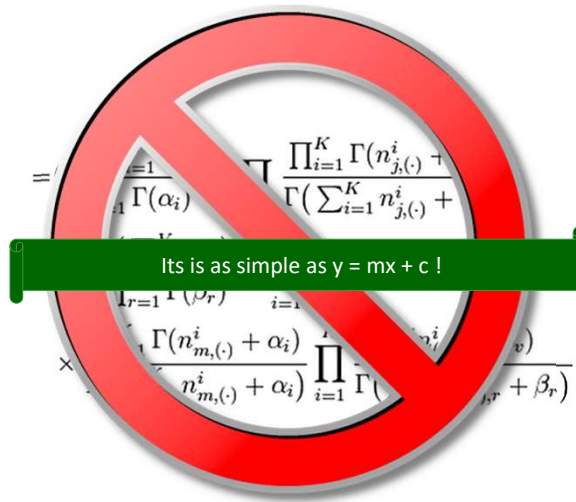


3

Solution to Equation of Perceptron



Ian Goodfellow
Yoshua Bengio
Aaron Courville

11/17/2023

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

11/17/2023

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

11/17/2023

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]

11/17/2023

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 $\sum x_i \cdot w_i < w_0 \Rightarrow \hat{y} = 0$

11/17/2023

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

11/17/2023

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

11/17/2023

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

11/17/2023

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)

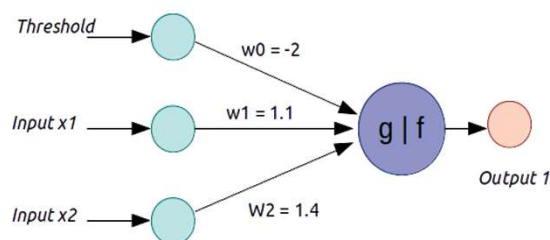
id	Threshold	Team Members		Ground		Calculations	Likely	Played	Loss
	w0	x1	w1	x2	w2	$w_0 + x_1 * w_1 + x_2 * w_2$	(y_hat)	(y)	$(y - y_hat)^2$
1	-2.00	1.00	1.10	1.00	1.40	0.50	1	1	0
2	-2.00	1.00	1.10	0.83	1.40	0.27	1	1	0
3	-2.00	1.00	1.10	0.67	1.40	0.03	1	1	0
4	-2.00	0.44	1.10	1.00	1.40	-0.11	0	0	0
5	-2.00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

11/17/2023

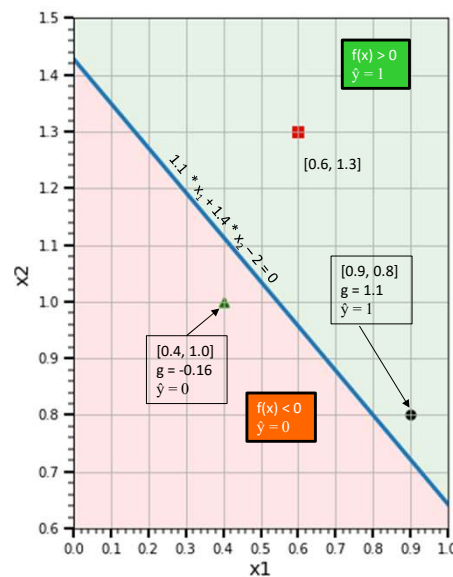
pra-sâmi

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 $w_0 = -2$
 - ❖ $g = [x_1, x_2] \cdot \begin{bmatrix} 1.1 \\ 1.4 \end{bmatrix} - 2$
 - ❖ $g = 1.1 * x_1 + 1.4 * x_2 - 2$



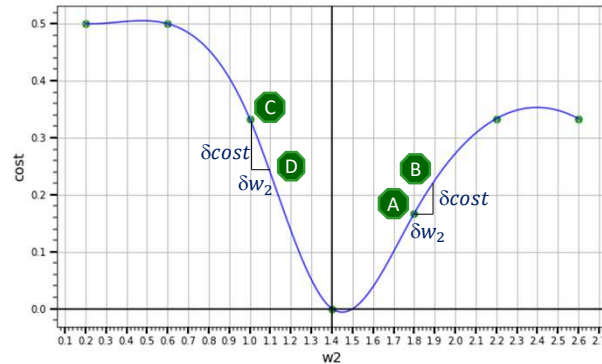
11/17/2023

pra-sâmi

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.



11/17/2023

pra-sâmi

15

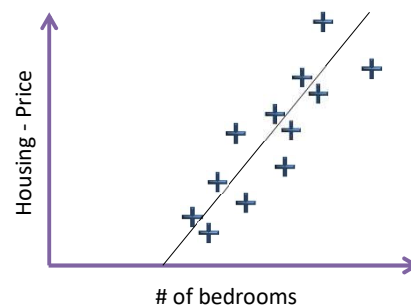
Perceptron – Activation Function

- 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?

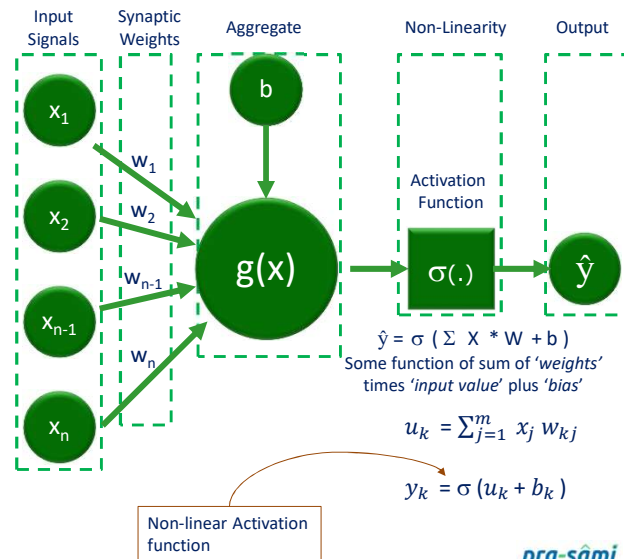
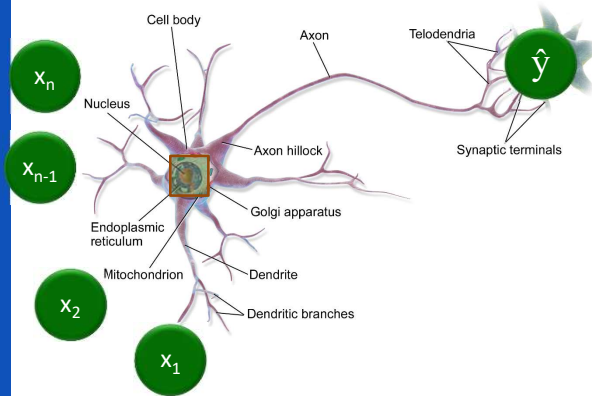


11/17/2023

pra-sâmi

16

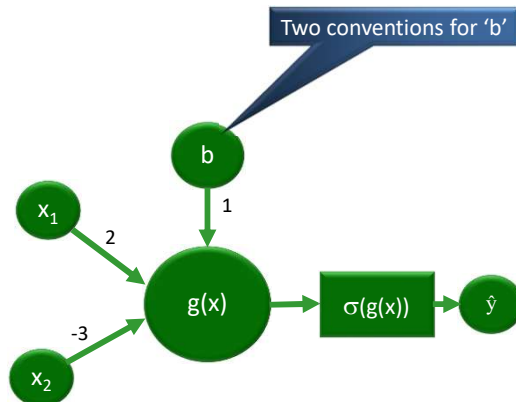
Non Linear Activation function



11/17/2023

17

Perceptron with non-linear activation function



□ Given:

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

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

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

z

□ $\hat{y} = \sigma(z)$;

□ Lets use sigmoid function for σ .

❖ $\hat{y} = \frac{1}{(1+e^{-z})}$

11/17/2023

pra-sâmi

18

Perceptron with non-linear activation function

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

For $X = [-3, 4]$

$$\hat{y} = \sigma(1 + 2 * (-3) - 3 * 4)$$

$$\hat{y} = \sigma(1 - 6 - 12)$$

$$\hat{y} = \sigma(-17)$$

$$\hat{y} = 0.0$$

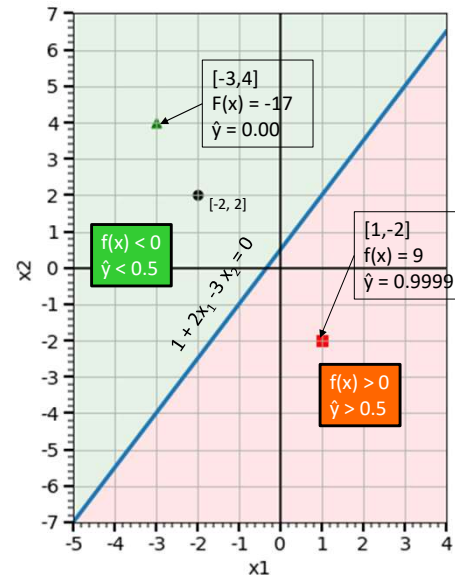
Similarly, for $X = [1, -2]$

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

$$\hat{y} = \sigma(1 + 2 - 6)$$

$$\hat{y} = \sigma(9)$$

$$\hat{y} = 1.0$$



11/17/2023

pra-sâmi

19

Perceptron with non-linear activation function

$$\hat{y} =$$

For

$$\hat{y} =$$

$$\hat{y} =$$

Are we there yet!

Lets learn some math too!!

Yeehaw!!!

$$f(x) > 0$$

$$\hat{y} > 0.5$$

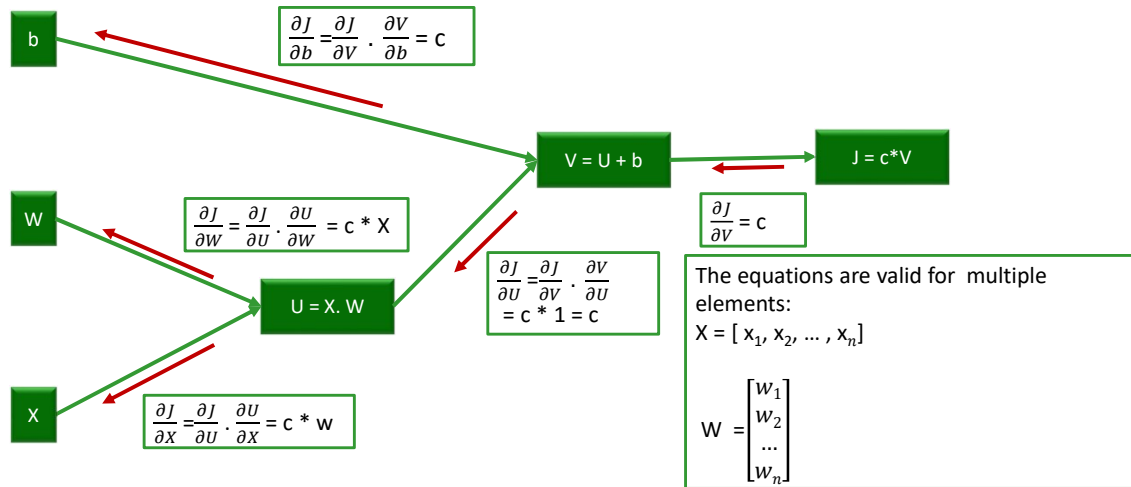
11/17/2023

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} ; $J = c * \hat{y}$



11/17/2023

pra-sâmi

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 :

❖ $b = 6$

❖ $w = 3$

❖ $x = 2$



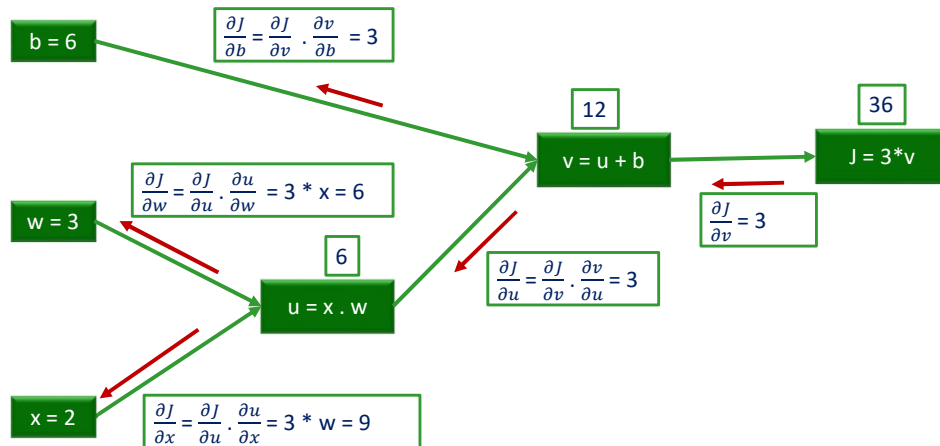
11/17/2023

pra-sâmi

22

Exercise - Solution

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



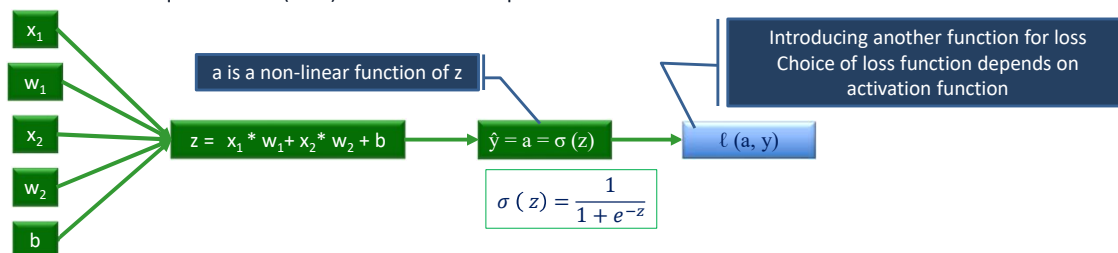
11/17/2023

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.



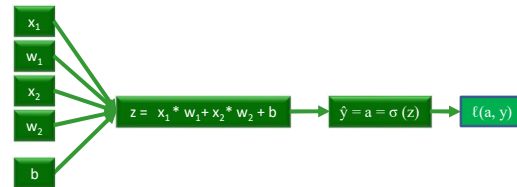
11/17/2023

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 $\ell(\hat{y}, y) = -[y * \log(\hat{y}) + (1 - y) * \log(1 - \hat{y})]$



11/17/2023

pra-sâmi

25

Cost Function

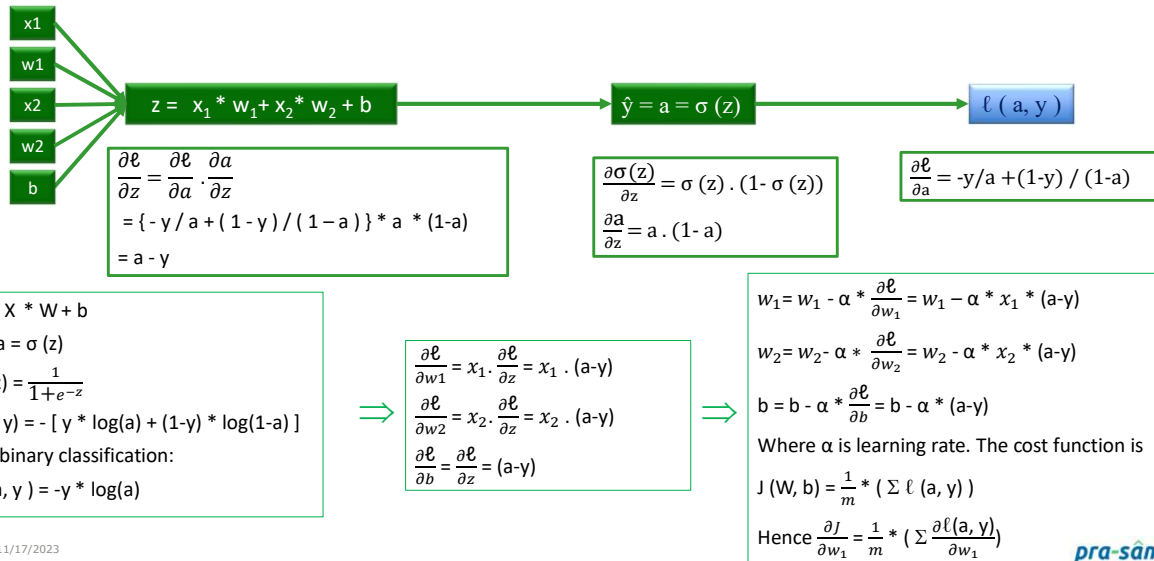
- ❑ $\hat{y} = \sigma(\sum W * X + b)$
- ❑ Where $\sigma(z) = \frac{1}{1 + e^{-z}}$
- ❑ Loss function:
 - ❖ A parameter which defines how good our outputs are i.e.
 - ❖ How far our predicted values ' \hat{y} ' (y hat) were from ground truth 'y'
- ❑ For logistic regression
 - ❖ $\text{Loss}(\hat{y}, y) = -(y \cdot \log \hat{y} + (1 - y) \cdot \log(1 - \hat{y}))$
 - ❖ 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:
 - ❖ $J(W, b) = \frac{1}{m} (\sum \text{Loss}(\hat{y}, y))$
 - ❖ $= -\frac{1}{m} (\sum (y \cdot \log \hat{y}))$

11/17/2023

pra-sâmi

26

Forward and Back Propagation



11/17/2023

pra-sâmi

27

So where are the hidden layers!!!

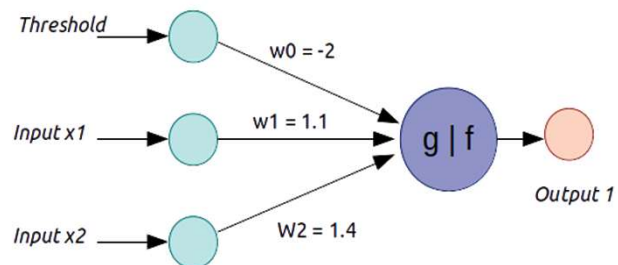
11/17/2023

pra-sâmi

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



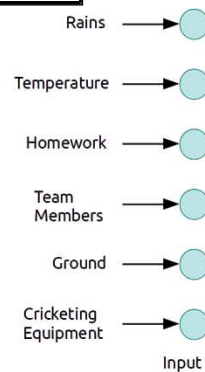
11/17/2023

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



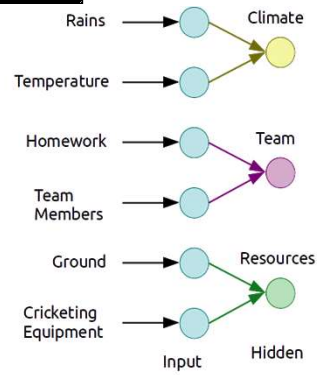
11/17/2023

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



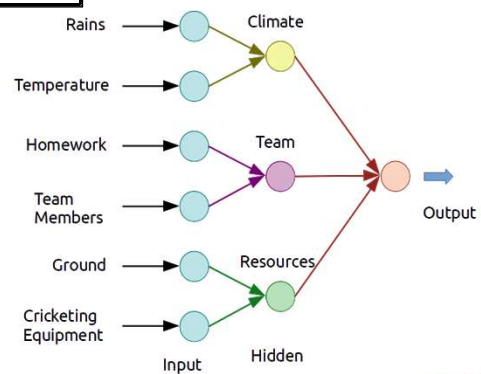
11/17/2023

pra-sâmi

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

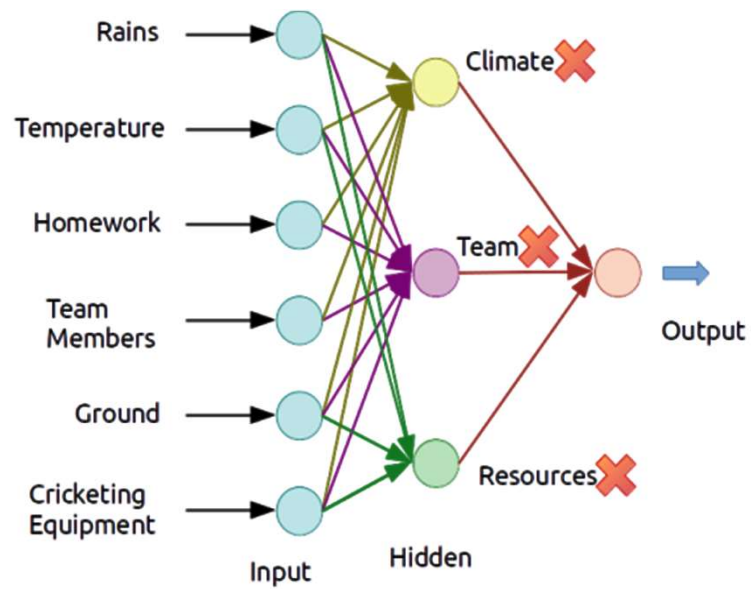


11/17/2023

pra-sâmi

32

Hidden Layers



11/17/2023

pra-sâmi

33

Hidden Layers

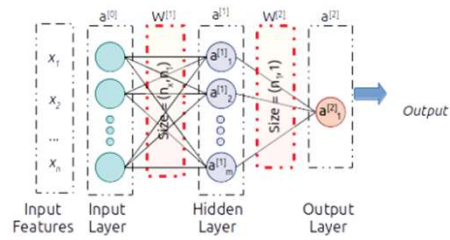


11/17/2023

pra-sâmi

34

Two Major Conventions

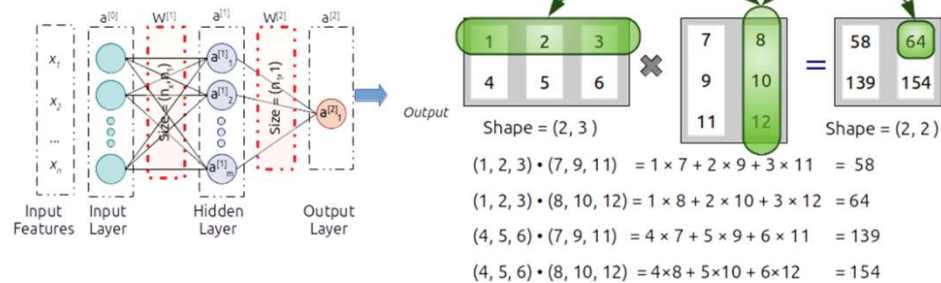


11/17/2023

pra-sâmi

35

Two Major Conventions

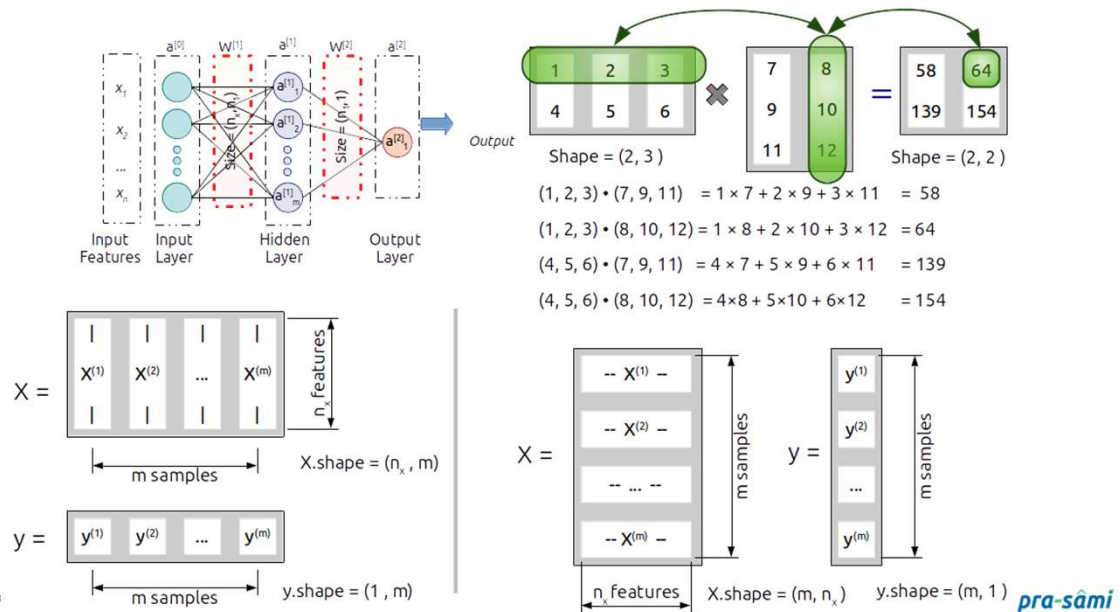


11/17/2023

pra-sâmi

36

Two Major Conventions



37

Two M



38

Reflect...

❑ How many type of layers Deep Learning Algorithms have?

- a. 2
- b. 3
- c. 4
- d. 5

❑ Answer : b

❑ The first layer is called the?

- a. Input Layer
- b. Output Layer
- c. Hidden Layer
- d. None of The Above

❑ Answer : a

❑ Which of the following is/are Limitations of deep learning?

- a. Data labeling
- b. Obtain huge training datasets
- c. Both A and B
- d. None of the above

❑ Answer : c

❑ Deep learning algorithms are _____ more accurate than machine learning algorithm in image classification.

- a. 33%
- b. 37%
- c. 40%
- d. 41%

❑ Answer : d

11/17/2023

pra-sâmi

39

Reflect...

❑ In which of the following applications can we use deep learning to solve the problem

- a. Protein structure prediction
- b. Prediction of chemical reactions
- c. Detection of exotic particles
- d. All of the above

❑ Answer : d

❑ The number of nodes in the input layer is 10 and the hidden layer is 5. The maximum number of connections from the input layer to the hidden layer are:

- a. 50
- b. less than 50
- c. more than 50
- d. It is an arbitrary value

❑ Answer : a

❑ What is a perceptron?

- a. A type of neural network
- b. A reinforcement learning algorithm
- c. A clustering algorithm
- d. A regression algorithm

❑ Answer : a

❑ Who is credited with the invention of the perceptron?

- a. Geoffrey Hinton
- b. Yann LeCun
- c. Frank Rosenblatt
- d. Andrew Ng

❑ Answer : c

11/17/2023

pra-sâmi

40

Reflect...

- ❑ What is the basic building block of a perceptron?
 - a. Neuron
 - b. Weight
 - c. Activation function
 - d. Bias
- ❑ Answer: a

- ❑ In a perceptron, what is the purpose of the activation function?
 - a. To compute the weighted sum of inputs
 - b. To introduce non-linearity
 - c. To adjust the weights during training
 - d. To add a bias to the output
- ❑ Answer: b

- ❑ What is the primary purpose of training a perceptron?
 - a. To optimize the activation function
 - b. To minimize the error in the output
 - c. To increase the number of neurons
 - d. To add more layers to the network
- ❑ Answer: b

- ❑ In a binary classification problem, what is the output of a perceptron?
 - a. Real number
 - b. Probability
 - c. Binary value (0 or 1)
 - d. Vector
- ❑ Answer: c

11/17/2023

pra-sâmi

41

Reflect...

- ❑ What is the perceptron learning rule used for?
 - a. a. Updating weights to reduce prediction error
 - b. b. Adjusting the learning rate during training
 - c. c. Initializing weights in the network
 - d. d. Selecting the appropriate activation function
- ❑ Answer: a

- ❑ What happens if a perceptron is unable to learn a linearly separable function?
 - a. a. It converges quickly
 - b. b. It converges slowly
 - c. c. It never converges
 - d. d. It always converges
- ❑ Answer: c

- ❑ Which of the following statements about the perceptron is true?
 - a. a. It can only be used for linearly separable problems
 - b. b. It is suitable for any type of problem
 - c. c. It can only have one layer
 - d. d. It has no weights
- ❑ Answer: a

- ❑ What is the main limitation of a single-layer perceptron?
 - a. a. It cannot learn non-linearly separable functions
 - b. b. It requires a large amount of training data
 - c. c. It is computationally expensive
 - d. d. It is not suitable for classification tasks
- ❑ Answer: a

11/17/2023

pra-sâmi

42

Next Session - Coding Perceptron Model in Python

11/17/2023

pra-sâmi

43



11/17/2023

pra-sâmi

EXTRA MATERIAL

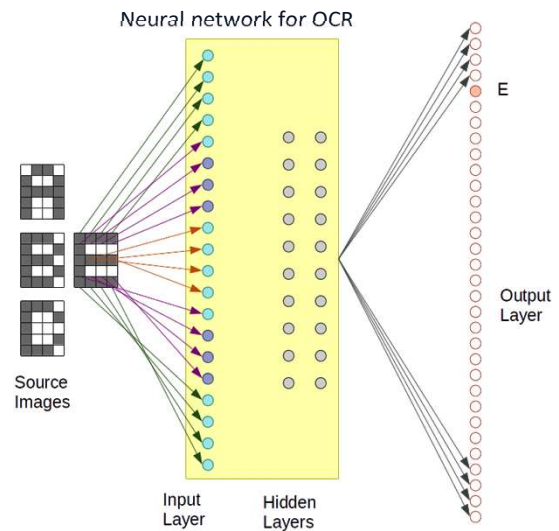
pra-sâmi

Applications

45

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



11/17/2023

pra-sâmi

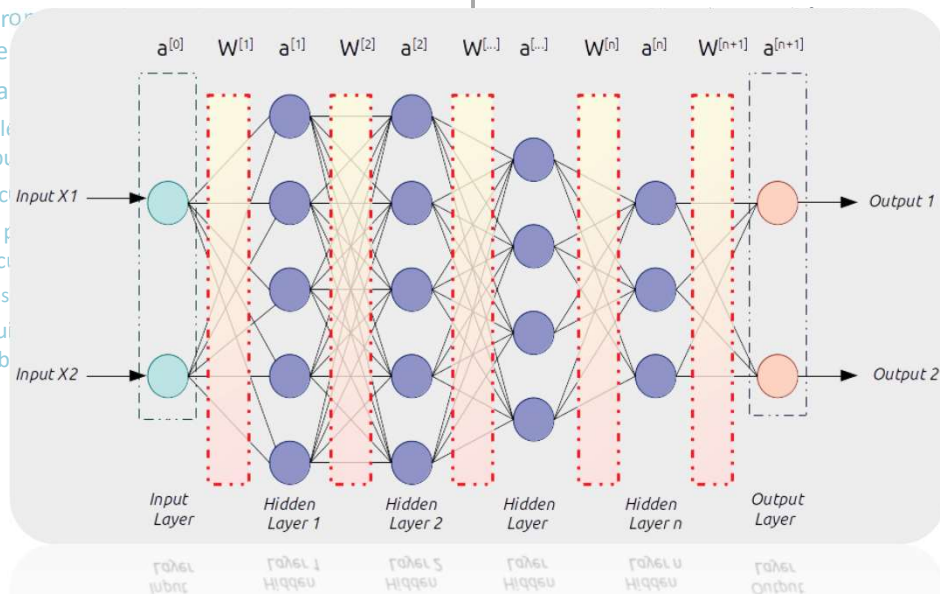
46

Applications

- The problem where

- Typical

- ❖ Can learn from output
- ❖ Difficult to learn from output
- ❖ firm performance difficult to learn from output
- Unsupervised learning
- ❖ Requires a lot of data can be used for



11/17/2023

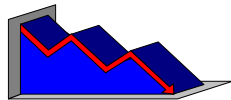
pra-sâmi

47

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



11/17/2023

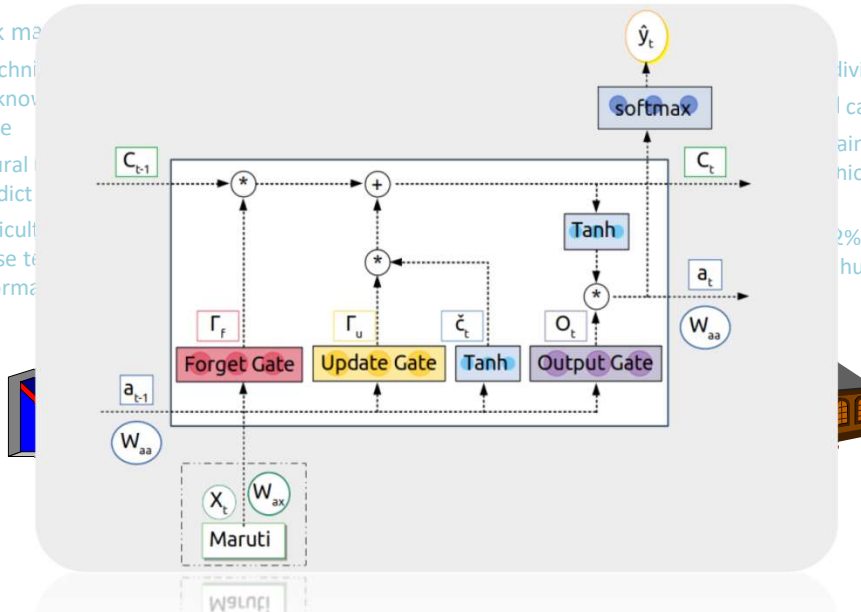
pra-sâmi

48

Applications

□ Stock market

- ❖ "Technical analysis" based on known price
- ❖ Neural networks used to predict
- ❖ Difficult to interpret these technical information



11/17/2023

pra-sâmi

49

Applications

□ ALVINN: Autonomous Land Vehicle In a Neural Network

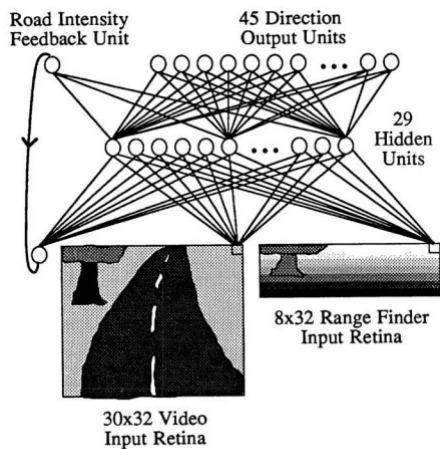
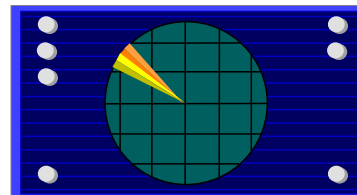


Figure 1: ALVINN Architecture

11/17/2023

□ 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.



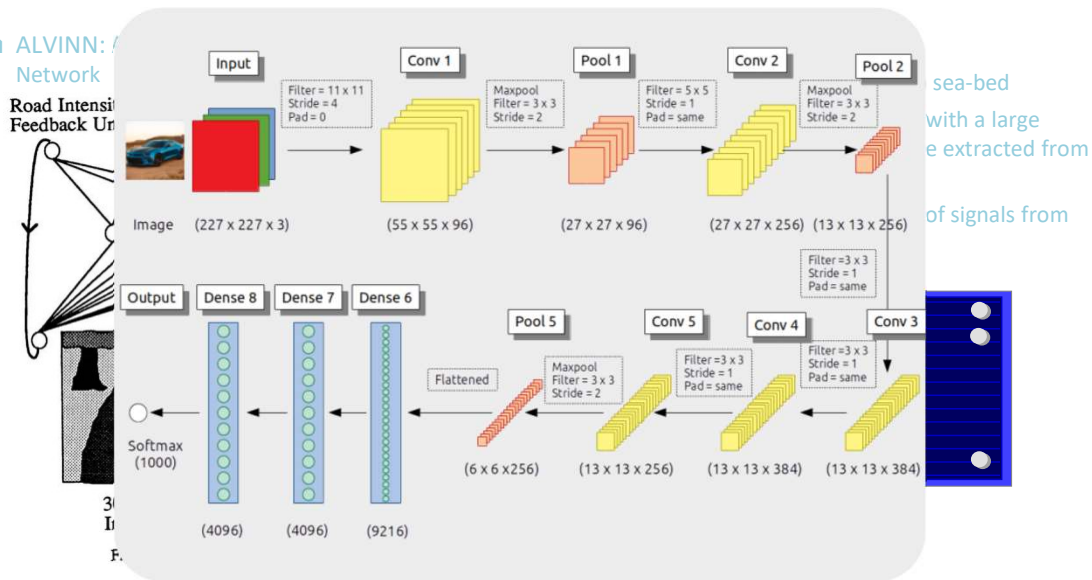
pra-sâmi

50

Applications

ALVINN: A Deep Learning Network

Road Intensity Feedback Unit



11/17/2023

pra-sâmi

51

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



11/17/2023

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.

Pelham

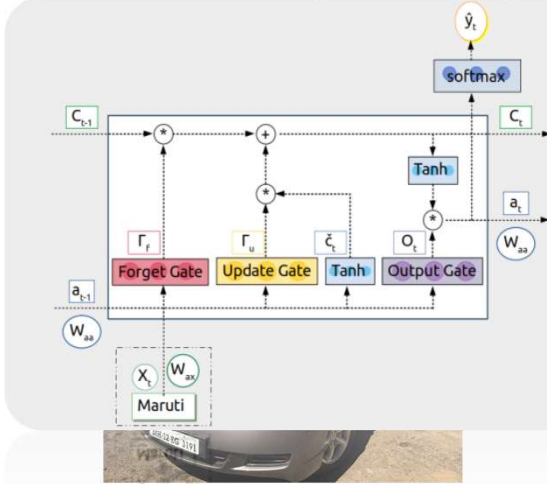
pra-sâmi

52

Applications

□ Engine management

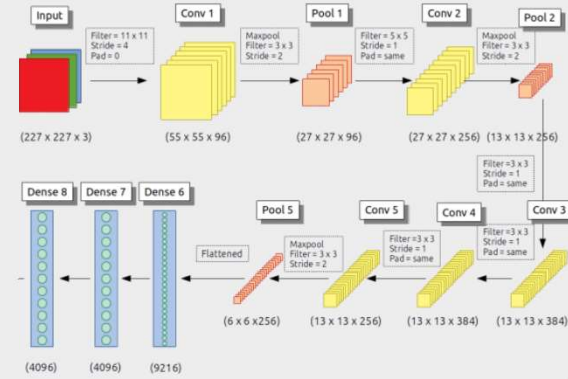
- ❖ The behavior of a car engine is influenced by a



11/17/2023

□ Signature recognition

- ❖ Each person's signature is different.



pra-sâmi

53

Derivation of Sigmoid

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

11/17/2023

pra-sâmi