

Representation

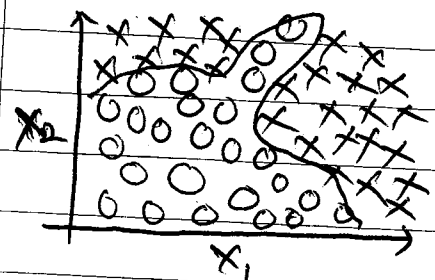
Non-linear hypotheses.

No. _____
Date _____

Example:

Supervised learning classification

Non-linear classification



logistic regression,
(with non-linear features)
 $g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1^2 x_2^2 + \dots)$

Can work well as only 2 features
- x_1 and x_2

If we have like 100 features

- $x_1 = \text{size}$
- $x_2 = \# \text{ bedrooms}$
- $x_3 = \# \text{ floors}$
- $x_4 = \text{age}$
- \dots
- x_{100}

$n=100$

will have a lot of term
example,

$x_1^2, x_1 x_2, x_1 x_3, x_1 x_4, \dots, x_1 x_{100}$
 $x_2^2, x_2 x_3, \dots$

\rightarrow will ≈ 5000 features

num. of quadratic features,
grow order as n^2 $O(n^2)$
closer $\approx \frac{n^2}{2}$

too lot
features may end up over-fitting

Can only include subset of n

$\Rightarrow x_1^2, x_2^2, \dots, x_{100}^2$

Smaller num. features

↳ not enough features and can't
fit like graph on top.

If cubic

$\rightarrow x_1 x_2 x_3, x_1^2 x_2, x_{10} x_{11} x_{12}, \dots$

$O(n^3)$ may end up 170,000 features.
 \hookrightarrow will blow up feature space

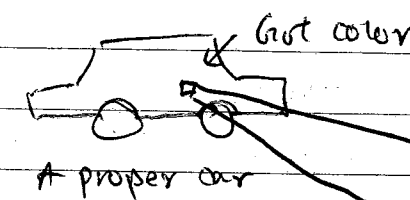
Not a good way to come up additional features
which to build non-linear classifier when n large. \Rightarrow But many ML problems, n is pretty large

Example

want to use ML to train a classifier to examine an
image and tell us whether the image is a car or not.

What is this?

You see this:



Camera see:
104 210 2...
:
:
:

Computer see this:
A matrix of this grid
pixel intensity values
that tell us the brightness
of each pixel in the image

Computer vision problem is to look at the matrix
of pixel intensity values, and tell us that
these numbers represent the door handle of a car.

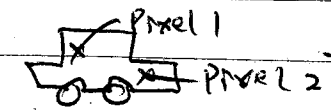
When we use ML to build a car detector, we come up with label
training set

a few label examples of cars

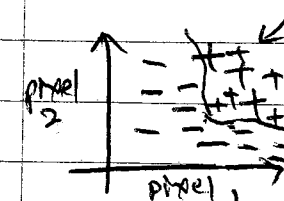
a few label examples of
things are not cars

Then give training set to learning algorithm & train classifier

Example:



Learning Algorithm



+ cars
- non-cars

Plot car at the point of depending
Intensity of pixel 1 and pixel 2
Put another example of car, it
has different intensity of pixel 1 and 2
So at different part of figure

Cars and non-cars
end up lying in different
regions, we need non-linear
hypotheses to separate 2 classes

if only 50 pixels
50x50 pixel images
 $\hookrightarrow 2500$ pixels

$n=2500$ (2500 if RGB)
if using
grayscale image
with colors
 $x = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix}$

learn non-linear hypothesis
by include all quadratic features
 $\hookrightarrow (x_i \times x_j) : \approx 3 \text{ million features}$

Neural Network : Representation

Neurons and the brain

No. _____
Date _____

No. _____
Date _____

Neural Networks

Origins: Algorithms that try to mimic the brain
Very widely used in 80s and ^{early} 90s; popularity diminished in late 90s.

Recent resurgence: State-of-the-art technique for many applications

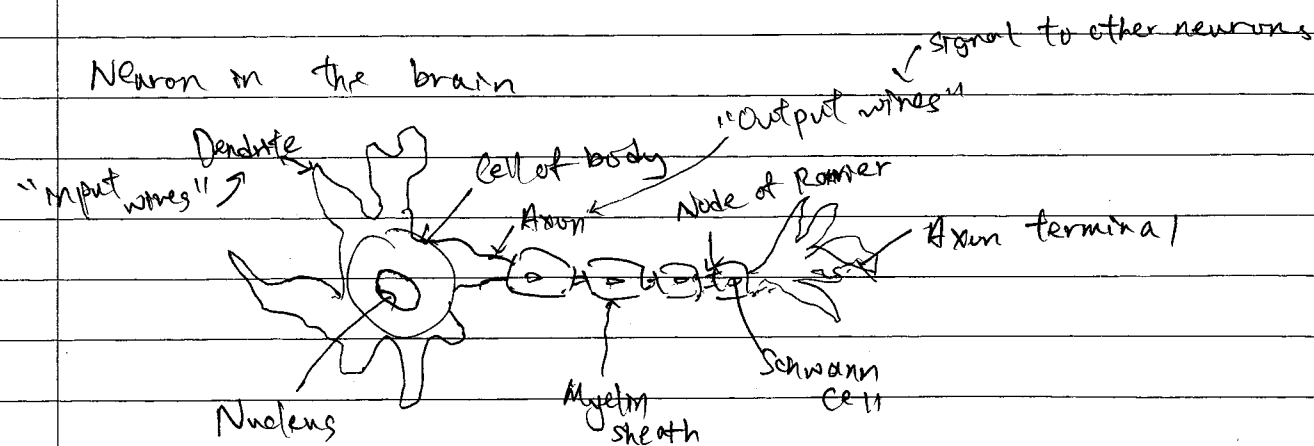
Brain can learn to do a lot of things.

↳ The Brain does it not by mimicking a lot software but just a single learning algorithm.

Model Representation 1

⇒ How we represent neural networks / how we represent hypothesis / represent our model when using neural networks.

Neuron in the brain



A computational unit get number of inputs through dendrite and does some computation and sent outputs via axon to other neurons in the brain.

In artificial neuron network that we implement in the computer, we going to use a simple model on what a neuron does.

Neuron model: Logistic unit

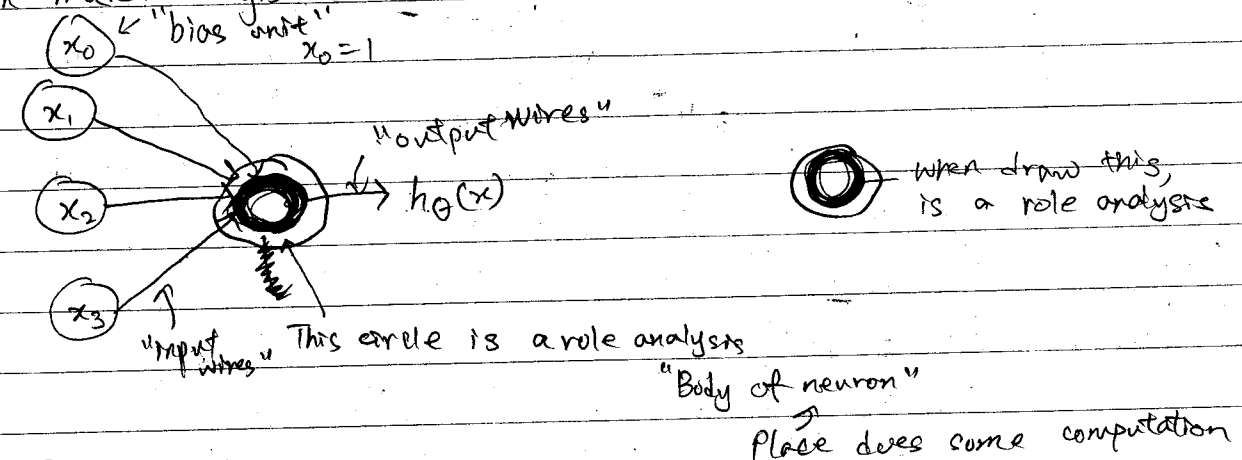
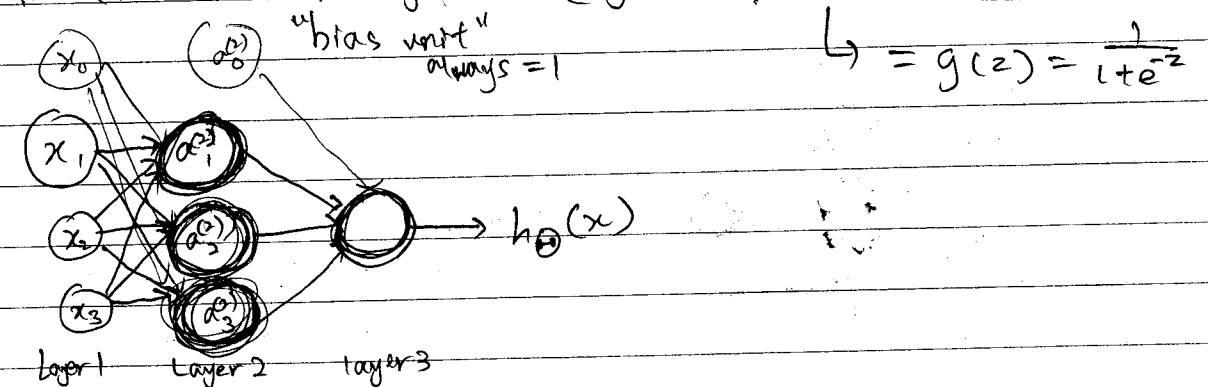


Diagram like this represent a computation or $h_\theta(x) = \frac{1}{1 + e^{-\theta x}}$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

"weight" (parameter) of the model

An artificial neuron with Sigmoid (logistic) activation function.



neural network is just a group of different neurons strong together

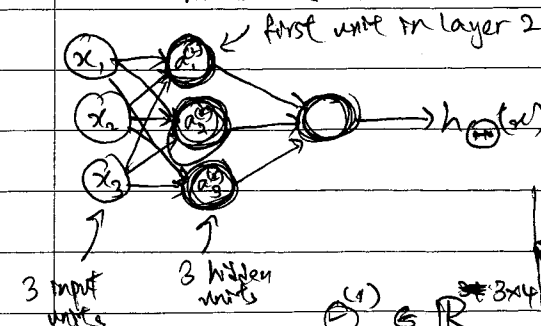
Layer 1 = input layer

Layer 2 = hidden layer

Layer 3 = output layer

all the layer that are not input/output then is hidden.

Neural Network



$a_i^{(j)}$ = "activation" of unit i in layer j

$\Theta^{(j)}$ = matrix of weight controlling function mapping from layer j to layer $j+1$

(example: mapping from layer 1 to layer 2 or layer 2 to layer 3)

neural network is parameterized by these matrices. matrix of parameter governing our mapping from 3 input units to 3 hidden units

Sigmoid activation function

$$a_1^{(2)} = g(\theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \theta_{13}^{(1)}x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(1)}x_0 + \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \theta_{23}^{(1)}x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)}x_0 + \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3)$$

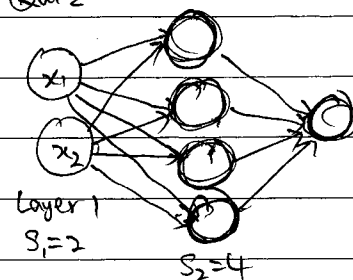
$$h(x) = a_1^{(3)} = g(\theta_{10}^{(2)}a_0^{(2)} + \theta_{11}^{(2)}a_1^{(2)} + \theta_{12}^{(2)}a_2^{(2)} + \theta_{13}^{(2)}a_3^{(2)})$$

Generally:

If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then

$\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$

Quiz



Dimension of

$$\Theta^{(1)} = s_{1+1} \times (s_1 + 1)$$

$$= s_2 \times (2 + 1)$$

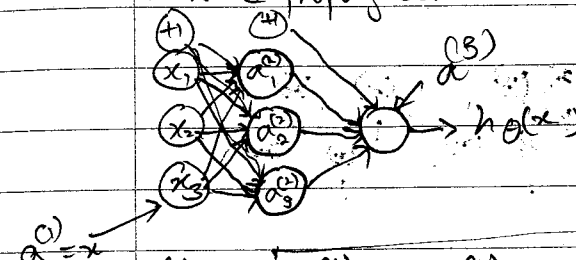
$$= 4 \times 3$$

-vectorized implementation

-intuition on why these neural network representation

from input \Rightarrow hidden \Rightarrow output

Forward propagation: Vectorized Implementation



$$a^{(1)} = g(\theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \theta_{13}^{(1)}x_3)$$

$$a^{(2)} = g(\theta_{20}^{(1)}x_0 + \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \theta_{23}^{(1)}x_3)$$

$$a^{(3)} = g(\theta_{30}^{(1)}x_0 + \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3)$$

$$h(x) = g(\theta_{10}^{(2)}a_0^{(2)} + \theta_{11}^{(2)}a_1^{(2)} + \theta_{12}^{(2)}a_2^{(2)} + \theta_{13}^{(2)}a_3^{(2)})$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

1

2

$$z^{(2)} = \Theta^{(1)} x = \begin{bmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$a^{(2)} = g(z^{(2)}) \quad [a^{(j)} = g(z^{(j)})]$$

activation of apply the sigmoid function element-wise to each of the $z^{(2)}$ elements

For the bias unit, have to

$$\text{Add } a_0^{(2)} = 1 \Rightarrow \text{After adding bias } \rightarrow a^{(2)} \in \mathbb{R}^4$$

$$z^{(3)} = \Theta^{(2)} a^{(2)} \quad [z^{(j+1)} = \Theta^{(j)} a^{(j)}]$$

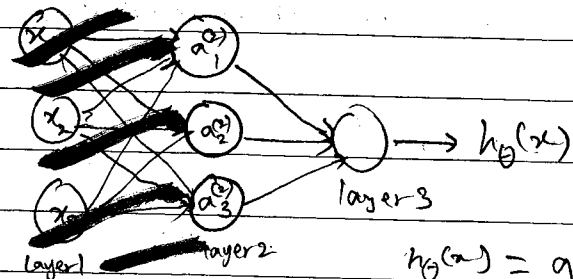
$$h(x) = a^{(3)} = g(z^{(3)}) \quad [h(x) = a^{(j+1)} = g(z^{(j+1)})]$$

\hookrightarrow more efficient way of computing $h(x)$

neural Network learning its own features

No. _____
Date _____

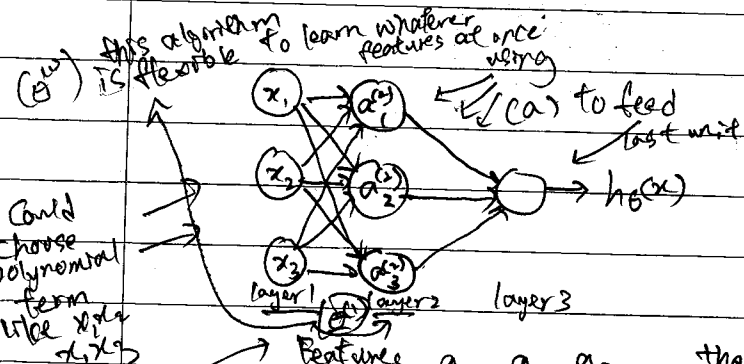
No. _____
Date _____



$$h_{\theta}(x) = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

cover up this part

is just like [logistic regression] exactly same but replace x_1, x_2, x_3 with a_1, a_2, a_3

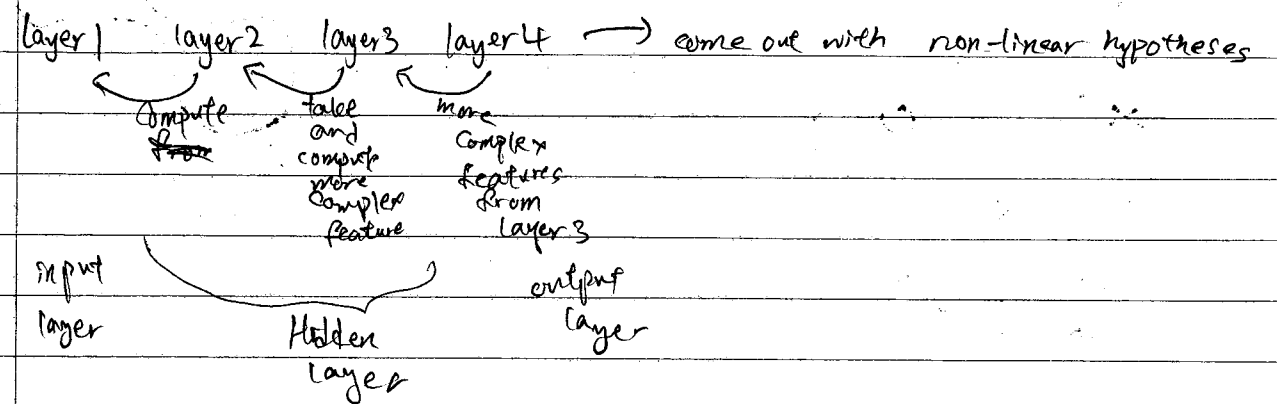
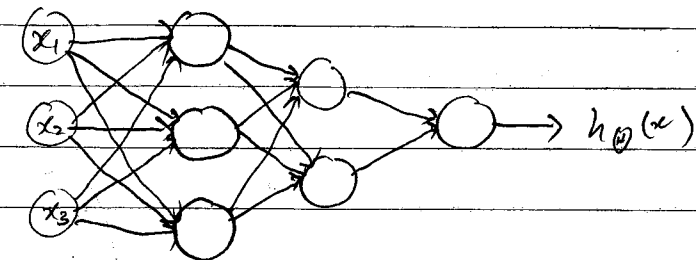


Features a_1, a_2, a_3 , they themselves learn as function of input
function that map from layer 1 \rightarrow layer 2
is determined by other set of parameters $\theta^{(1)}$

instead of being constrained to feed features x_1, x_2, x_3 into logistic regression, it learns its own features a_1, a_2, a_3 to feed into logistic regression
Depends the parameter for $\theta^{(1)}$, will learn some interesting and complex feature.
Parameter 决定你的 feature

The way that neural networks are connected \Rightarrow architecture

other network architectures



Representation

Examples and intuitions

No.

Date

No.

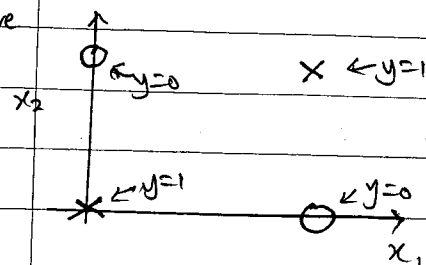
Date

Neural Network is used to learn complex non-linear hypotheses

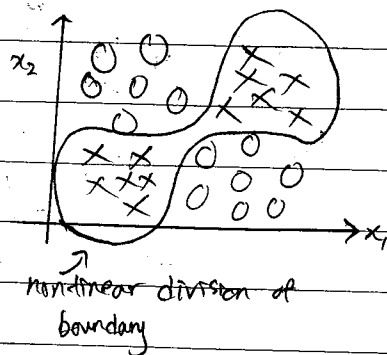
Non-linear classification example: XOR/XNOR

x_1, x_2 are binary (0 or 1)

x - positive
o - negative



Simplified version
of this →



$$y = x_1 \text{ XOR } x_2$$

$$\rightarrow x_1 \text{ XNOR } x_2$$

$$\rightarrow \text{NOT}(x_1 \text{ XOR } x_2)$$

Purpose:

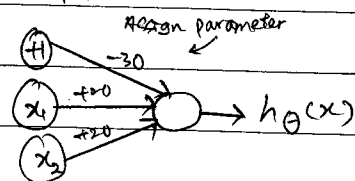
Get neural network
to fit this training set

Before getting to the hard part, we take easy example first.

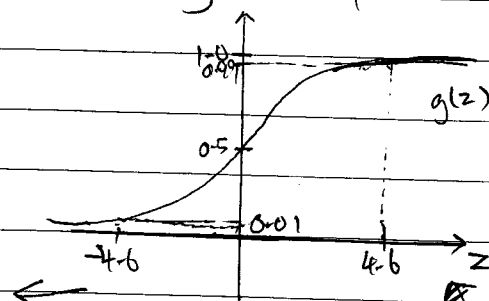
Simple example: AND

$$x_1, x_2 \in \{0, 1\}$$

$$y = x_1 \text{ AND } x_2$$



things
with this
single neuron
network
compute



refer to the
graph

$$h_\theta(x) = g(-30 + 20x_1 + 20x_2)$$

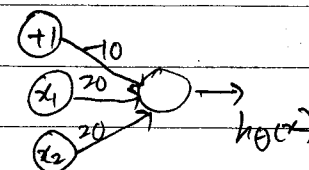
$$\theta_{10} \quad \theta_{11} \quad \theta_{12}$$

x_1	x_2	$h_\theta(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$$h_\theta(x) \approx x_1 \text{ AND } x_2$$

writing out this kind of truth table, we manage
to figure what's logical function our
neural network computes.

Suppose x_1 and x_2 are binary valued (0 or 1) - what boolean function does
the network shown below (approximately) compute?



x_1	x_2	$h_\theta(x)$
0	0	$g(-10) \approx 0$
1	0	$g(10) \approx 1$
0	1	$g(10) \approx 1$
1	1	$g(30) \approx 1$

$$\therefore x_1 \text{ OR } x_2$$

Neural Network Representation:

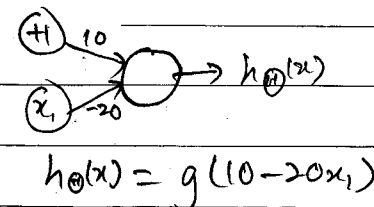
Examples and intuitions 11

No. _____
Date _____

* Neural Network can ~~not~~ be used to compute $x_1 \text{ AND } x_2$, $x_1 \text{ OR } x_2$ when x_1 and x_2 are binary $x_1, x_2 \in \{0, 1\}$

We can have a network to compute negation.

Negation: NOT x_1

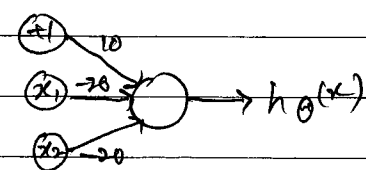


x_1	$h_\theta(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

negate x_1

General idea: Put large negative weight in front of variable you want to negate.

(NOT x_1) AND (NOT x_2)

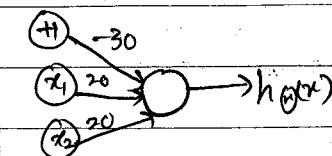
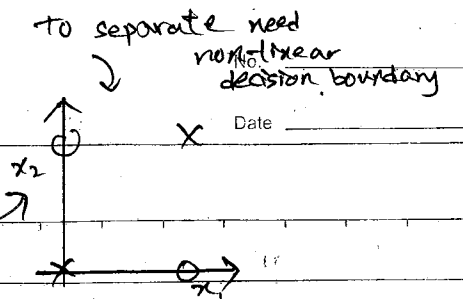


$$h_\theta(x) = g(10 - 20x_1 - 20x_2)$$

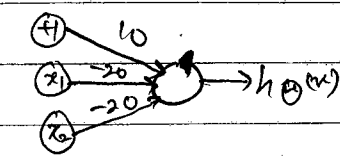
x_1	x_2	$h_\theta(x)$
0	0	$g(10) \approx 1$
1	0	$g(-10) \approx 0$
0	1	$g(-10) \approx 0$
1	1	$g(-30) \approx 0$

≈ 1 if and only if $x_1 = x_2 = 0$

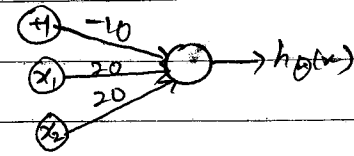
Putting it together: x_1 XNOR x_2



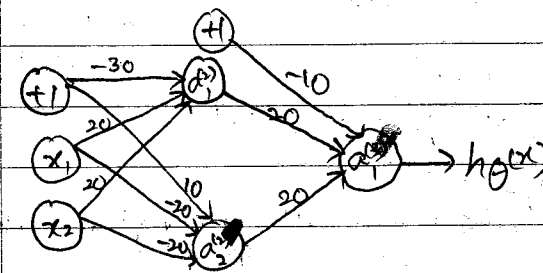
$x_1 \text{ AND } x_2$
 $\theta^{(1)} = [-30 \ 20 \ 20]$



(NOT x_1) AND (NOT x_2)
 $\theta^{(1)} = [10 \ -20 \ -20]$



$x_1 \text{ OR } x_2$
 $\theta^{(1)} = [-10 \ 20 \ 20]$



input layer hidden layer output layer

x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_\theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	1	1

combine

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} a^{(3)} \end{bmatrix} \rightarrow h_\theta(x)$$

transition between 1st and 2nd layer

$$\theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix}$$

transition between 2nd and 3rd layer

$$\theta^{(2)} = [-10 \ 20 \ 20]$$

values for all nodes:

$$a^{(2)} = g(\theta^{(1)} \cdot x)$$

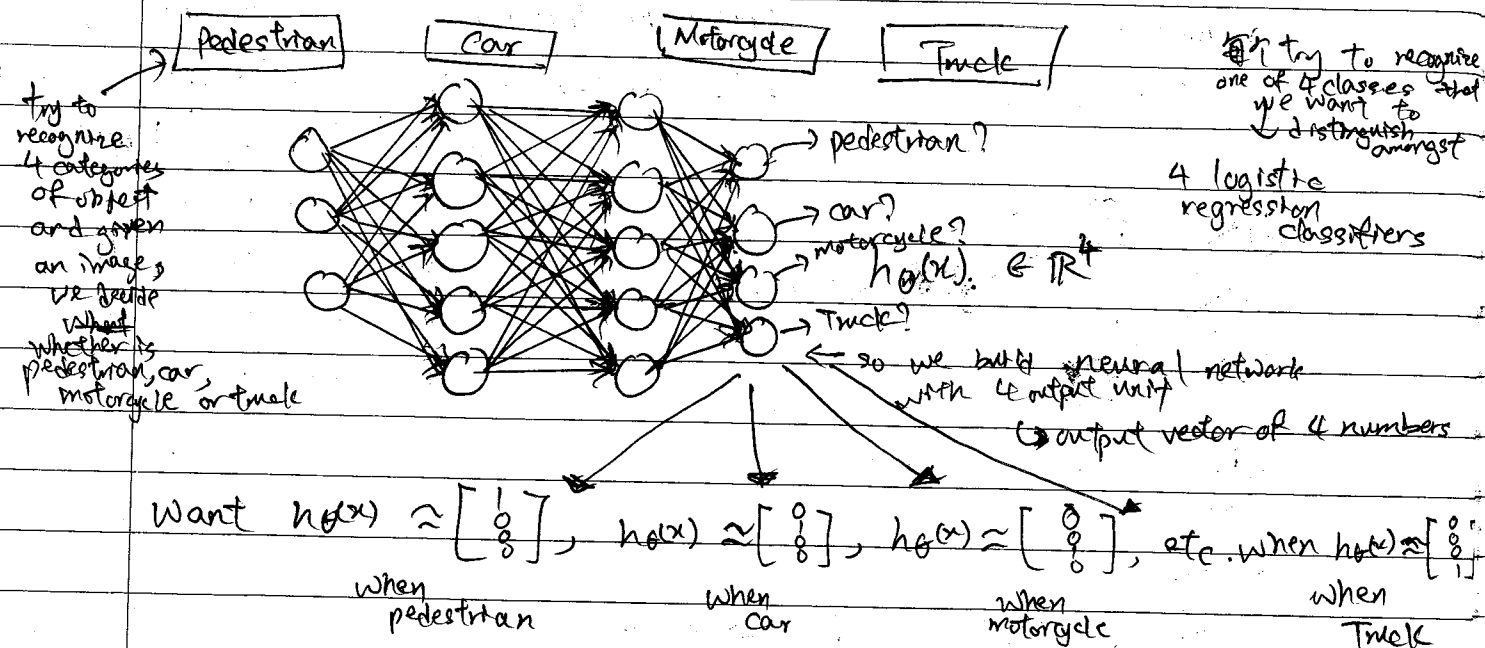
$$a^{(3)} = g(\theta^{(2)} \cdot a^{(2)})$$

$$h_\theta(x) = a^{(3)}$$

Handwritten digit recognition \Rightarrow multi-class classification problem

10 possible categories to recognize the digits from 0-9

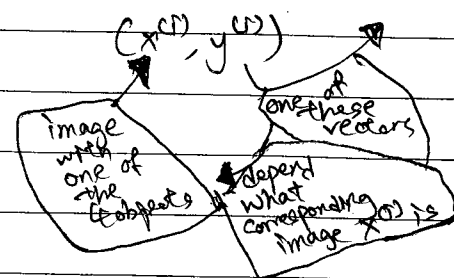
Multiple output units: One-vs-all



Like one-vs-all method in logistic regression

Training set: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

$y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$
pedestrian car motorcycle truck



Not using this

Get a way to get neural network to output value where $h(x^{(i)}) \approx y^{(i)}$

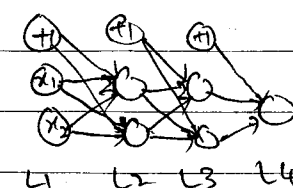
\mathbb{R}^4

Quiz: you have multi-class classification problem with 10 classes, your neural network has 3 layers, and the hidden layer has 5 units, using one-vs-all method, how many element does $\theta^{(2)}$ have?

Answer: 60 (+ bias unit)

Quiz

Consider neural network below, which equation computes activation $a_1^{(3)}$? Note: $g(z)$ is sigmoid activation function.

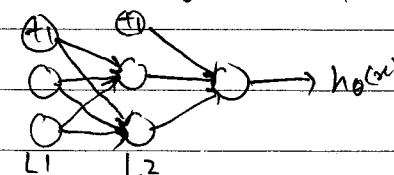


$$a_1^{(3)} = g(\theta_{1,0}^{(2)} a_0^{(2)} + \theta_{1,1}^{(2)} a_1^{(2)} + \theta_{1,2}^{(2)} a_2^{(2)})$$

You are using neural network below and have learned parameters.

$\theta^{(1)} = \begin{bmatrix} 1 & -1.5 & 3.7 \\ 1 & 5 & 1 & 2.3 \end{bmatrix}$ (used to compute $a^{(2)}$) and $\theta^{(2)} = \begin{bmatrix} 1 & 0.6 & -0.8 \end{bmatrix}$ (used to compute $a^{(3)}$)

Suppose you swap parameter for the first hidden layer between its two units as a function of $a^{(2)}$.
so $\theta^{(1)} = \begin{bmatrix} 1 & 5 & 1 & 2.3 \\ 1 & -1.5 & 3.7 \end{bmatrix}$ and also swap output layer so $\theta^{(2)} = \begin{bmatrix} 1 & -0.8 & 0.6 \end{bmatrix}$. How will this change value of output $h(x)$?



It will stay the same.

Setup looks like:

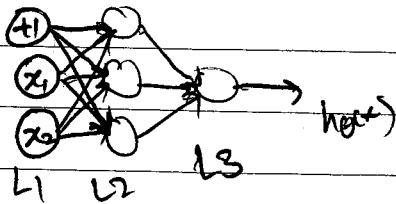
$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ \vdots \\ a_n^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(3)} \\ a_1^{(3)} \\ a_2^{(3)} \\ \vdots \\ a_n^{(3)} \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} h(x)_1 \\ h(x)_2 \\ h(x)_3 \\ h(x)_4 \end{bmatrix}$$

our resulting hypothesis for one set look like this:

$$h(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

in which case our resulting class is the third one down, or $h(x)_3$, which represents motorcycle.

Quiz



You want to compute activations of hidden layer $a^{(2)} \in \mathbb{R}^3$.

One way to do so is the following octave code.

"

%Theta1 is Theta with superscript "(1)" from lecture

%i.e, the matrix of parameters for the mapping from layer 1 (input) to layer 2

%Theta1 has size 3x3

% Assume 'sigmoid' is a built-in function to compute $1/(1+\exp(-z))$

$a2 = \text{zeros}(3,1);$

for $i=1:3$

for $j=1:3$

$a2(i) = a2(i) + x(j) * \text{Theta1}(i,j);$

end

$a2(i) = \text{sigmoid}(a2(i));$

end

"

which one vectorized implementation, correctly compute $a^{(2)}$.

$z = \text{Theta1} * x; a2 = \text{sigmoid}(z);$

Correct statement in neural network.

- 1) If a neural network is overfitting the data, can increase the regularization parameter λ . [larger value of λ will shrink the magnitude of parameter Θ , thereby reducing the chance of overfitting data.]
- 2) Output of a neural network are not probabilities, their sum $(a_1^{(3)} + a_2^{(3)} + a_3^{(3)})$ need not be 1.
- 3) The activation values of the hidden units in a neural network, with sigmoid activation function applied at every layer, are always in range $(0,1)$. The activation function $g(z) = \frac{1}{1+\exp(-z)}$ has range of $(0,1)$.