

# Interview Prep Session

[Probability & Statistics]

Recap through Questions

Q: Count the # times model will be trained, if we want to do a grid search with the following params!

→ Automated | EASY

```
params = {'depth': [3, 4, 5],  
          'n_estimators': [10, 100, 1000],  
          'min_samples': np.arange(1, 5)]
```

a) 10

b) 36

c) 11

d) 45

Ans: 3 × 3 × 4 = 36

Q: If  $n$  customers come in, I will buy on avg. Find the probability of selling at-least 1 unit, if  $n$  customers come in.

- a) 1
- b)  $1/n$
- c)  $\frac{175}{256}$
- d)  $\frac{81}{256}$

$$\text{Ans} \rightarrow P(\text{Buy}) = 1/n$$

$$P(\text{Don't Buy}) = 1 - 1/n = 3/n$$

$$n \text{ cust not buying} = \left(\frac{3}{n}\right)^n = \frac{81}{256}$$

$$\text{At least 1 purchase} \\ = 1 - \frac{81}{256} =$$

$$\boxed{\frac{175}{256}}$$

Q: I have a binary classification model that has 90% accuracy.

Find the prob that it makes exactly  $k$  errors out of 20 predictions

→ Ans: 1 classification

↳ Bernoulli event

$$p = 1 - 0.9 = 0.1$$

20 predictions

↳ Binomial event

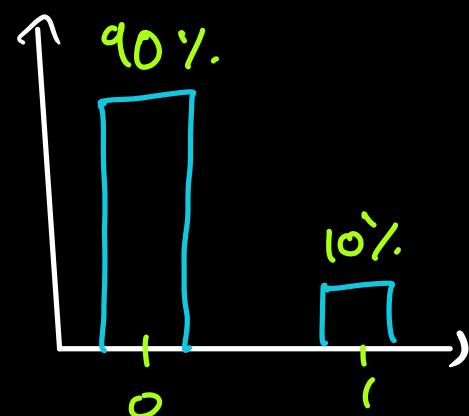
$$d = \text{Bin}(20, 0.1)$$

$$P(n) = d \cdot \text{pdf}(n)$$

$$= 20C_n \times (0.1)^n \cdot (0.9)^{16}$$

Toda Recap binomial

Q How to solve this??



Binomial dist of 1 trial.

But we have 20 trials!!

Q! What are the possible o/c in 20 trials?

- 1 pass
- 2, pass

{ → we need atleast 4

→ 20 pass → to pass

Q: What is the prob of exactly 1 error?

Soln → Total possible combinations??

$$\frac{P/F}{1} \frac{P/F}{2} \frac{P/F}{3} = \dots \frac{P/F}{20} = 2^{20}$$

Number of ways in which 1 error

$$\frac{P}{F} \frac{F}{F} \frac{F}{F} \dots \rightarrow 20C_1 = 20 \text{ ways}$$
$$(0.1) \frac{(0.7)}{F} \frac{(0.9)}{F} \dots \binom{ }{F}$$

$$P(1 \text{ pass}) = \frac{20}{2^{20}} \rightarrow \times \text{ Events are not}$$

equally likely  
Imagine getting HHT with biased coin

$$= p \times (1-p)^{19} = (0.1)(0.9)^{19} \times 20$$

Q: Prob of getting 2 passes?

$$\# \text{ of ways} = 20C_2 = \frac{20 \times 19}{2} = 190$$

$$\text{Prob} \rightarrow (0.1)^2 (0.9)^{18} \times 190$$

Q: Notice a pattern yet?

Probability of getting  $n$  success in 'n' trials

$$P(n_{\text{succ}}; n_{\text{trial}}) = {}^nC_n p^n (1-p)^{n-n}$$

Q: There are 3 servers behind a local balancer. Traffic is diverted to each server with given probability. The down time prob for each server is given. If one user received server not found error, what is the prob it was server 1.

$P_{\text{server}}$

$$S_1 = 20\%$$

$$S_2 = 40\%$$

$$S_3 = 40\%$$

$P_{\text{down}}$

$$S_1 = 3\%$$

$$S_2 = 2\%$$

$$S_3 = 2\%$$

- a) 27%    b) 35%    c) 50%    d) 10%

$$P(A) = 0.2$$

$$P(E_1/A) = 0.03$$

$$P(B) = 0.4$$

$$P(E_2/B) = 0.02$$

$$P(C) = 0.4$$

$$P(E_3/C) = 0.02$$

Bayes theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$\therefore P(A/E_1) = \frac{P(E_1/A) \cdot P(A)}{P(E_1/A) \cdot P(A) + P(E_2/B) \cdot P(B)}$$

. . .

. . .

$$= \frac{(0.03)0.4}{(0.03)(0.2) + (0.02)(0.9) + (0.02)(0.6)}$$

$$= \frac{6}{6 + 8 + 8} = \frac{6}{22} = 27\%$$

# All Equations

King . Truth + or King . Lie

$$P = \frac{\# \text{ favourable outcomes}}{\# \text{ total possible outcomes}}$$

$$P(A') = 1 - P(A) \quad \# \text{ complementary}$$

$$P(A \cap B) = P(A) \cdot P(B) \quad \# \text{ mutually independent}$$

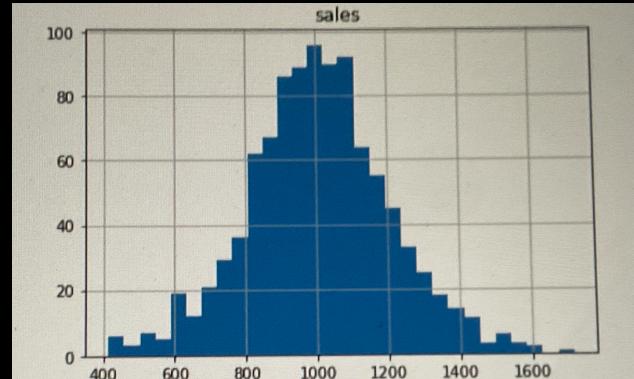
$$P(A \cup B) = 1 \quad \# \text{ exhaust} \quad P(A \cap B) = 0 \quad \# \text{ exclusive}$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A|B) = P(A \cap B) / P(B)$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Q. Toothpaste sales  
Retail →



**Quiz:** Which among these mostly likely represent the weekly sales data of toothpaste?

- 1. Bernoulli
- 2. Binomial
- 3. Gaussian
- 4. Geometric

\*\*\*

Suppose the store has **beginning on-hand (BOH)** inventory of 1300. If the demand is more than the BOH, there is a need for stock **replenishment**. Which distribution best characterizes the need for stock replenishment?

**Quiz:** Which among these mostly likely represent need for replenishment?

- 1. Bernoulli
- 2. Binomial
- 3. Gaussian
- 4. Geometric

\*\*\*

Suppose there are 2000 stores, each with BOH of 1300. The **distribution center (DC)** needs to calculate the number of stores which might need replenishment. Which distribution best characterizes this number?

**Quiz:** Suppose there are 2000 stores, each with BOH of 1300. The distribution center (DC) needs to calculate the number of stores which might need replenishment. Which distribution best characterizes this number?

- 1. Bernoulli
- 2. Binomial
- 3. Gaussian
- 4. Geometric

\*\*\*

The manager of the DC sequentially calls a few store managers to see if they need replenishment. We need to monitor how many calls are needed till the first time a store manager asks for replenishment. Which distribution best characterizes this?

**Quiz:** The manager of the DC sequentially calls a few store managers to see if they need replenishment. We need to monitor how many calls are needed till the first time a store manager asks for replenishment. Which distribution best characterizes this?

- 1. Bernoulli
- 2. Binomial
- 3. Gaussian
- 4. Geometric

Q: What do you understand by p-value?

- a)  $P(\text{rejecting } H_0)$
- b)  $P(\text{rejecting } H_A)$
- c)  $P(H_0 \mid \text{data})$
- d)  $P(\text{data} \mid H_0)$

## Examples

→ Crime suspect

Null  $\rightarrow H_0$ : Innocent (until proven guilty)

Alternative  $H_A$ : Guilty

D: (Evidence / witness)

Data:

- 1) he has a knife
- 2) knife has blood staining
- 3) blood matches victim's
- 4) No credible alibi.

$$P(\text{Data} | H_0) = ?? \text{ Low}$$

$\underbrace{\quad}_{\text{p-value}}$   $\hookrightarrow \text{Reject} \rightarrow \underline{\text{Jail}}$

→ Suppose you are the 3<sup>rd</sup> empire. On field  
empire has signalled "out"

$H_0 \rightarrow$  On field empire is correct

Data: Slow motion video + audio

If we have sufficient data, we can change  
on-field empire decision

$$P(\text{Data} \mid \text{Out}) = \begin{cases} \text{High} \rightarrow \text{Fail to reject (out)} \\ \text{Low} \rightarrow \text{Reject} \\ \quad (\text{Not out}) \end{cases}$$

Q: Which of the following is the best way to detect a normal distribution?

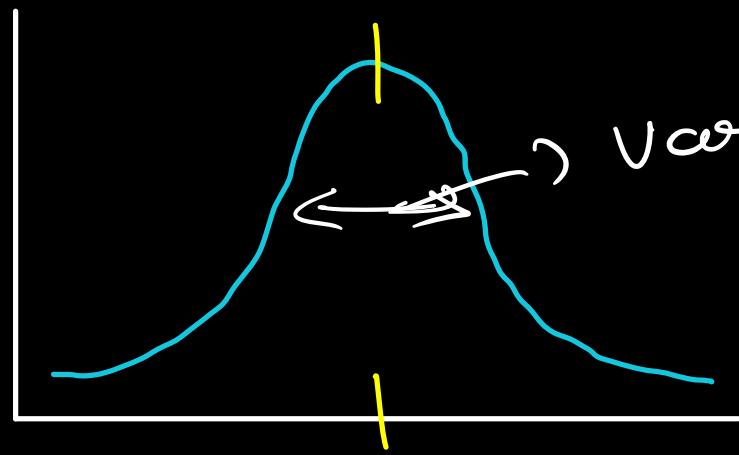
- a) 68 - 95 - 99.7 rule
- b) Q Q plot
- c) KS test
- d) area under curve = 1

Mathematically.

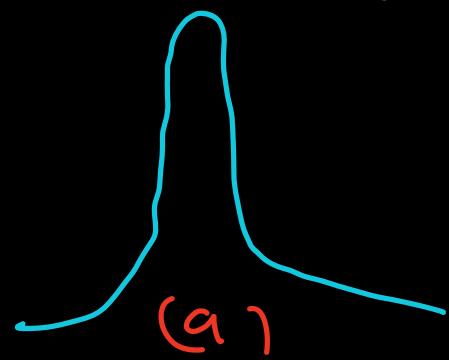
$$X \sim N(\mu, \sigma)$$

normal      mean      std

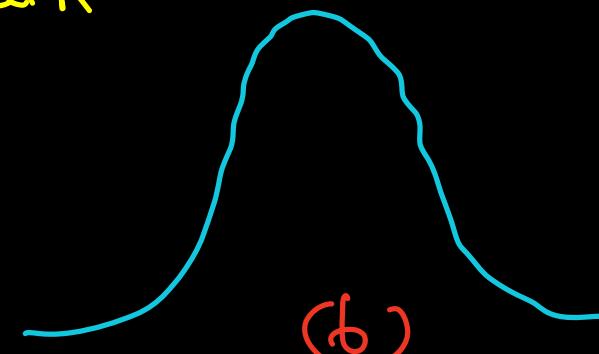
$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Variance / std



$\mu = \text{Mean}$



(b)

Not Imp

$\text{std}(a) < \text{std}(b) < \text{std}(c)$



(c)

## Properties of Normal Dist

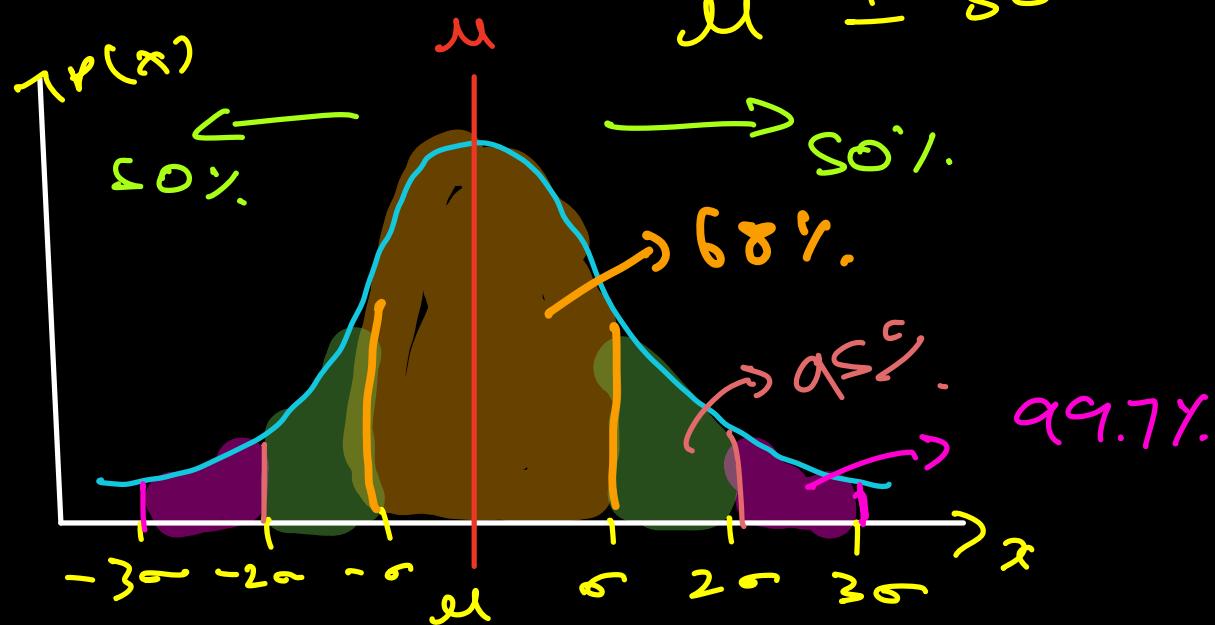
→ Symmetric

→ mean =  $\mu$  , std =  $\sigma$  , var =  $\sigma^2$

→ area with  $M \pm \sigma = 68\%$ .

$$M \pm 2\sigma = 95\%$$

$$\mu \pm 3\sigma = 99\%.$$



Q: How do calc the CI of the mean of a population if one sample is given to you?

Calc for sample  $\rightarrow \bar{x} = 12$

$$s = 3$$

$\therefore$  95%). CI  $n = 100$

$$= 12 \pm \frac{2 \times 3}{\sqrt{100}} = \underbrace{11.4 \rightarrow 12.6}$$

Q: Match test to application  $\rightarrow$

- |             |  |
|-------------|--|
| a) T        | Gender affects churn                   |
| b) $\chi^2$ | all subjects are equally difficult     |
| c) Z        | Indian customers shopping profile = US |
| d) KS       | $U_1 > U_2$ with 10 samples            |
| e) F        | $S \propto u_1 > u_2$                  |

- $\rightarrow$  Ans:
- a)  $\rightarrow 4$
  - b)  $\rightarrow 1$
  - c)  $\rightarrow 5$
  - d)  $\rightarrow 3$
  - e)  $\rightarrow 2$

Q: Diff b/w expectation and average  
↳ calc the expectation of dice roll

Q: Historic avg monthly car sales are known to be 15000 units. Using last years sales data, can you design a Z-test to validate if the average has changed?

↳ No, because this sample is

not independent.

Q: Why do investors want to diversify their portfolio?

↳ the variance of the group is lower than the variance of individuals

↳ similar to what happens in bagging methods