

Perceptron, ANN, BackProp

Deep Neural Networks
Session 02
Pramod Sharma
pramod.sharma@prasami.com

2 Agenda

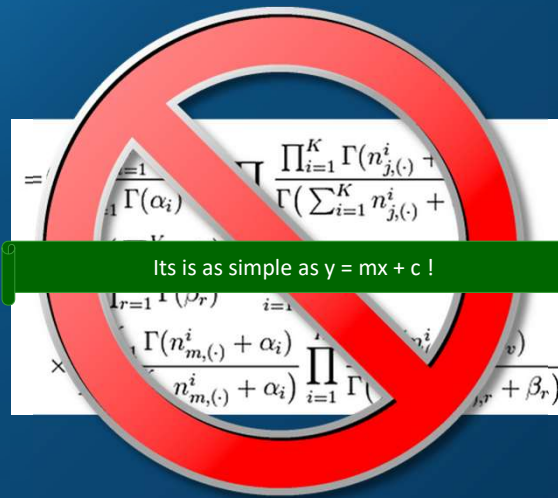
- Perceptron
- Single Layer Neural Network
- Overview of back propagation of errors

11/18/2024

pra-sami

3

Solution to Equation of Perceptron



Frank Rosenblatt

11/18/2024

pra-sâmi

4

To play or not to play...

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0	38	1	15	0	600	1
2	0	25	1	15	1	800	1
3	0	26	1	15	1	1000	1
4	5	27	1	10	1	600	0
5	20	23	0	8	1	1800	0
6	30	22	0	6	0	600	0

□ Features:

- ❖ Rains in millimeter
- ❖ Temperature in ° C
- ❖ Homework completed? – 0 : No; 1: Yes
- ❖ Team members : How many team members are ready to play?
- ❖ Is cricket equipment available?
- ❖ Ground: per hour rent in Rupees/hour

11/18/2024

pra-sâmi

5

Weights

- ❑ Each of the feature has different importance
- ❑ To assign importance to each of the feature, we use weights!
- ❑ Values of each features are in different order of magnitude
 - ❖ Summation is not going to work
 - ❖ Scale the features between 0 and 1

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0	38	1	15	0	600	1
2	0	25	1	15	1	800	1
3	0	26	1	15	1	1000	1
4	5	27	1	10	1	600	0
5	20	23	0	8	1	1800	0
6	30	22	0	6	0	600	0

- ❑ Note:
 - ❖ Variation in features have different bearing on the results
 - ❖ Team members → higher the better
 - ❖ Ground cost → lower the better

11/18/2024

pra-sâmi

6

Perceptron

- ❑ In MP Neuron Model,
 - ❖ All inputs had same weights
 - ❖ Threshold ' w_0 ' could take limited values
 - ❖ Every feature needed to be [0,1]
- ❑ Perceptron model introduced different weights to different inputs features
- ❑ Real values are also accepted
 - ❖ Temperatures are in tens and ground rent is in hundreds.
 - ❖ Min – Max – Scaler to compensate for huge difference in values
- ❑ Threshold ' w_0 ' can take any value
- ❑ Outputs are still [0, 1]

11/18/2024

pra-sâmi

7

Perceptron

□ Loss Function:

- ❖ A correction is applied on the outputs
- ❖ To adjust values of ' w_i ' to reach right results
- ❖ It would also give us indications of what weights to be fixed to arrive at the solution

□ Activation function $g(x)$ is applied as follows:

- ❖ If $\sum x_i \cdot w_i \geq w_0 \Rightarrow \hat{y} = 1$
- ❖ If $\sum x_i \cdot w_i < w_0 \Rightarrow \hat{y} = 0$

11/18/2024

pra-sâmi

8

Perceptron – Data Preprocessing

- Lets consider "Ground" and "Team Members" as features and its associated weights to arrive at the solution.

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0	38	1	15	0	600	1
2	0	25	1	15	1	800	1
3	0	26	1	15	1	1000	1
4	5	27	1	10	1	600	0
5	20	23	0	8	1	1800	0
6	30	22	0	6	0	600	0

11/18/2024

pra-sâmi

9

Perceptron – Data Preprocessing

- Scaled Data (all columns to be between 0 and 1)

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0.00	0.00	1.00	1.00	0.00	1.00	1
2	0.00	0.81	1.00	1.00	1.00	0.83	1
3	0.00	0.75	1.00	1.00	1.00	0.67	1
4	-0.17	0.69	1.00	0.44	1.00	1.00	0
5	-0.67	0.94	0.00	0.22	1.00	0.00	0
6	-1.00	1.00	0.00	0.00	0.00	1.00	0

- What about reverse correlation
- Two option to address reverse correlation
 - ❖ Take negative of values
 - ❖ Use negative weight

11/18/2024

pra-sâmi

10

Perceptron – Weights

- Weights – consider importance of each of the feature

id	Threshold	Team Members		Ground		Calculations	Likely	Played	Loss
	w0	x1	w1	x2	w2	$w0+x1*w1+x2*w2$	(y_hat)	(y)	(y-y_hat)^2
1	-1.00	1.00	1.10	1.00	1.00	1.10	1	1	0
2	-1.00	1.00	1.10	0.83	1.00	0.93	1	1	0
3	-1.00	1.00	1.10	0.67	1.00	0.77	1	1	0
4	-1.00	0.44	1.10	1.00	1.00	0.49	1	0	1
5	-1.00	0.22	1.10	0.00	1.00	-0.76	0	0	0
6	-1.00	0.00	1.10	1.00	1.00	0.00	1	0	1

11/18/2024

pra-sâmi

11

Perceptron – Weights and Loss

- Our best solution would be where ground truth and predicted values are same
- Loss is some function of ground truth and predicted values
- And we want it to be cumulative, Square of difference looks promising
 - ❖ $\ell(\hat{y}, y) = (y - \hat{y})^2$
 - ❖ Our overall loss was 2.
- By adjusting weights (w_1, w_2) and threshold (w_0) we can bring the loss to minimum (zero in this case)

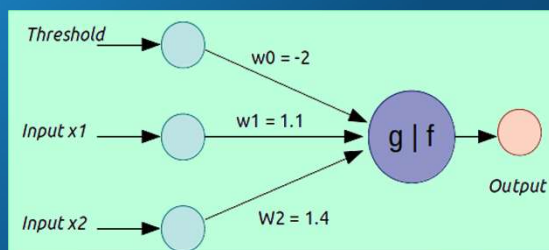
id	Threshold	Team Members		Ground		Calculations	Likely	Played	Loss
	w0	x1	w1	x2	w2	$w_0 + x_1 * w_1 + x_2 * w_2$	(y_hat)	(y)	$(y - y_hat)^2$
1	-2.00	1.00	1.10	1.00	1.40	0.50	1	1	0
2	-2.00	1.00	1.10	0.83	1.40	0.27	1	1	0
3	-2.00	1.00	1.10	0.67	1.40	0.03	1	1	0
4	-2.00	0.44	1.10	1.00	1.40	-0.11	0	0	0
5	-2.00	0.22	1.10	0.00	1.40	-1.76	0	0	0
6	-2.00	0.00	1.10	1.00	1.40	-0.60	0	0	0

11/18/2024

pra-sami

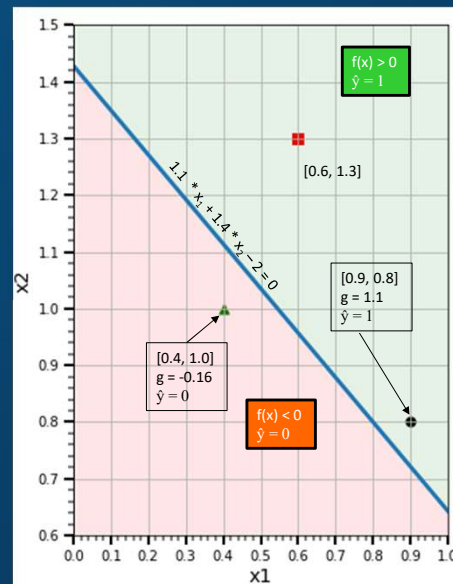
13

Perceptron



- We can represent : $g = w_0 + x_1 * w_1 + x_2 * w_2$
 - ❖ As $g = [x_1, x_2] \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + w_0$
- Given: $W = \begin{bmatrix} 1.1 \\ 1.4 \end{bmatrix}$ and $w_0 = -2$
 - ❖ $g = [x_1, x_2] \cdot \begin{bmatrix} 1.1 \\ 1.4 \end{bmatrix} - 2$
 - ❖ $g = 1.1 * x_1 + 1.4 * x_2 - 2$

11/18/2024

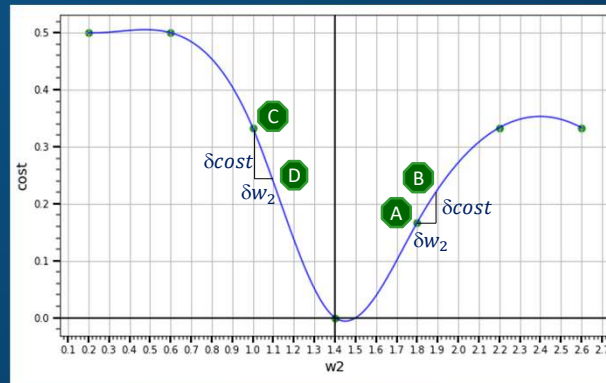


pra-sami

14

Perceptron – Gradient Descent

- w_0, w_1, w_2 need to be adjusted to arrive at most optimal solution i.e. lowest point on the graph.
- Assume that w_0 is fixed at -2, and w_1 at 1.1 and w_2 varies from 0 to 3 (only one variable considered to make plotting simple)
- From point A to B, slope is positive hence w_2 value needs to be decreased
- From point C to D slope is negative hence w_2 needs to be increased.



11/18/2024

pra-sami

15

Perceptron – Activation Function

- So we based our entire calculations on:

$$z = w_0 + x_1 * w_1 + x_2 * w_2$$



But that's an equation of straight line! 😊
What happened to all those 'inhibitory' features?

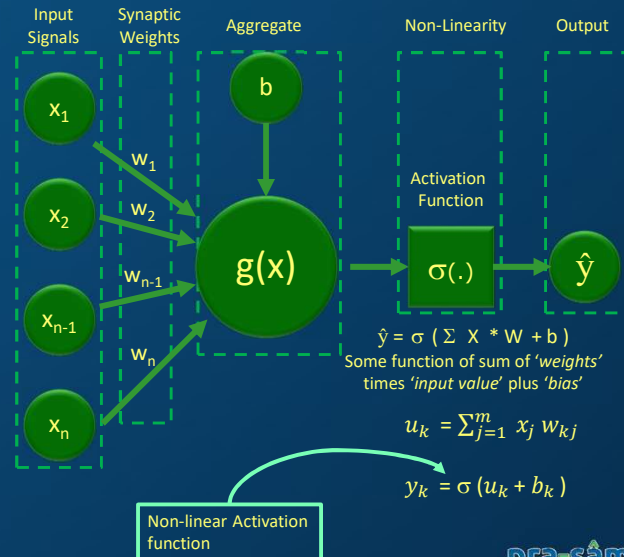
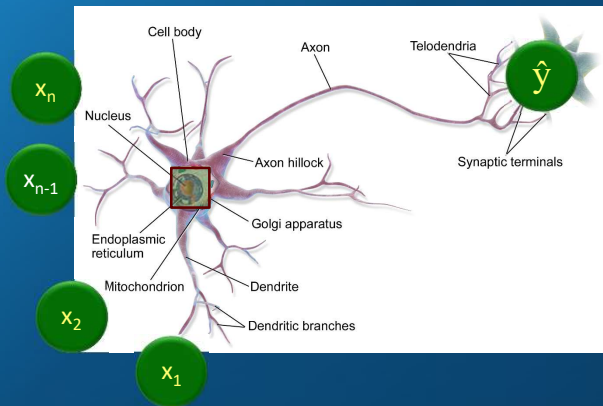


11/18/2024

pra-sami

16

Non Linear Activation function

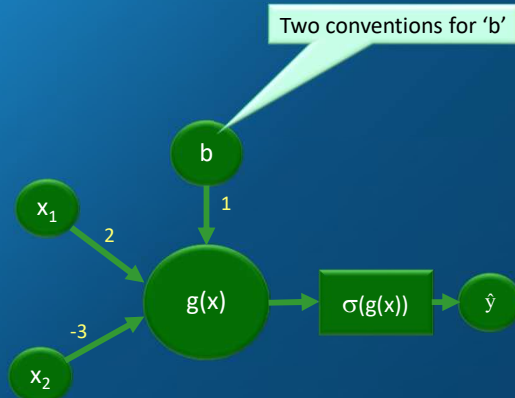


11/18/2024

pra-sami

17

Perceptron with non-linear activation function



□ Given:

$$\diamond W = \begin{bmatrix} 2 \\ -3 \end{bmatrix} \text{ and } b = 1$$

$$\diamond \hat{y} = \sigma \left([x_1, x_2] \cdot \begin{bmatrix} 2 \\ -3 \end{bmatrix} + 1 \right)$$

$$\diamond \hat{y} = \sigma \left(\underbrace{1 + 2 * x_1 - 3 * x_2}_z \right)$$

$$\square \hat{y} = \sigma(z);$$

□ Lets use sigmoid function for σ .

$$\diamond \hat{y} = \frac{1}{(1+e^{-z})}$$

11/18/2024

pra-sami

18

Perceptron with non-linear activation function

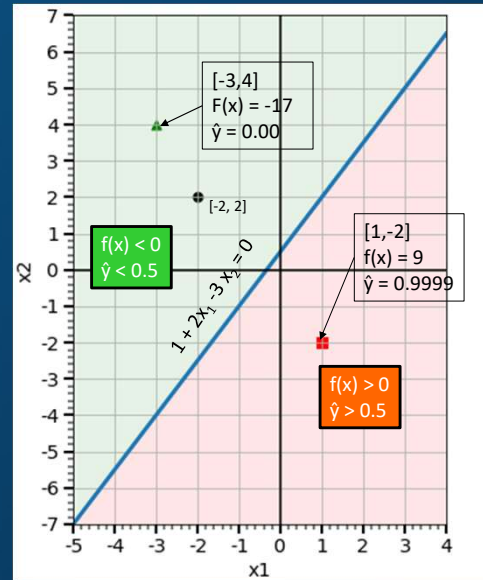
$$\hat{y} = \sigma(1 + 2 * x_1 - 3 * x_2)$$

For $X = [-3, 4]$

- ❖ $\hat{y} = \sigma(1 + 2 * (-3) - 3 * 4)$
- ❖ $\hat{y} = \sigma(1 - 6 - 12)$
- ❖ $\hat{y} = \sigma(-17)$
- ❖ $\hat{y} = 0.0$

Similarly, for $X = [1, -2]$

- ❖ $\hat{y} = \sigma(1 + 2 * 1 - 3 * (-2))$
- ❖ $\hat{y} = \sigma(1 + 2 - 6)$
- ❖ $\hat{y} = \sigma(9)$
- ❖ $\hat{y} = 1.0$



11/18/2024

pra-sami

19

Perceptron with non-linear activation function

$$\hat{y} = \sigma(1 + 2 * x_1 - 3 * x_2)$$

For $X = [-3, 4]$

- ❖ $\hat{y} = \sigma(1 + 2 * (-3) - 3 * 4)$
- ❖ $\hat{y} = \sigma(1 - 6 - 12)$
- ❖ $\hat{y} = \sigma(-17)$
- ❖ $\hat{y} = 0.0$

Are we there yet!

Lets learn some math too!!

Yeehaw!!!

$f(x) > 0$
 $\hat{y} > 0.5$

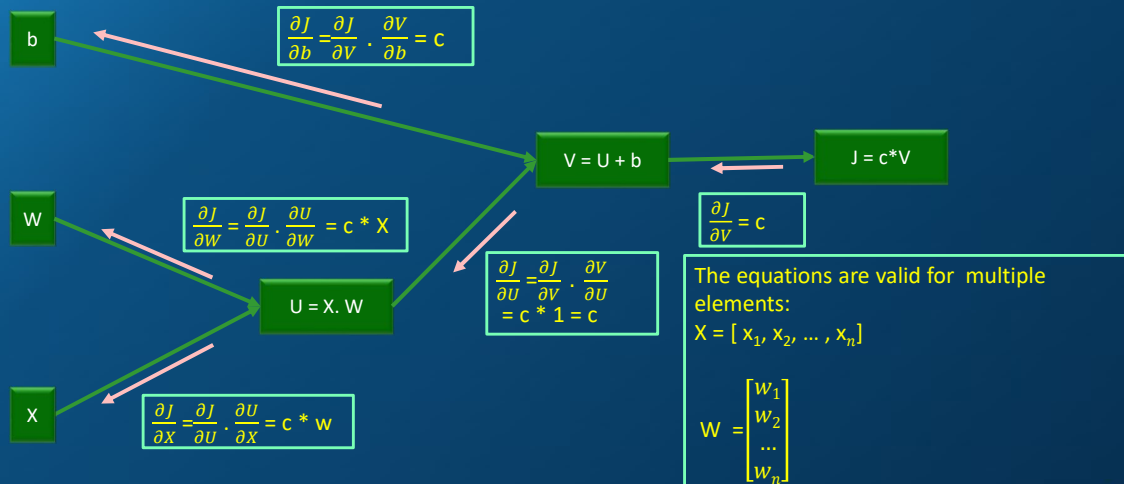
11/19/2024

20

Computational Graph

□ Consider following hypothetical case, basic equation for single neuron :

❖ $\hat{y} = X \cdot W + b$ and Cost is some constant times \hat{y} ; $J = c * \hat{y}$



11/18/2024

pra-sami

21

Exercise 2 : Computational Graph

□ Given a Cost Function J

❖ $J(w, x, b) = 3 * (b + x * w)$

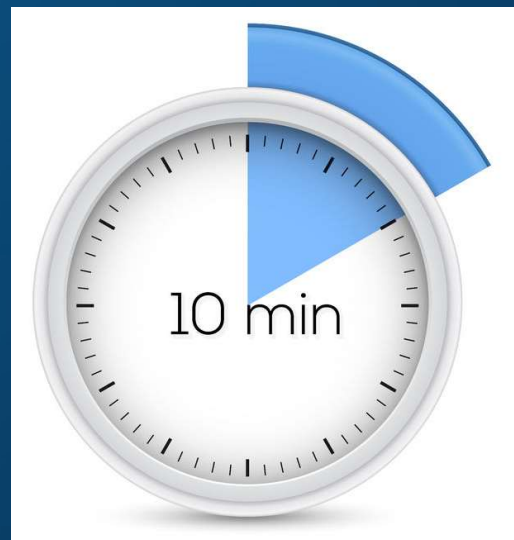
□ Calculate $\frac{\partial J}{\partial w}$, $\frac{\partial J}{\partial x}$ and $\frac{\partial J}{\partial b}$

□ Calculate slope at point :

❖ $b = 6$

❖ $w = 3$

❖ $x = 2$



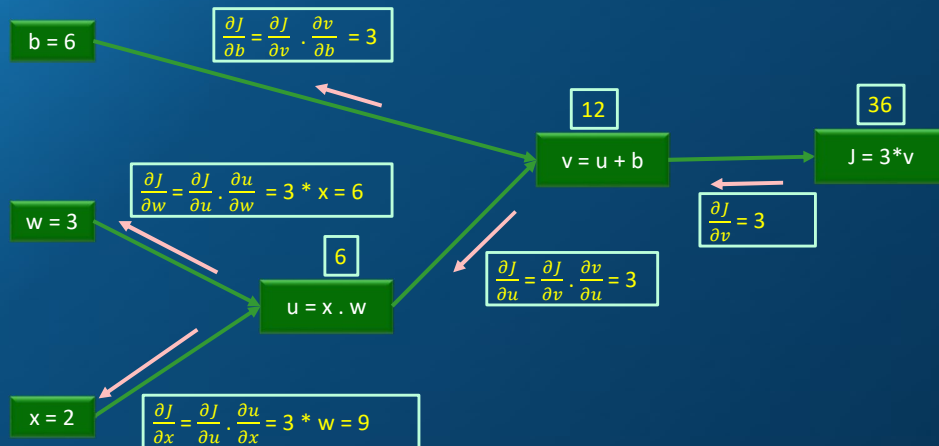
11/18/2024

pra-sami

22

Exercise - Solution

- Given a Cost Function $J(w, x, b) = 3 * (b + w * x)$



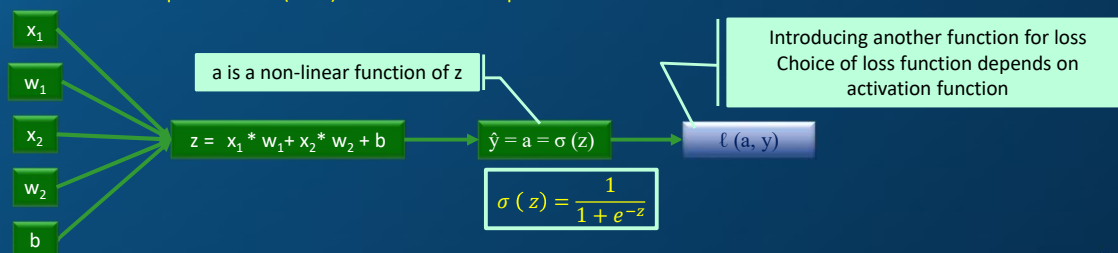
11/18/2024

pra-sâmi

23

Consider Single Path... MLE

- Maximum likelihood estimation, or MLE, is a framework for inference for finding the best statistical estimates of parameters from historical training data
 - ❖ Exactly what we are trying to do with the neural network
- In Classification, output is probability of it belonging to a class
 - ❖ Maximum likelihood estimation, seeks a set of model weights that minimize the difference between the predicted probability distribution and the Ground Truth [cross-entropy]
- In Regression problems:
 - ❖ Use the mean squared error (MSE) loss function or equivalent.



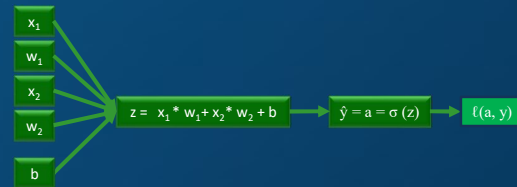
11/18/2024

pra-sâmi

24

Consider Single Path... Loss Function

- A function used to evaluate a candidate solution
- Helps to maximize or minimize the objective function
- Estimates how closely the distribution of predictions made by a model matches the ground truth (maximum likelihood)
- Under maximum likelihood framework, the error between two probability distributions is measured using cross-entropy
 - ❖ Hence $\ell(\hat{y}, y) = - [y * \log(\hat{y}) + (1 - y) * \log(1 - \hat{y})]$



11/18/2024

pra-sami

25

Cost Function

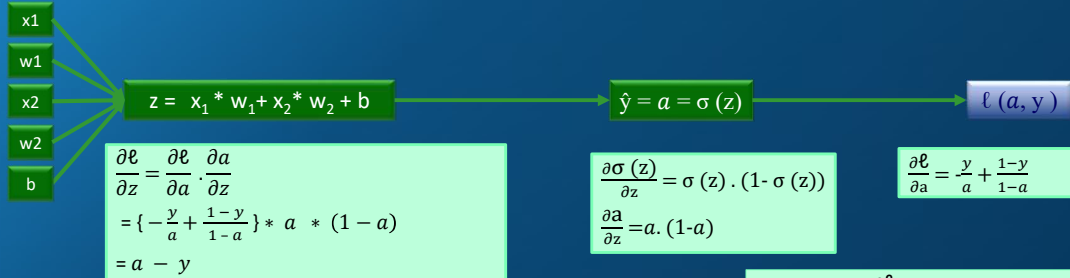
- $\hat{y} = \sigma(\sum W * X + b)$
- Where $\sigma(z) = \frac{1}{1+e^{-z}}$
- Loss function:
 - ❖ A parameter which defines how good our outputs are i.e.
 - ❖ How far our predicted values ' \hat{y} ' (y hat) were from ground truth 'y'
- For logistic regression
 - ❖ $\text{Loss}(\hat{y}, y) = - (y \cdot \log \hat{y} + (1 - y) \cdot \log (1 - \hat{y}))$
 - ❖ Loss function is for an instance
 - ❖ In case of binary classification, $\text{Loss}(\hat{y}, y) = - y \cdot \log \hat{y}$
- Cost Function: Its a sum of losses for all instances
 - ❖ $J(W, b) = \frac{1}{m} (\sum \text{Loss}(\hat{y}, y))$
 - ❖ $= - \frac{1}{m} (\sum (y \cdot \log \hat{y} + (1 - y) \cdot \log (1 - \hat{y})))$
- For binary classification:
 - ❖ $J(W, b) = \frac{1}{m} (\sum \text{Loss}(\hat{y}, y))$
 - ❖ $= - \frac{1}{m} (\sum (y \cdot \log \hat{y}))$

11/18/2024

pra-sami

26

Forward and Back Propagation



$$z = X * W + b$$

$$\hat{y} = a = \sigma(z)$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\ell(a, y) = -[y * \log(a) + (1-y) * \log(1-a)]$$

For binary classification:

$$\ell(a, y) = -y * \log(a)$$

$$\Rightarrow \frac{\partial \ell}{\partial w_1} = x_1 \cdot \frac{\partial \ell}{\partial z} = x_1 \cdot (a-y)$$

$$\Rightarrow \frac{\partial \ell}{\partial w_2} = x_2 \cdot \frac{\partial \ell}{\partial z} = x_2 \cdot (a-y)$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial z} = (a-y)$$



$$w_1 = w_1 - \alpha * \frac{\partial \ell}{\partial w_1} = w_1 - \alpha * x_1 * (a-y)$$

$$w_2 = w_2 - \alpha * \frac{\partial \ell}{\partial w_2} = w_2 - \alpha * x_2 * (a-y)$$

$$b = b - \alpha * \frac{\partial \ell}{\partial b} = b - \alpha * (a-y)$$

Where α is learning rate. The cost function is

$$J(W, b) = \frac{1}{m} * (\sum \ell(a, y))$$

$$\text{Hence } \frac{\partial J}{\partial w_1} = \frac{1}{m} * (\sum \frac{\partial \ell(a, y)}{\partial w_1})$$

11/18/2024

pra-sami

27

So where are the hidden layers!!!

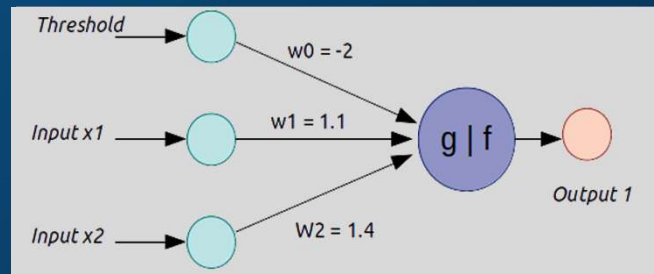
11/18/2024

pra-sami

28

Hidden Layers

id	Threshold	Team Members		Ground	
	x0	x1	w1	x2	w2
1	-2.00	1.00	1.10	1.00	1.40
2	-2.00	1.00	1.10	0.83	1.40
3	-2.00	1.00	1.10	0.67	1.40
4	-2.00	0.44	1.10	1.00	1.40
5	-2.00	0.22	1.10	0.00	1.40
6	-2.00	0.00	1.10	1.00	1.40



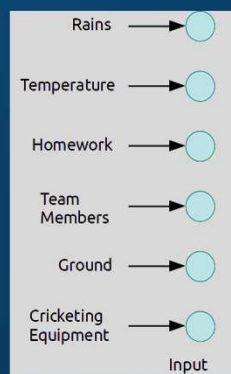
11/18/2024

pra-sâmi

29

Hidden Layers

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0.00	0.00	1.00	1.00	0.00	1.00	1
2	0.00	0.81	1.00	1.00	1.00	0.83	1
3	0.00	0.75	1.00	1.00	1.00	0.67	1
4	-0.17	0.69	1.00	0.44	1.00	1.00	0
5	-0.67	0.94	0.00	0.22	1.00	0.00	0
6	-1.00	1.00	0.00	0.00	0.00	1.00	0



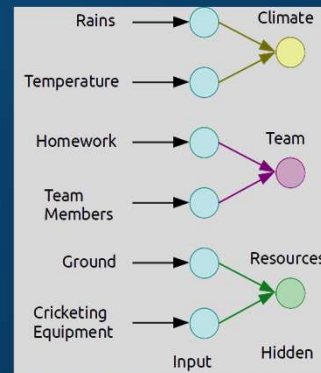
11/18/2024

pra-sâmi

30

Hidden Layers

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0.00	0.00	1.00	1.00	0.00	1.00	1
2	0.00	0.81	1.00	1.00	1.00	0.83	1
3	0.00	0.75	1.00	1.00	1.00	0.67	1
4	-0.17	0.69	1.00	0.44	1.00	1.00	0
5	-0.67	0.94	0.00	0.22	1.00	0.00	0
6	-1.00	1.00	0.00	0.00	0.00	1.00	0



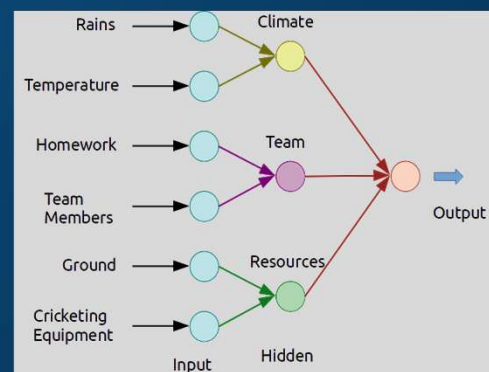
11/18/2024

pra-sami

31

Hidden Layers

id	Rains	Temp	Homework	Team Members	Equipment	Ground	Played
1	0.00	0.00	1.00	1.00	0.00	1.00	1
2	0.00	0.81	1.00	1.00	1.00	0.83	1
3	0.00	0.75	1.00	1.00	1.00	0.67	1
4	-0.17	0.69	1.00	0.44	1.00	1.00	0
5	-0.67	0.94	0.00	0.22	1.00	0.00	0
6	-1.00	1.00	0.00	0.00	0.00	1.00	0

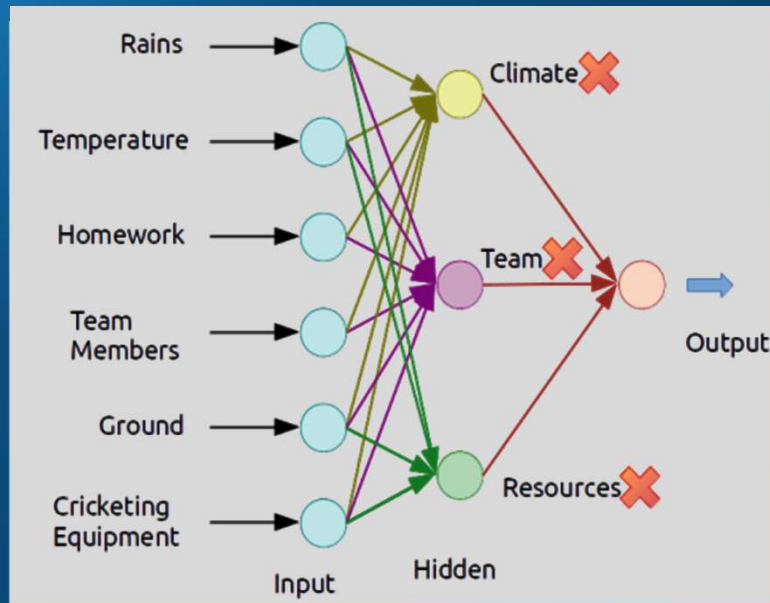


11/18/2024

pra-sami

32

Hidden Layers



11/18/2024

pra-sâmi

33

Hidden Layers

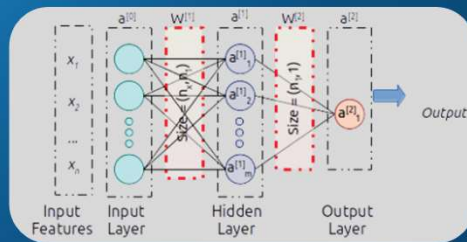


11/18/2024

pra-sâmi

34

Two Major Conventions

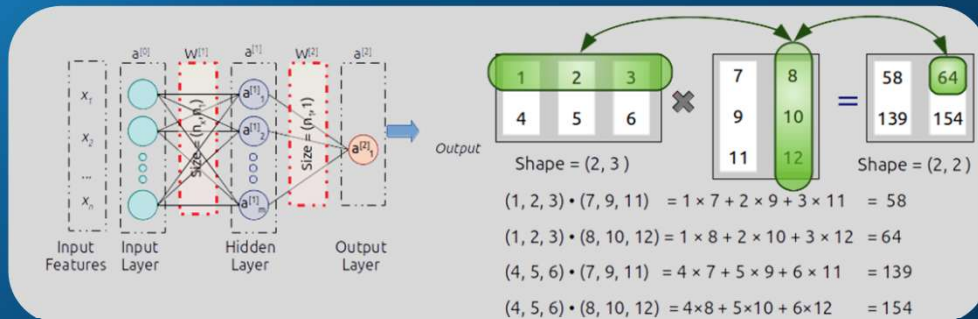


11/18/2024

pra-sami

35

Two Major Conventions

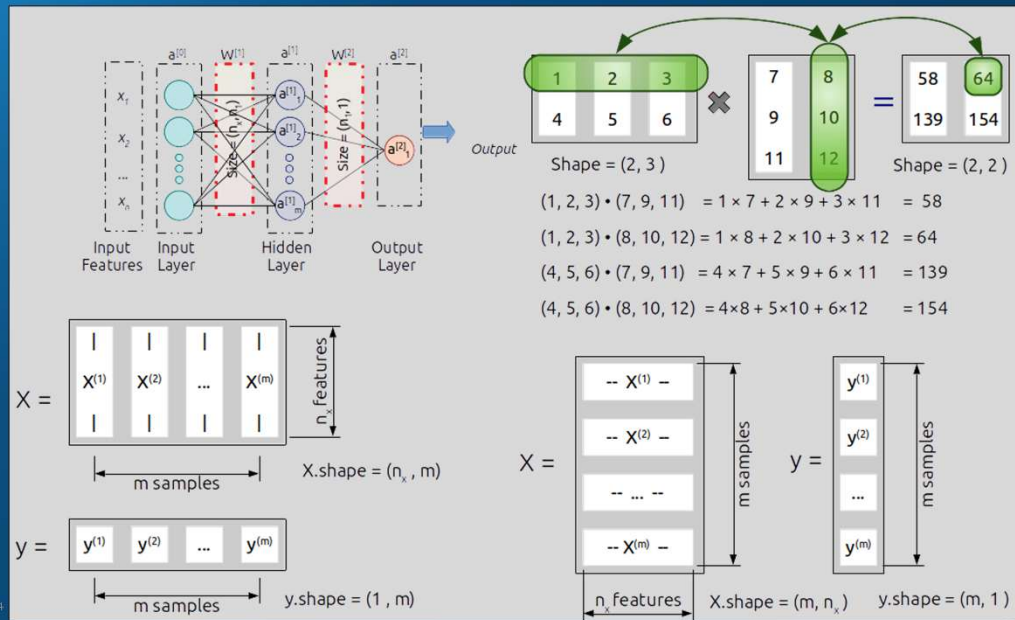


11/18/2024

pra-sami

36

Two Major Conventions

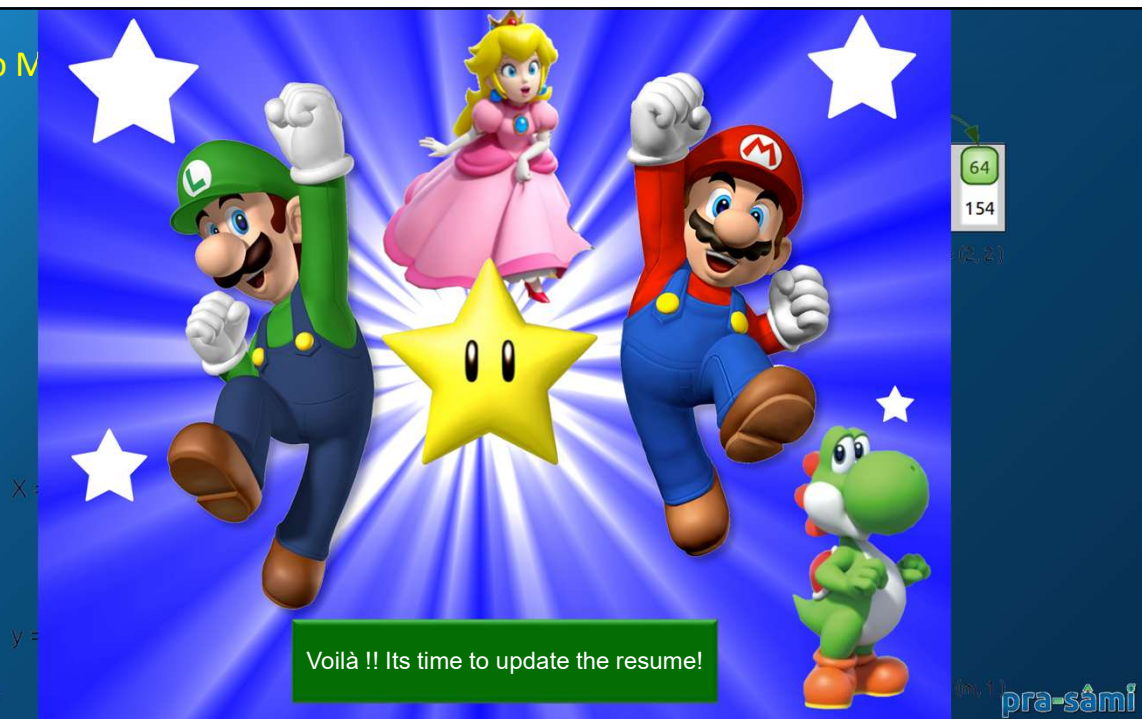


11/18/2024

ra-sâmi

37

Two M



11/18/2024

pra-sâmi

38

Next Session - Coding Perceptron Model in Python

11/18/2024

pra-sâmi

39

THANK YOU

11/18/2024

pra-sâmi

ADDITIONAL MATERIAL

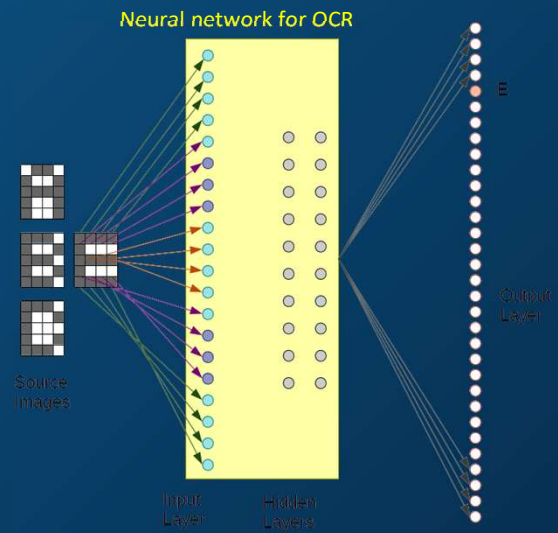


Applications

41

Applications

- The properties of neural networks define where they are useful
- Typical Network
 - ❖ Can learn complex mappings from inputs to outputs, based solely on samples
 - ❖ Difficult to analyse
 - ❖ Firm predictions about neural network behaviour difficult;
 - Unsuitable for safety-critical applications.
 - ❖ Require limited understanding from trainer, who can be guided by heuristics



11/18/2024



42

Applications

- The problem is where

- Typical

- ❖ Can be used for output

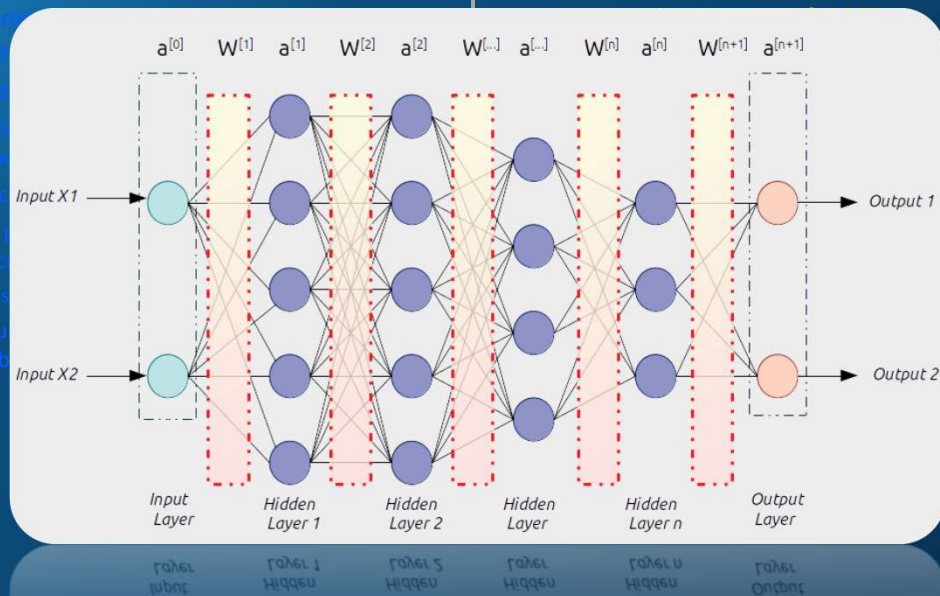
- ❖ Difficult to

- ❖ firm prediction difficult

- Unsatisfactory

- ❖ Requires a lot of data

- can be used for



11/18/2024

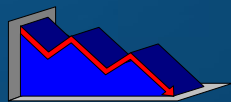
pra-sami

43

Applications

- Stock market prediction

- ❖ "Technical trading" refers to trading based solely on known statistical parameters; e.g. previous price
- ❖ Neural networks have been used to attempt to predict changes in prices.
- ❖ Difficult to assess success or otherwise
 - Since companies using these techniques are reluctant to disclose information.



- Mortgage assessment

- ❖ Assess risk of lending to an individual
- ❖ Difficult to decide on marginal cases
- ❖ Neural networks have been trained to make decisions, based upon the opinions of expert underwriters
- ❖ Neural network produced a 12% reduction in delinquencies compared with human experts



11/18/2024

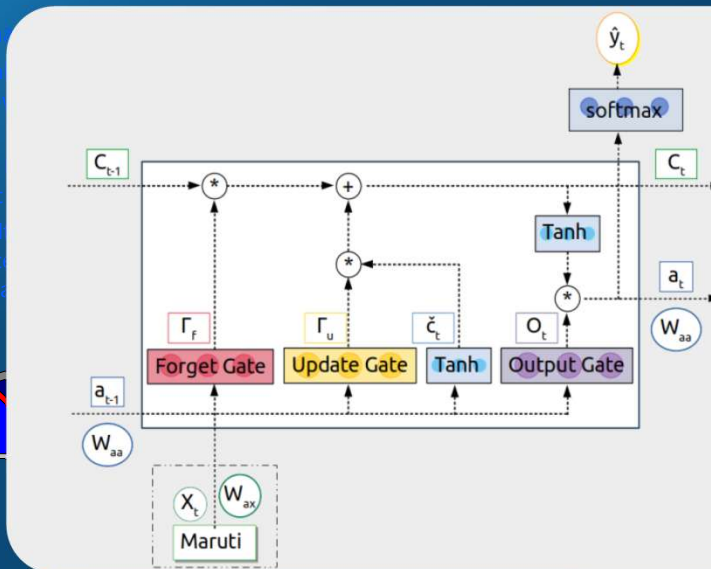
pra-sami

44

Applications

Stock market

- ❖ Technical analysis on known price
- ❖ Neural networks predict
- ❖ Difficult to get these types of information



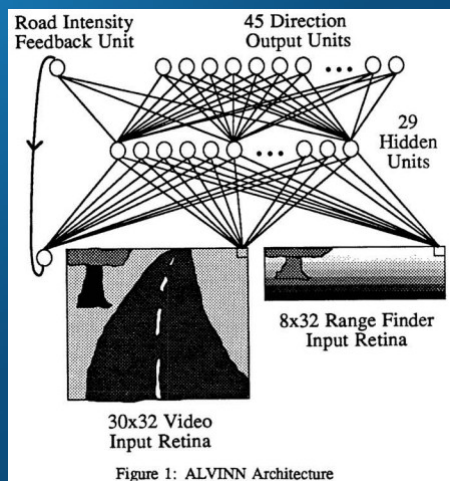
11/18/2024

pra-sami

45

Applications

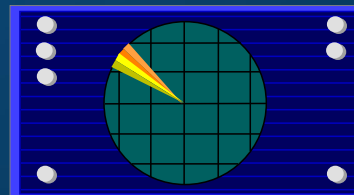
ALVINN: Autonomous Land Vehicle In a Neural Network



11/18/2024

Sonar target recognition

- ❖ Distinguish mines from rocks on sea-bed
- ❖ The neural network is provided with a large number of parameters which are extracted from the sonar signal.
- ❖ The training set consists of sets of signals from rocks and mines.



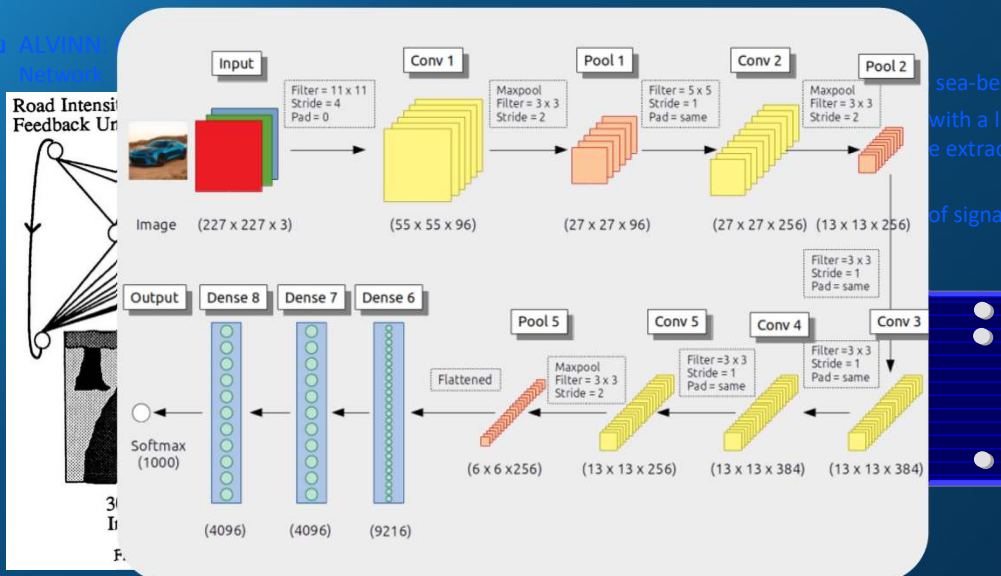
pra-sami

46

Applications

ALVINN: Network

Road Intensity Feedback Unit



11/18/2024

pra-sami

47

Applications

Engine management

- ❖ The behavior of a car engine is influenced by a large number of parameters
 - temperature at various points
 - fuel/air mixture
 - lubricant viscosity.
- ❖ Major companies have used neural networks to dynamically tune an engine depending on current settings



11/18/2024

Signature recognition

- ❖ Each person's signature is different.
- ❖ There are structural similarities which are difficult to quantify.
- ❖ Recognizes signatures to a high level of accuracy.
- ❖ Considers speed in addition to gross shape
- ❖ Makes forgery even more difficult.

pra-sami

48

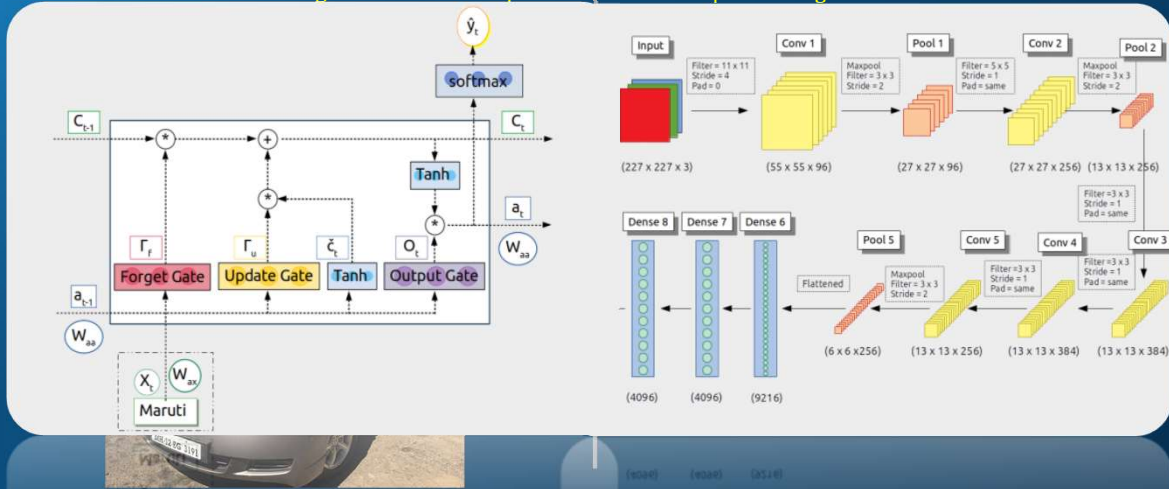
Applications

Engine management

- ❖ The behavior of a car engine is influenced by a

Signature recognition

- ❖ Each person's signature is different.



11/18/2024

pra-sami

49

Derivation of Sigmoid

$$\begin{aligned}
 \partial a &= \partial \sigma(z) \\
 &= \frac{\partial}{\partial z} \left[\frac{1}{1 + e^{-z}} \right] \\
 &= \frac{\partial}{\partial z} (1 + e^{-z})^{-1} \\
 &= -(1 + e^{-z})^{-2} (-e^{-z}) \\
 &= \frac{e^{-z}}{(1 + e^{-z})^2} \\
 &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} \\
 &= \frac{1}{1 + e^{-z}} \cdot \frac{(1 + e^{-z}) - 1}{1 + e^{-z}} \\
 &= \frac{1}{1 + e^{-z}} \cdot \left[\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right] \\
 &= \frac{1}{1 + e^{-z}} \cdot \left[1 - \frac{1}{1 + e^{-z}} \right] \\
 &= \sigma(z) \cdot (1 - \sigma(z)) \\
 &= a \cdot (1 - a)
 \end{aligned}$$

11/18/2024

pra-sami