

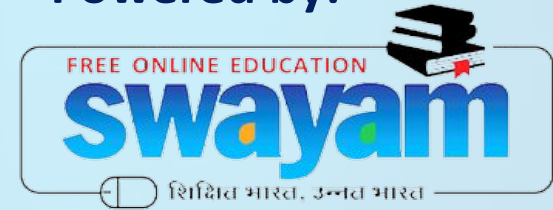
NPTEL Week-3 Live Session

on Machine Learning and Deep Learning - Fundamentals and Applications (noc24_ee146)

A course offered by: Prof. Manas Kamal Bhuyan, IIT Guwahati

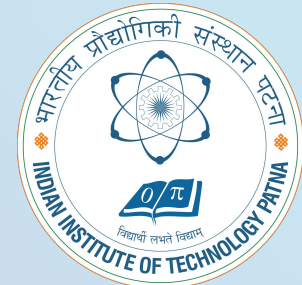
NPTEL Quiz Solution: ^{ve} week-1, ^W week-2

Powered by:



PMRF

Prime Minister's Research Fellows
Ministry of Education
Government of India



By

Arka Roy

NPTEL PMRF TA

Prime Minister's Research Fellow

Department of Electrical Engineering, IIT Patna

Web: <https://sites.google.com/view/arka-roy/home>



In a binary classification problem, the confusion matrix is a

① _____ matrix.

1x1

~~2x2~~

3x3

1x2

Actual Pos
= Neg

	Pos	Neg	
Pos	TP	FN	(2x2)
Neg	FP	TN	
	Predicted		

② Precision is defined as

$$TP / (TP + TN)$$

$$TP / (TP + FN)$$

$$\checkmark TP / (TP + FP)$$

$$TN / (TN + FP)$$

$$Precision = \frac{TP}{TP + FP}$$

③ In a binary classification problem, a classifier correctly predicts 90 instances as positive, incorrectly predicts 15 instances as positive when they are negative, correctly predicts 90 instances as negative, and incorrectly predicts 10 instances as negative when they are positive. What is the accuracy of the classifier?

80

85

~~87.8~~

95

$$\rightarrow TP = 90; FP = 15; TN = 90; FN = 10$$

$$\checkmark Acc = \frac{TP + TN}{TP + FN + FP + TN} = \frac{90 + 90}{90 + 10 + 15 + 90} = 87.8\%$$

For the above question find the F1 score?

78.2%

85%

~~87.8%~~

$$F1 \text{ score} = \frac{2 \times Pre \times Rec}{Pre + Rec}$$

$$= \left(\frac{2 \times 0.857 \times 0.9}{0.857 + 0.9} \right) = 0.878$$

Sensitivity

$$Rec = \frac{TP}{TP + FN} = \frac{90}{90 + 10} = 0.9$$

$$Prec = \frac{TP}{TP + FP} = \frac{90}{90 + 15} = \frac{90}{105} = 0.857$$

Q Consider a dataset with actual values (Y) and predicted values (Y_pred) given below:

✓ Y = [5, 8, 12, 10, 15],

✓ Y_pred = [4, 7, 10, 11, 13].

What is the bias of the model?

Linear regression:- weights

$$\hat{y} = ax + b$$

$\hat{y} \rightarrow$ Prediction values.

$\{x_i, y_i\}$
 \uparrow Input
 \rightarrow Output / Actual values.

✓ Bias = $E(\hat{y}_{pred} - y_{true})$

= Avg ($\hat{y}_{pred} - y_{true}$)

= Avg (-1, -1, -2, 1, -2) = (-1) Bias estimate

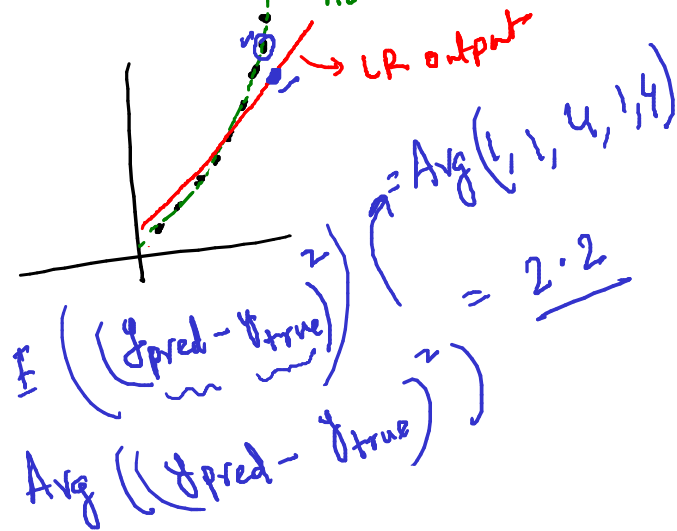
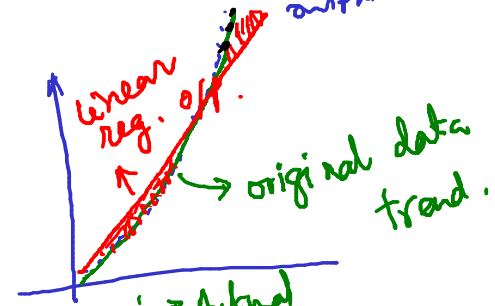
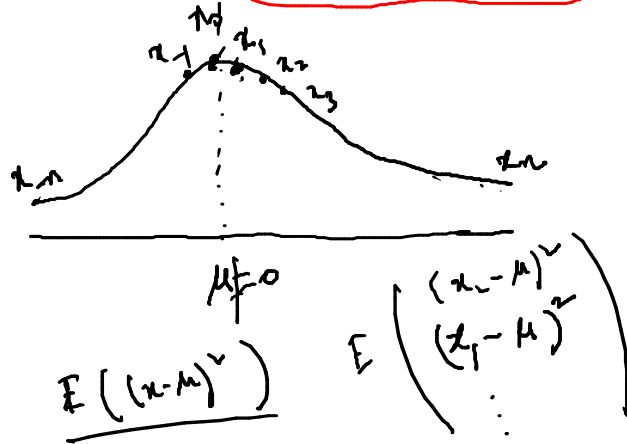
What is the variance of the model for the data given in the above question?

0

-1

✓ 2.2

None of the above



Consider a dataset with actual values (Y) and predicted values (Y_pred) given below:

Y = [5, 8, 12, 10, 15], \rightarrow Actual
 Y_pred = [4, 7, 10, 11, 13]. \rightarrow Predict

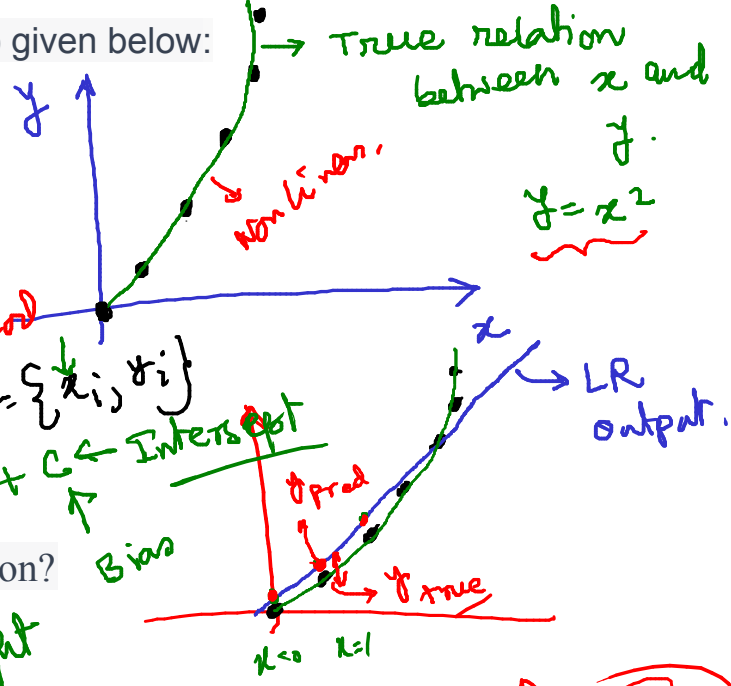
What is the bias of the model?

Model performance
 Bias estimate
 = Avg (Error type)
 = $E(L_1 \text{ norm})$

True estimates = $\{(1,1), (2,4), (3,9), (4,16), (5,25)\}$
 $y = x^2$

Linear regression method

Fit a line using $D = \{x_i, y_i\}$
 $\hat{y} = mx + c$
 \uparrow weight \uparrow Intercept



What is the variance of the model for the data given in the above question?

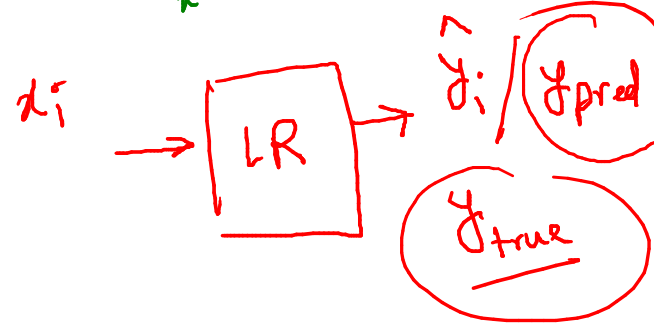
0

-1

2.2

None of the above

Variance estimate
 = Avg $\left(\frac{(y_{pred} - y_{true})^2}{2} \right)$
 = $E(L_2 \text{ norm})$



Given $X = \{-2, -1, 0, 1, 2, 3, 4, 5, 6, 7\}$ and the corresponding $Y = \{-0.5267, 1.3517, 3.8308, 5.5853, 7.5497, 9.9172, 11.2858, 13.7572, 15.7537, 17.3804\}$. Find the parameters of the linear regression model.

2.0065, 4.0312

2.0065, 3.5722

1.9214, 3.5722

None of the above

$y + c$

$D = \{x_i, y_i\} \rightarrow \boxed{LR} \rightarrow \hat{y} = wx + b \leftarrow \begin{matrix} w \\ m \end{matrix} x + c$

Parameters = $\{w, b\}$

$y = wx + b \quad \text{--- (A)}$

$E(y) = E(wx + b)$

$\Rightarrow E(y) = E(wx) + E(b)$

$\Rightarrow E(y) = w E(x) + b$

$\Rightarrow \left(\frac{1}{n} \sum y_i\right) = w \left(\frac{1}{n} \sum x_i\right) + b$

$\Rightarrow 8.58851 = 2.5w + b \quad \text{--- (1)}$

Solve (1) & (2) $\Rightarrow w = 2.0065$
 $b = 3.5722$

(A) $\times x$ on both side

$xy = wx^2 + b$

$E(xy) = w E(x^2) + E(b)$

$E(xy) = w E(x^2) + b$

$\Rightarrow \frac{1}{n} \sum x_i y_i = w \frac{1}{n} \sum x_i^2 + b \quad \text{--- (2)}$

$\Rightarrow \sum x_i y_i = \frac{w}{n} \sum x_i^2 + nb \quad \text{--- (2a)}$

$\Rightarrow 380.2522 = w \cdot 145 + 25b \quad \text{--- (2b)}$

LR:-

$\hat{y} = 2.0065x + 3.5722$

$\sum y_i (x^2) =$
 $\sum y_i (x-1) =$

weights
 w
 b (intercept of line)

Find the MSE for the above question.

0.05783

0.04247

0.04876

None of the above

$$MSE = \frac{1}{n} \sum (y_{pred} - y_{true})^2$$

$$E \left(\left[y_{pred} \right] - \left[y_{true} \right] \right)^2 =$$

y_{pred}
 $\hat{y} = 2.0065x + 3.5722$
 $y_{pred} = \{ -0.4408, \dots \}$
 $x = \{ -2, -1, 0, 1, 2, 3, 4, 5, 6, 7 \}$
 $y = \{ -0.5267, \dots \}$

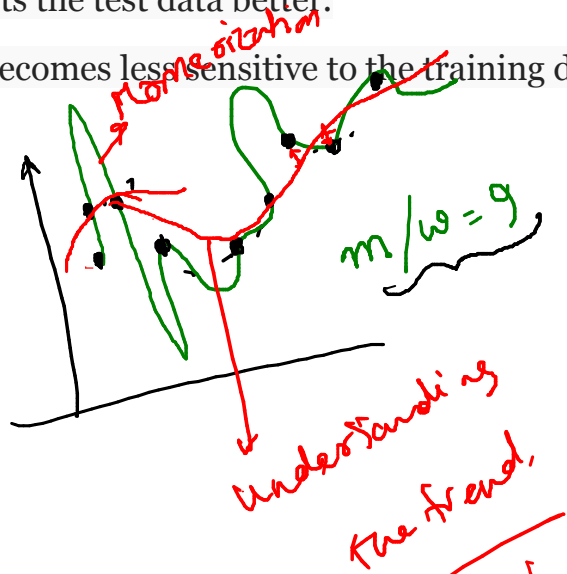
A model with high variance and low bias means

It can be too simple to understand the patterns of the data used in the training.

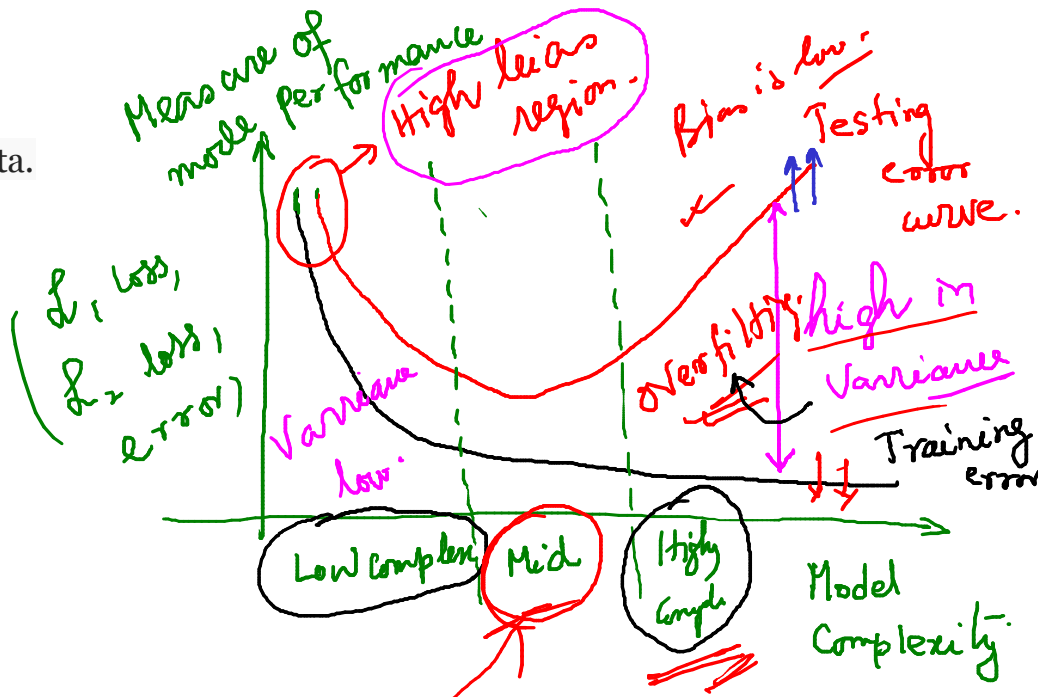
An excellent performance in the training data, but has a significant decrease in performance when evaluating the test data.

The model fits the test data better.

The model becomes less sensitive to the training data.



$m/w = 9$



Which of the following techniques is used to prevent overfitting in machine learning?

☐ To create complex machine learning models.

☐ Train the model for more epochs. (you increase the training data).

☒ Using a regularization to the model.

☐ To increase the variance of the model.

L_2 reg, L_1 reg Weight norm reg.

Consider a binary classification problem with two classes, A and B with prior probability $P(A)=0.6$, $P(B)=0.4$. Let X be a single binary feature that can take values 0 or 1. Given: $P(X=1|A)=0.8$ and $P(X=0|B)=0.7$. Determine which class the classifier will classify when $X=1$

Class A

Class B

Equiprobable for Class A and Class B

Not enough information

Feature is distributed Bernoulli

$$P(A) = 0.6$$

$$P(B) = 0.4$$

$$X \in [0, 1]$$

Bernouli dist

$$p(x=x) = p^x (1-p)^{1-x}$$

class conditional

$$P(X=1|A) = 0.8$$

$$P(X=0|B) = 0.7$$

$$P(X=0|A) = 0.2$$

$$P(X=1|B) = 0.3$$

$$P(A/X=1) = \frac{P(A) \cdot P(X=1|A)}{P(A) \cdot P(X=1|A) + P(B) \cdot P(X=1|B)} = \frac{0.8 \times 0.6}{0.8 \times 0.6 + 0.4 \times 0.3}$$

$$\Rightarrow \frac{0.48}{0.6}$$

$$= 0.8$$

$$P(B/X=1) = \frac{P(B) \cdot P(X=1|B)}{0.6} = \frac{0.4 \times 0.3}{0.6} = 0.2$$

$$P(A/X=1) > P(B/X=1)$$

$$X \in A$$

Bayes' decision theory assumes that:

The feature vectors are dependent on each other. ~~X~~

The feature vectors are normally distributed. ~~X~~

The feature vectors are identically distributed.

The feature vectors are uniformly distributed.

many ✓

$x = \{x_1, x_2, x_3, \dots, x_{d+1}\}$
Multivariate iid
 $\propto \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$

$\frac{1}{\sqrt{2\pi}|\Sigma|^{1/2}}$
Multivariate gaussian distribution

Bayes' Theory
Discrimination function

$g(x)$
 $g(x) = P(w/x)$

$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \end{bmatrix}$
independent

$\Sigma = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$
Dependent

$x \sim$ uniform
 $x \sim$ Bern.
 $x \sim$ Gaussians.
 $x \sim$ Any arbitrary distribution

Assume that the word 'offer' occurs in 80% of the spam messages in my account. Also, let's assume 'offer' occurs in 10% of my desired e-mails. If 30% of the received e-mails are considered as a spam, and I will receive a new message which contains 'offer', what is the probability that it is spam?

0.778

☒ 0.774

0.668

0.664

$$P(\text{offer}/\text{spam}) = 0.8$$

$$P(\text{spam}/\text{offer}) = ?$$

$$P(\text{offer}/\text{nonspam}) = 0.1$$

$$P(\text{spam}) = 0.3$$

$$P(\text{non spam}) = 0.7$$

$$\begin{aligned}
 \overset{\text{Posterior}}{\uparrow} P(\text{spam}/\text{offer}) &= \frac{\overset{\text{Prior}}{\uparrow} P(\text{spam}) \cdot P(\text{offer}/\text{spam})}{P(\text{spam}) \cdot P(\text{offer}/\text{spam}) + \overset{\text{class Cond.}}{\uparrow} P(\text{nonspam}) P(\text{offer}/\text{non spam})} \\
 &= \frac{0.8 \times 0.3}{(0.8 \times 0.3) + (0.7 \times 0.1)} = \underline{0.774}
 \end{aligned}$$

The optimal decision in Bayes Decision Theory is the one that

Minimizes the error rate.

Maximizes the error rate.

Minimizes the loss function.

Maximizes the loss function.

$$P(w/x) = \frac{P(w)P(x/w)}{P(x)}$$

$$P(w/x) \propto P(w)$$

Optimal Bayesian Theory.

Generalized classifiers.

Penalty due to uncertainty

$$BRC \approx BME$$

$$\lambda_{ij} = \begin{cases} 0 & ; i=j \\ 1 & ; i \neq j \end{cases}$$

$$x \in w_1 :-$$

$$R(x_1/x) = \lambda_{11}P(w_1/x) + \lambda_{12}P(w_2/x)$$

$$R(x_1/x) = P(w_2/x)$$

$$x \in w_1 \rightarrow BRC :- R(x_1/x) \downarrow$$

$$x \in w_1 \rightarrow BME :- P(w_2/x) \downarrow$$

Bayesian

Minimum error classifiers.

Bayesian

Minimum Risk classifiers.

Minimizes loss/Risk.

The risk function in Bayesian decision theory combines:

The prior probabilities and the likelihood function. ✗

The decision boundaries and the feature vectors. ✗

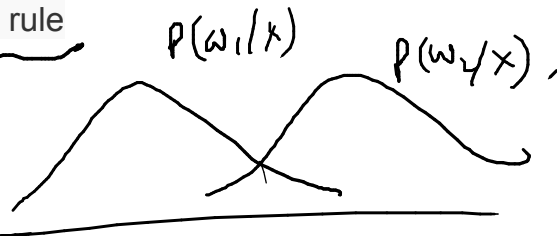
The training set and the test set. ✗

The loss function and the decision rule

$$R(a_1/x) = \lambda_{11} P(w_1/x) + \lambda_{12} P(w_2/x)$$

Annotations:

- λ_{11} : Uncertainty / Penalty / loss.
- $P(w_1/x)$: Conditional / Prior.
- $P(w_2/x)$: Posterior prob. likelihood.
- Decision rule



The loss function used in risk-based Bayesian decision theory:

Quantifies the cost of different types of errors.

Is equal to the likelihood function. $P(w/x)$ ✗

Ignores the prior probabilities of the classes. ✗

Is not used in the decision-making process. ✗

$R \propto \lambda$

The risk-based Bayesian decision rule accounts for the consequences of different decisions by considering the:

- Number of features in the dataset
- The complexity of the classifier
- ✓ Uncertainty in the data and the associated losses
- Mean and standard deviation of the feature vectors

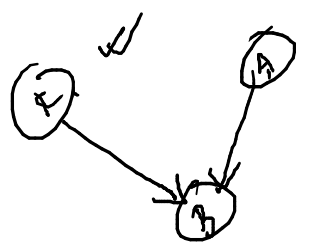
uncertainty
loss function

$$R(a/x) = \underbrace{d_{11}}_{\text{risk}} P(w_1/x) + d_{12} P(w_2/x)$$

Post.
prior *class cond.*

The generalized form of a Bayesian network that represents and solves decision problems under uncertain knowledge is known as an?

- Directed Acyclic Graph
- Table of conditional probabilities
- ✓ Influence diagram
- None of the above



Consider the following Bayesian network, where F = having the flu and C = coughing:

$$P(F) = 0.1$$



$$P(C | F) = 0.8$$

$$P(C | \neg F) = 0.3$$

Determine the probability $P(F|C)$ for the following Bayesian network so that it specifies the same joint probabilities as the given network.



$$0.23$$

$$0.03$$

$$0.35$$

None of the above.

$$P(F) = 0.1 ; \quad P(\bar{F}) = 1 - P(F) = 1 - 0.1 = 0.9$$

$$P(C/F) = 0.8$$

$$P(C/\bar{F}) = 0.3$$

$$P(F/C) = ?$$

$$P(F/C) = \frac{P(F) \cdot P(C/F)}{P(F) \cdot P(C/F) + P(\bar{F}) \cdot P(C/\bar{F})}$$

Flu part

$$= \frac{0.1 \times 0.8}{(0.1 \times 0.8) + (0.9 \times 0.3)}$$

$$= \frac{0.2286}{\approx 0.23}$$

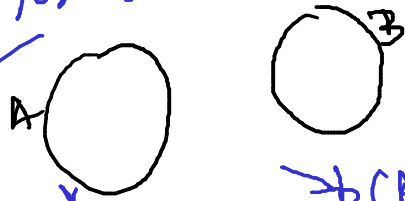
For the above question, Are C and F independent in the given Bayesian network?

- Yes.
- ☒ No.
- Can't say.
- Insufficient information.

$\cancel{P(F/C) = P(F)}$ if C, F → independent

$\cancel{P(F/C) \neq P(F)}$
 $P(F) = 0.1$

$P(F/C) = 0.23$

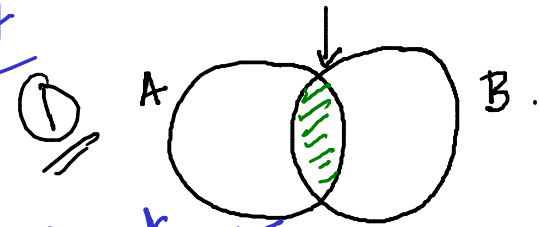


A & B are independent on each.

$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$

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

$$P(A/B) = P(A)$$

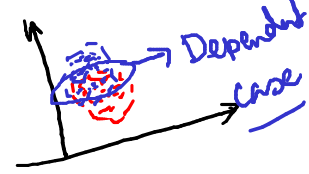


A & B are dependent events

$P(A/B)$ = Proportion of event A & B when they occur simultaneously

Dependent

$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$



Probability of B.

Bayes Rule