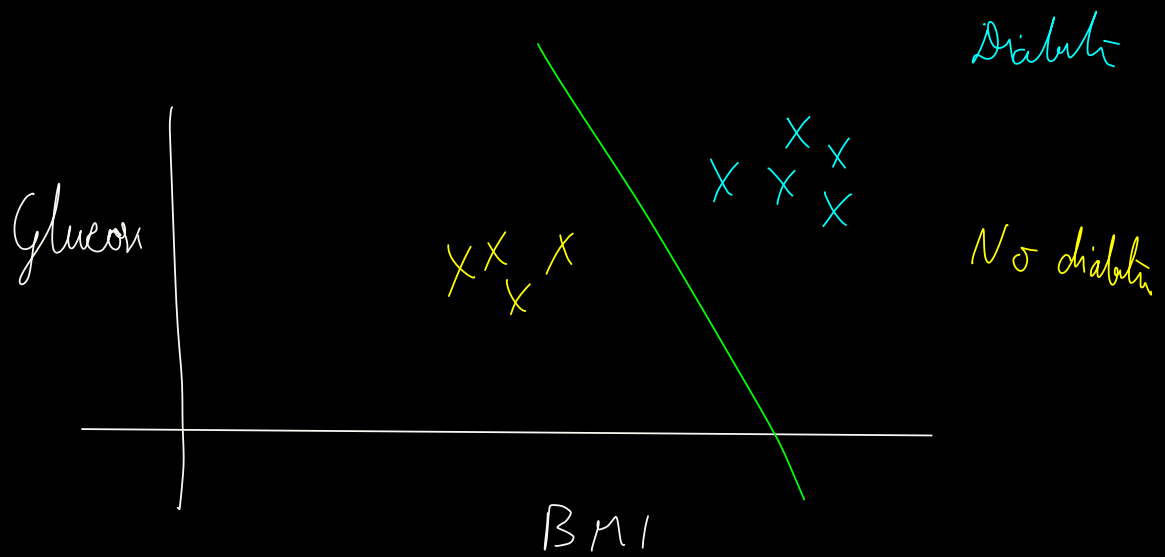
IPL    Machine Learning

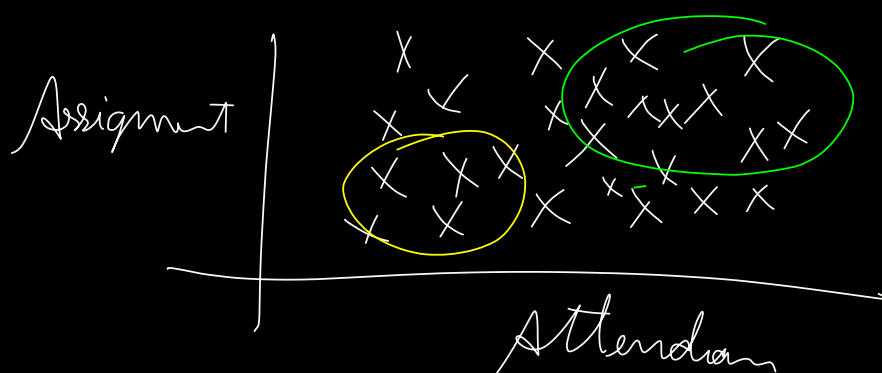6 overs ⟶ 90 runs        "Happy"

6 overs ⟶ 20 runs        "Sad"

15 overs ⟶ 90 runs       "Sad"



Learning

"Training"

□ Happy

□ Sad

$y = mx + c$

250

200

Runs

100
90

20

O        6        10        20

Overs

**Glucose**

Diabetic

No diabetic

BMI

**Lot of Students**

**Assignments**

Attendance
Live / Recorded

**Assignment**

Attendance

"Clustering"

"Unsupervised"

K - means
cluster

Recy Sys

Recommends
↙

increase sales

2012 $\longrightarrow$ `` AlexNet '' ImageNet



`` Deep Learning ''

Segment Anything

Entry level

A

No python
Not in SQL     51%
No Math

B    11%

Yes python / Java / C++

Not in SQL

No math

C    16%

Yes python
Yes SQL
No Math

D    22%

Yes python
Yes SQL
Yes Math

A:   Beginner  ⟶  Target a Data Analyst role

B:   Beginner w/o python

C:   Intermediate
D:   Advanced

$\underline{IPL}$ :   6 overss $\longrightarrow$ 90 runs

6 overs $\longrightarrow$ 20 runs

Happy

Sad



$y = mx + c$

250

Runs

180
90

0      6    10        20 overs

2012 $\longrightarrow$ AlexNet

"Deep Learning"

① "Computer vision"

VGG / ResNet

"Neural networks"

$\{25, 30, 40 \ldots \}$
$\{ \ldots - - - - \} \Longrightarrow$ "car"

② "Natural Language Proc"

NLP

NLP





your plan do I really look like a guy with a plan you

Audio $\longrightarrow$ text    NLP

Computer vision $\longrightarrow$ Deep Learning

House prices    loc.   size, 1 BHK $\left.\begin{array}{c} \\ \vdots \\ 1 \\ 2 \\ 3 \end{array}\right\}$ $\longrightarrow$ price

Machine Learning

| Size | loc | Bed | Price |
|------|-----|-----|-------|
| 1000 | Ba  | 3   | 1 Cr  |

$\longleftarrow$ tabular data

Eg: Linear regression
Decision Tree

Data science

| | Run | ball | 4 | 6 | SR | . . . | Win |
|---|---|---|---|---|---|---|---|
| ① | | | | | | | ✓ |
| ② | | | | | | | ✗ |
| ③ | | | | | | | ✗ |
| ⋮ | | | | | | | ✓ |

↓
cricinfo data

Prob of Kohli scoring century?
Prob of winning when Kohli scores century

Deep learning





A simple neural network
input layer   hidden layer   output layer

→ classify
    Draw

## Word-cloud

Null hypothesis

Test - stat

$\boxed{P\text{-value}} \longrightarrow$ if $p < \alpha \nearrow$ reject Ho
(Choose Ha)

"$\alpha$" Significance level $(0.05, 0.01)$

Confusion matrix $(FP, FN)$

$P\left[\begin{array}{c} \text{Data as extreme} \\ \text{as that observed} \end{array} \middle| \begin{array}{c} \text{Ho is} \\ \text{true} \end{array}\right]$

Burger company
Its burger weighs 200 grams
An unsatisfied customer, who is still hungry, wants to disprove this claim

Customer should prove that on average, the burger weighs < 200

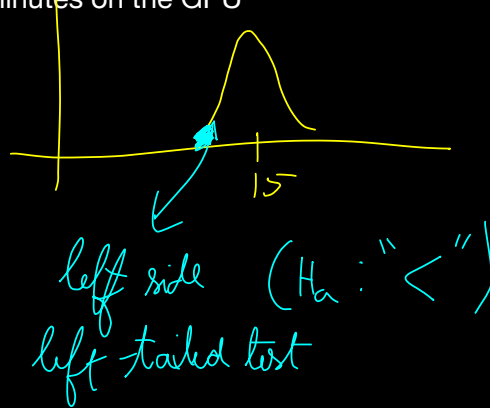H0: average = 200
Ha: average < 200

AI chip company
This company wants to claim that it is better than GPU
The training time for ResNet is 15 minutes on the GPU
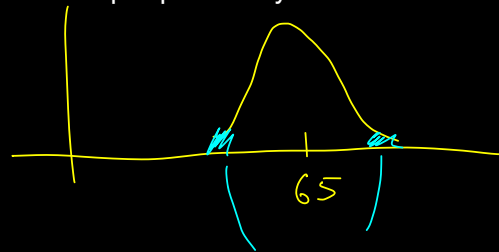
H0: training time $=$ 15
Ha: training time < 15

left side $\quad (H_a : "<")$
left tailed test

The average height of Indians is 65 inches
You want to verify whether this is true for people from your state

H0: height = 65
Ha: height != 65

65

two-tailed test

Retail example

Recap of CLT
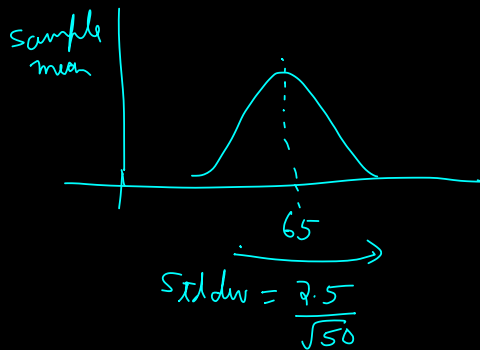Average height is 65 inches, std dev 2.5
We sample 50 people
Let "m" denote sample mean
Is "m" a random variable? yes
What is its distribution? Gaussian (normal)
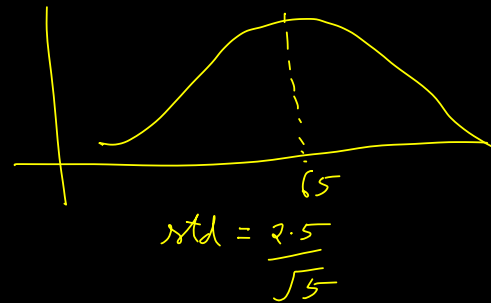What is E[m]?  65
What is the std dev of m?  $\frac{2.5}{\sqrt{50}}$

sample
mean

65

$Stddv = \frac{2.5}{\sqrt{50}}$

Sample 5 people

$m \rightarrow$ sample mean

$E[m] = 65$

$Std\ dev = \frac{2.5}{\sqrt{5}}$

65

$std = \frac{2.5}{\sqrt{5}}$

Retail outlet with 2000 stores.
Weekly sales: Shampoo bottles: mean = 1800, std dev = 100

1)
We want to hire a marketing team to improve sales

Test them in 50 stores
In these 50 stores, our average sales is 1850
Want 99% confidence --> alpha = 0.01 (significance level)

2) Another team is deployed, and their average is 1900, number of stores was 5

What is the null hypothesis

H0: average = 1800 (marketing has no effect)
Ha: average > 1800 (marketing has effect)

Data: 50 stores, average here is 1850
Test statistic "m": sample mean of 50 stores
Distribution of test statistic? Gaussian
What is E[m] = 1800 ( under H0)
What is std dev of m? 100/root(50)

$$P\left[m \geq 1850 \mid H_0 \text{ is true}\right]$$

$$1 - \text{norm.cdf}(3.53) = 0.0002 < 0.01$$

Reject $H_0$ (Marketing had effect)

$H_0$ : avg = 1800

$H_a$ : avg > 1800

5 - stores, Sample mean obs way 1900
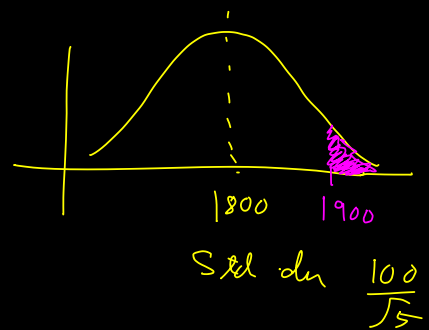
Test statistic $m =$ sample mean

$$E[m] = 1800$$

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

$$P\left[m > 1900 \mid H_0 \text{ is true}\right]$$

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



1800   1900

Std dev $\frac{100}{\sqrt{5}}$

$$1 - \text{norm.cdf}(2.23) = 0.012 > 0.01$$

Effect is not statistically significant (stick $H_0$)

$\mu = 1800$

$\sigma = 100$

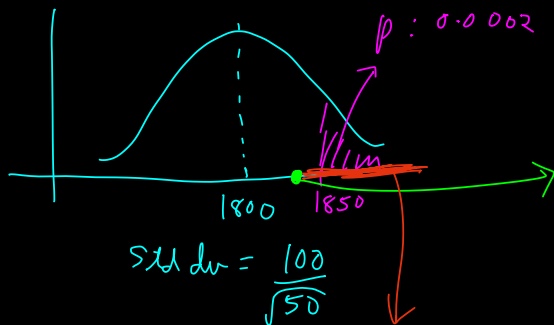$H_0 :$ avg $= 1800$

$H_a :$ avg $> 1800$

50 story

$\alpha = 0.01$

$p : 0.0002$

1800    1850

Std dv $= \dfrac{100}{\sqrt{50}}$

3-score    norm. $pdf\,(0.99) = 2.32$

$1800 + \dfrac{100}{\sqrt{50}}\,(2.32) \;=\; 1832.8$

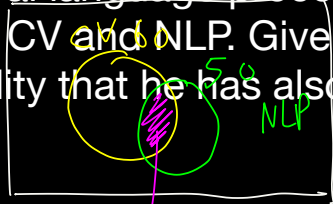"critical value"

$> 1832.8$

"critical region"

Interview  Style Class

Among 100 students, 60 have taken the computer vision (CV) module, 50 have taken natural language processing (NLP). Also, it is seen that 20 have taken both CV and NLP. Given that a person has taken NLP, what is the probability that he has also taken CV?

$$P[CV] = \frac{60}{100}$$

① $P[CV \cap NLP]$

② $P[CV \mid NLP]$

③ $P[NLP \mid CV]$

$$P[NLP] = \frac{50}{100}$$

$$\frac{20}{100}$$

$$P[CV \mid NLP] = \frac{P[CV \cap NLP]}{P[NLP]} = \frac{20/100}{50/100} = \frac{20}{50}$$