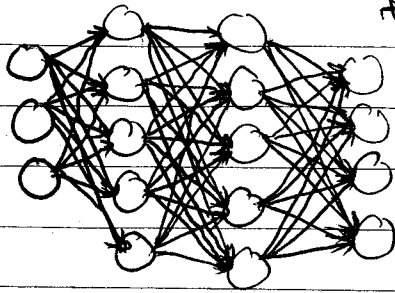


Neural Network:

Learning algorithm
Cost function

→ To fit parameters of a neural network given a training set

(application of neural network to classification problem)



training set $\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

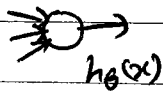
L = total no. of layers in network $L=4$ here

S_l = no. of units (not counting bias unit) in layer l . $S_1=3, S_2=5, S_4=4=S_L$

Binary classification

$y = 0$ or 1

1 output unit



$h_{\theta}(x) \in \mathbb{R}$ (real number)

To simplify,
 $K=1$

$S_L = 1$

Index of final layer
(num. of layer we have in network)

num. of units in output layer

Multi-class classification (K classes)

$y \in \mathbb{R}^K$ E.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$
pedestrian car motorcycle truck

K output units

$h_{\theta}(x) = \mathbb{R}^K$

$S_L = K$ (num. of output unit)

($K \geq 3$)

one-vs-all method when $K \geq 3$

Cost function we use for neural network is generalization of the one we use for logistic regression.

Cost Function

Note:

Double sum adds up the logistic regression costs calculated for each cell in the output layer

Logistic regression:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Neural Network:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[\sum_{k=1}^K y_k^{(i)} \log (h_{\theta}(x^{(i)}))_k + (1-y_k^{(i)}) \log (1-(h_{\theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

Apply to y_k and $(h_{\theta}(x))_k$

Select out i th element of vector that is output by neural network

$(h_{\theta}(x))_i = i$ th output

num K output, we have 4 output units:

For example:
 $y_k = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

en this is sum from $k=1$ to 4 of basically the logistic regression algorithm's cost function but summing that cost function over each of 4 output units in turn.

Regularization term
→ summing all terms

$\theta_{ji}^{(l)}$ for all values of i, j, l except we don't sum over the terms corresponding to bias values (which is 0) because we don't regularize and shrink their values as zero.

because we take K output unit and compare that to the value of y_k which is that one of those vectors saying what cost it should be.

Note: Triple sum adds up the square of all the individual θ s in entire network. This is to make sure that our parameters are not too large.

Quiz:

Suppose we want to minimize $J(\theta)$ as a function of θ , using one of the advanced optimization methods (fminunc, conjugate gradient, BFGS, L-BFGS, etc). What do we need to supply code to compute (as a function of θ)?

Answer: $J(\theta)$ and the (partial) derivative terms $\frac{\partial J(\theta)}{\partial \theta_{ij}}$ for every i, j, L .

Neural Networks:

Learning

Backpropagation algorithm (minimize cost function)

Gradient computation

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h_{\theta}(x^{(i)})_k + (1 - y_k^{(i)}) \log (1 - h_{\theta}(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{i=1}^m \sum_{k=1}^K \sum_{j=1}^n (\theta_{ij}^{(L)})^2$$

try to find parameter theta to

$$\min_{\theta} J(\theta)$$

Want to use gradient descent / advance optimization algorithms,

Need code to compute:

$$-J(\theta)$$

$$-\frac{\partial}{\partial \theta_{ij}^{(L)}} J(\theta)$$

parameter in neural network:

$$\theta_{ij}^{(L)} \in \mathbb{R} \text{ (real number)}$$

Gradient computation

Given one training example (x, y) :

Forward propagation:

$$a^{(1)} = x$$

$$z^{(2)} = \theta^{(1)} a^{(1)}$$

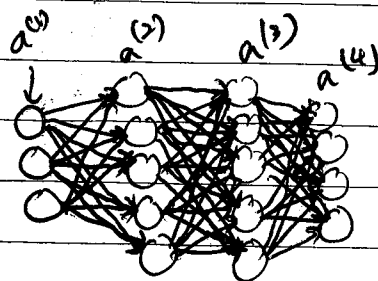
$$a^{(2)} = g(z^{(2)}) \text{ (add } a_0^{(2)})$$

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

$$a^{(3)} = g(z^{(3)}) \text{ (add } a_0^{(3)})$$

$$z^{(4)} = \theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\theta}(x) = a(z^{(4)})$$



L1 L2 L3 L4

In order to compute the derivatives, use algorithm called backpropagation

Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(L)} = \text{"error" of node } j \text{ in layer } L$.

capture error in activation of that neural net

For each output unit (layer $L=4$)

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad (h_{\theta}^{(4)})$$

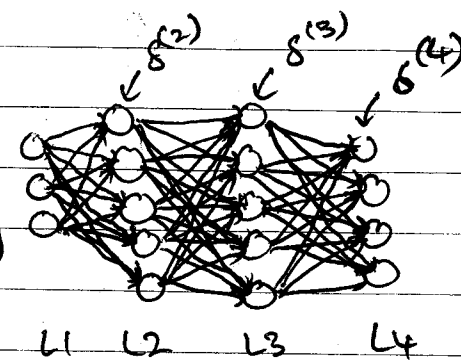
$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)})$$

(No $\delta^{(1)}$)

because the first layer corresponds to input layer, which are just features in training set, so that doesn't have any error associated with that.

vector: $\delta^{(4)} = a^{(4)} - y$
vector whose dimension = num. of output units in network



$a^{(3)} * (1 - a^{(3)})$
vector of activation values for the layer
vector of ones
vector of activation

Backpropagation \Rightarrow come from the fact we start by computing δ for output layer and then go back a layer to compute δ , so back propagating errors from output layer to layer 3 to layer 2

Just few steps of computation, it is possible to prove

$$\frac{\partial}{\partial \theta_{ij}^{(L)}} J(\theta) = a_j^{(L)} \delta_i^{(L+1)} \text{ (ignoring } \lambda \text{ ; if } \lambda=0 \text{)}$$

regularization

Implement backpropagation to compute derivatives with respect to your parameters, when we have large training set.

Backpropagation algorithm

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

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j) (use to compute $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\Theta)$)

capital Greek alphabet delta (δ) → These delta going to be used as accumulators that add things in order to compute

Next, loop through training set

For $j=1$ to m ← working with $(x^{(j)}, y^{(j)})$

set $a^{(1)} = x^{(j)}$

Perform forward propagation to compute $a^{(l)}$ for $l=2, 3, \dots, L$

Using $y^{(j)}$, compute $\delta^{(L)} = a^{(L)} - y^{(j)}$

backpropagation algorithm

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

vectorized notation

$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

After For loop, we compute

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$ if $j \neq 0$

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$ (for bias term, that's why no extra regularization term)

$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

Partial derivative of cost function with respect to each of the parameter

So can use those in gradient descent/advanced optimization algorithm

Quiz: Suppose you have 2 training examples $(x^{(1)}, y^{(1)})$ and $(x^{(2)}, y^{(2)})$. Which of the following is a correct sequence of operations for computing gradient (FP = forward propagation, BP = back propagation).

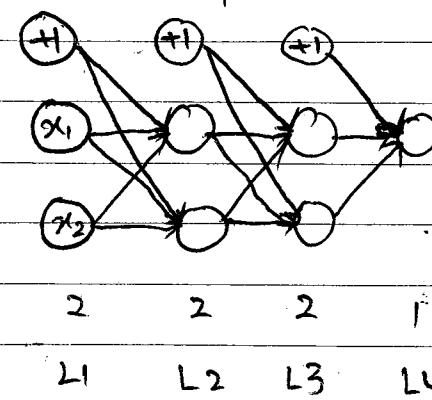
Answer: FP using $x^{(1)}$ followed by BP using $y^{(1)}$. Then FP using $x^{(2)}$ followed by BP using $y^{(2)}$.

Neural Networks:

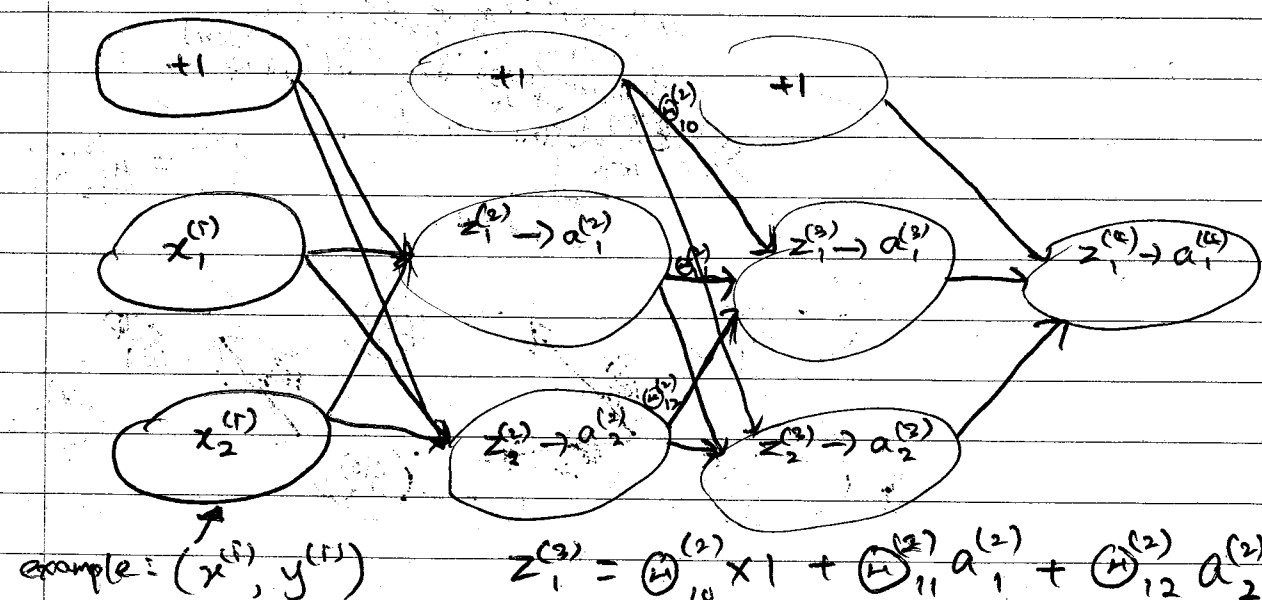
Learning

Backpropagation intuition

Forward Propagation



To illustrate it,



If more than 1 output unit
 $\sum_{k=1}^K$

What is backpropagation doing?

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]$$

Only one output unit
cost function

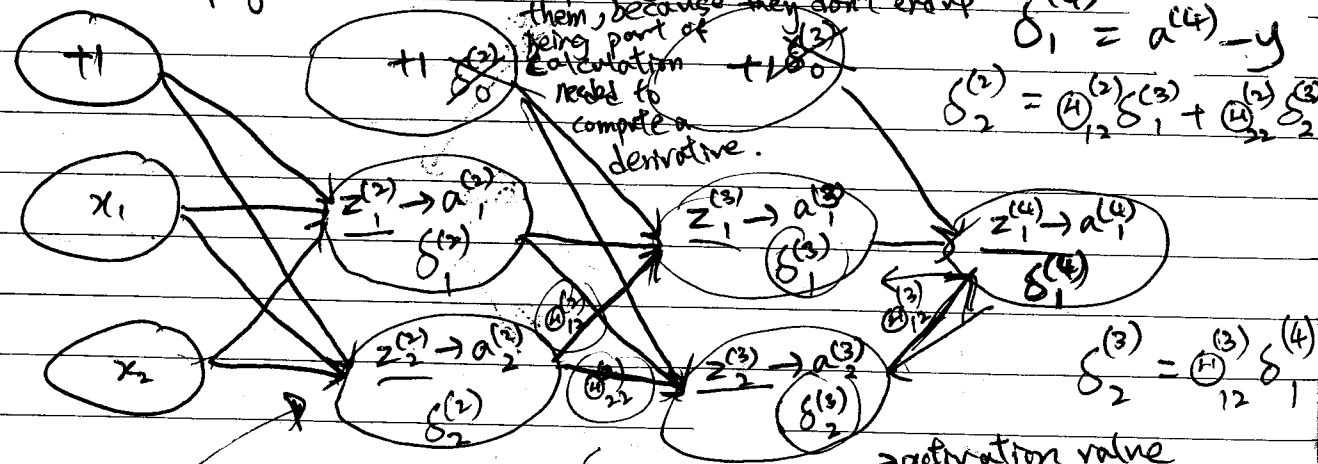
Focusing on a simple example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda=0$)

$$\text{cost}(i) = y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$

(Think of $\text{cost}(i) \approx (h_{\theta}(x^{(i)}) - y^{(i)})^2$)
 i.e. (how well is the network doing on example i ?)

How close is the output to the actual observed label $y^{(i)}$?

Forward Propagation

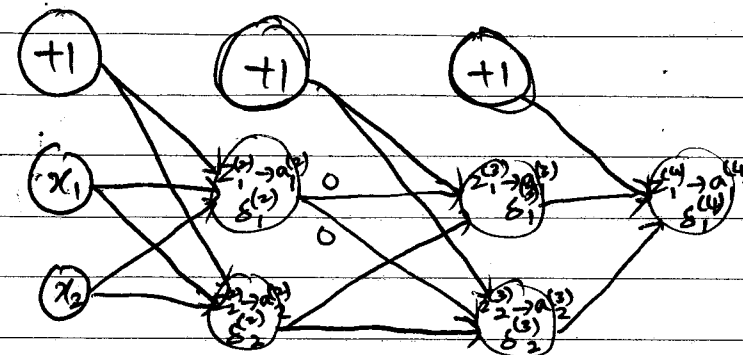


$\delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l)

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ (for $j \geq 0$), where

$$\text{cost}(i) = y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$

Measure of how much would we like to change the neural network's weight in order to affect these intermediate values or computation and affect final output of neural network h_{θ} and affect overall cost



Suppose both of the weights shown in red ($\theta_{11}^{(2)}$ and $\theta_{21}^{(2)}$) are equal to 0. After running BP, what can we say about value of $\delta_1^{(3)}$?

Answer: This is insufficient information to tell.

Neural Networks:

Learning

Implementation note:

Unrolling parameters (from matrices into vectors)

Advanced optimization, return cost function and derivatives
 function $[J, \text{Val}, \text{gradient}] = \text{costFunction}(\text{theta})$ ← input parameters theta
 \dots $\hookrightarrow \mathbb{R}^{n+1}$ $\hookrightarrow \mathbb{R}^{n+1}$ (vectors)

$\text{optTheta} = \text{fminunc}(@\text{costFunction}, \text{initialTheta}, \text{options})$

Neural Network ($L=4$):

$\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ - matrices ($\text{Theta1}, \text{Theta2}, \text{Theta3}$)
 $D^{(1)}, D^{(2)}, D^{(3)}$ - matrices ($D1, D2, D3$)

"Unroll" into vectors (unroll matrices into vector)

so the format could suit and pass into advanced optimization algorithm.

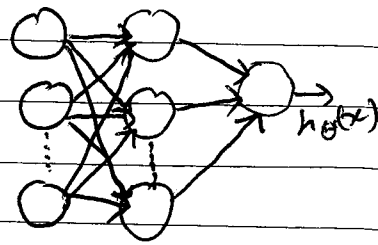
vector \leftarrow Unroll \rightarrow Matrices

Example

$$S_1 = 10, S_2 = 10, S_3 = 1$$

$$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

$$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$$



In octave, if want to convert btwn matrices and vectors

Unroll
into
big long
vector

$$\text{thetaVec} = [\Theta^{(1)}; \Theta^{(2)}; \Theta^{(3)}];$$

$$\text{DVec} = [D^{(1)}; D^{(2)}; D^{(3)}];$$

To reshape: back to matrices

$$\Theta^{(1)} = \text{reshape}(\text{thetaVec}(1:10), 10, 11);$$

$$\Theta^{(2)} = \text{reshape}(\text{thetaVec}(11:20), 10, 11);$$

$$\Theta^{(3)} = \text{reshape}(\text{thetaVec}(21:23), 1, 11);$$

Quiz: Suppose $D^{(1)}$ is a 10×6 matrix and $D^{(2)}$ is a 1×11 matrix,

$$\text{DVec} = [D^{(1)}; D^{(2)}];$$

Which of the following would get $D2$ back from DVec ?

Answer: $\text{reshape}(\text{DVec}(61:71), 1, 11)$

Learning algorithm (How we use unrolling idea to implement learning algorithm)

Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$. (unroll into long vector)

Unroll to get initialTheta to pass to

$\text{fminunc}(@\text{costFunction}, \text{initialTheta}, \text{options})$

function $[\text{Jval}, \text{gradientVec}] = \text{costFunction}(\text{thetaVec})$

From thetaVec , get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$. (reshape)

Use forward prop/backprop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$.

Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec .

Advantage of matrices representation:

- more convenient in forward/back propagation
- easier for vectorized implementations.

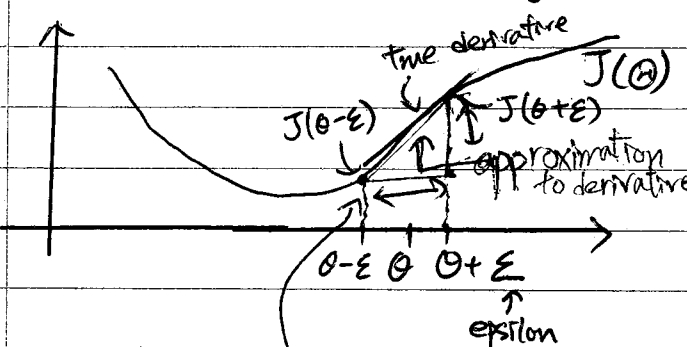
Advantage of vectors representation:

- in advanced optimization algorithms

May get bug in the implementation of backpropagation, use gradient checking to eliminate almost all the problems.

↳ can give implementation of forward/back prop 100% correct.

Numerical estimation of gradients



$$\text{Vertical height} = J(\theta + \epsilon) - J(\theta - \epsilon)$$

$$\text{Horizontal} = 2\epsilon$$

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \leftarrow \text{Two-sided difference (more accurate)}$$

if really small, $\epsilon = 10^{-4}$ (usually make epsilon small)
in mathematically it becomes derivatives $\epsilon \rightarrow 0$, may run into numerical problem.

OR formula like this

$$\frac{J(\theta + \epsilon) - J(\theta)}{\epsilon} \quad \leftarrow \text{call one-sided difference (less accurate)}$$

Implement in Octave:

$$\text{approximate derivative} \quad \text{gradApprox} = (J(\text{theta} + \text{EPSILON}) - J(\text{theta} - \text{EPSILON})) / (2 * \text{EPSILON})$$

Give a numerical estimate of gradient at that point.

Quiz: Let $J(\theta) = \theta^3$. Furthermore, let $\theta = 1$ and $\epsilon = 0.01$, You use formula $\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$ to approximate derivative. What value do you get using this approximation? (When $\theta = 1$, the true, exact derivative is $\frac{d}{d\theta} J(\theta) = 3$).

Answer: 3.0001

Parameter vector θ

$\theta \in \mathbb{R}^n$ (E.g. θ is unrolled version of $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$)
 $\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$

Give you a way to numerically approximate the partial derivative of J respect to θ_i

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon}$$

In Octave $n = \text{dimension of parameter or vector theta.}$

for $i = 1:n$,

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + EPSILON;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - EPSILON;

gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON);

end;

$$\frac{\partial}{\partial \theta_i} J(\theta)$$

Use this in neural network implementation

Check that gradApprox \approx DVec

derivative get from backprop

to check implementation of backprop is correct

Use this into gradient descent / advanced optimization algorithm

make us more confident that we're computing the derivatives correctly

Implementation Note:

- Implement backprop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$)
- Implement numerical gradient check to compute grad Approx.
- Make sure they give similar values.
- Turn off gradient checking. Using backprop code for learning. ^{First} ^{then} after checking, turn off g-checking and do BP

Important:

- Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of costFunction(...)) your code will be very slow.

Our

Main reason we use BP ~~code~~ rather than numerical gradient computation method during learning.

- The numerical gradient algorithm is ~~very~~ very slow.

Neural Networks:

Learning

Random initialization

Initial value of Θ

For gradient descent and advanced optimization method, need initial value for Θ

optTheta = fminunc(@costFunction, initialTheta, options)

Consider gradient descent

Set initialTheta = zeros(n, 1)?

initialize Θ as all vector zeros only work ok in logistic regression but not in training the network

Neural Networks:

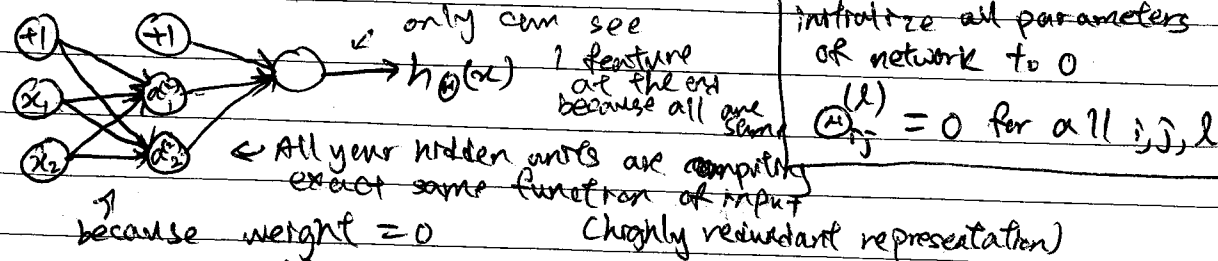
Learning Putting it together

No.

Date

Consider training following neural network

Zero initialization



$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$$

Problem is
Problem of symmetry
the ways are being the same

after 1 gradient descent update, will still same as each other

$$\theta_{01}^{(1)} = \theta_{02}^{(1)}$$

After each update, parameters corresponding to inputs going into each of two hidden units are identical

To get around this problem, the way we initialize parameter of a neural network is with random initialization.

Random initialization: Symmetry breaking

Initialize each $\theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$

If dimension of θ is 1x11, then θ is

E.g. (octave)

random 10x11 matrix (between 0 and 1)

End up $[-\epsilon, \epsilon]$

$$\theta_{01} = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

$$\theta_{02} = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

$\text{rand}(x, y)$ is a function that returns random real numbers between 0 and 1.

Consider this procedure for initializing parameters of a neural network.

1. Pick a random number $r = \text{rand}(1, 1) * (2 * \text{INIT_EPSILON})$

2. Set $\theta_{ij}^{(l)} = r$ for all j, l . $-\text{INIT_EPSILON};$

Does this work?

Answer: No, this fails to break symmetry.

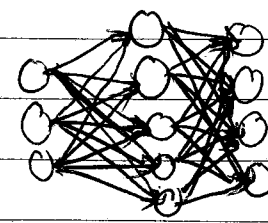
Summary: To create a neural network, have to randomly initialize the weights to small values close to 0, then implement BP, do gradient checking and use gradient descent / advanced optimization algorithm to minimize $J(\theta)$

Overall summary

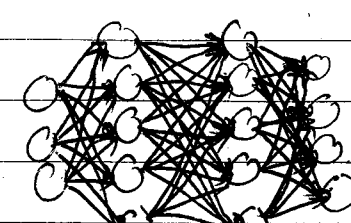
First -

Training a neural network

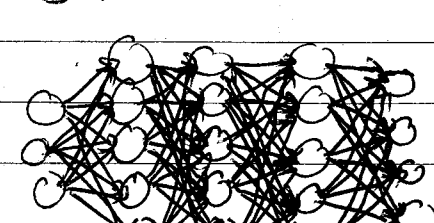
Pick a network architecture (Connectivity pattern between neurons)



3 5 4



3 5 5 4



3 5 5 5 4

Choices of how many hidden layer and how many hidden units in each layer, are architecture choices.

No. of input units: Dimension of features $x^{(1)}$

No. of output units: Number of classes

Reasonable default: 1 hidden layer, or if > 1 hidden layer, have same # of hidden units in every layer (usually the more the better)

Let say multi-classes \Rightarrow

$$y \in \{1, 2, 3, \dots, 10\}$$

if $y = 5$

don't forget to rewrite output y as these vectors.

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

show to neural network

Second - Training a neural network

1. Randomly initialize weights (small values near zero)

2. Implement forward propagation to get $h_{\theta}(x^{(i)})$ for any $x^{(i)}$

3. Implement code to compute cost function $J(\theta)$

4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$

for $i = 1:m$ { After FP and BP on $(x^{(i)}, y^{(i)})$ then FP, BP on $(x^{(2)}, y^{(2)})$ until $(x^{(m)}, y^{(m)})$

Perform FP and BP using example $(x^{(i)}, y^{(i)})$

(Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$)

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T \text{ (accumulation terms)}$$

$$\text{Compute } \frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$$

Thurs

Training a neural network

- 5- Use gradient checking to compare $\frac{\partial}{\partial \theta_{ijk}} J(\theta)$ computed using BP vs. using numerical estimate of gradient of $J(\theta)$. Then disable gradient checking code.
- 6- Use gradient descent or advanced optimization method with BP to try to minimize $J(\theta)$ as a function of parameters θ .
- $\frac{\partial}{\partial \theta_{ijk}} J(\theta)$

For neural network,

 $J(\theta)$ — non-convex

Gradient descent very good in minimizing $J(\theta)$ and get a very good local minimum, even it doesn't get to the global optimum

Contour plot of $J(\theta)$ using gradient descent on $J(\theta)$ by using 2 parameters — BP is computing direction to the minimum point

Quiz Suppose you are using gradient descent together with BP to try to minimize $J(\theta)$ as a function of θ . Which of the following would be a useful step for ~~very~~ verifying that the learning algorithm is running correctly?

Answer: Plot $J(\theta)$ as a function of the number of iterations and make sure it is decreasing (or at least non-increasing) with every iteration.

Quiz:

Correct statement:

- 1- Suppose we have a correct implementation of BP, training a neural network using g-descent. Suppose we plot $J(\theta)$ as a function of the number of iterations, and find that it is increasing rather than decreasing. One possible cause of this is learning rate α is too large.
- 2- If we are training a neural network using gradient descent, one reasonable "debugging" step is to make sure it is ~~not~~ working is to plot $J(\theta)$ as a function of the number of iterations, and make sure it is decreasing (or at least non-increasing) after each iteration.
- 3- If our neural network overfits the training set, one reasonable step is to take to increase the regularization parameter λ .
- 4- Using gradient checking can help verify if one's implementation of BP is bug-free.