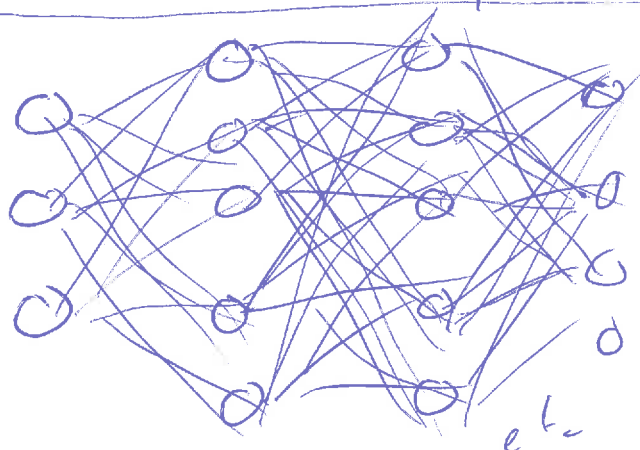


Notes: Week 5

Neural Networks learning: Cost function

Neural Networks \rightarrow for Classification



$$L = 4$$

$$s_1 = 3, s_2 = 5, s_4 = 4$$

Binary classification

$y = 0$ or 1 (binary)

1 output unit

$h_{\theta}(x) \rightarrow \mathbb{R}$ (real #)

$$s_L = 1 \quad / \quad K = 1$$

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

$\rightarrow L =$ total # of layers in network

$s_L =$ no. of units (not counting bias unit) in layer L

Multi-class classification

$$y \in \mathbb{R}^K \text{ e.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

prediction car motor cycle truck

" K " output units

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

$$s_L = K$$

$$(K \geq 3)$$

Cost function:

Logistic regression:

instead of having 1 (logistic regression) unit, we may in fact have K of them. ie

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

Notes: Weeks

Neural Networks Learning: Cost function

Neural Network error function for k output elements.

$h_{\theta}(x) \in \mathbb{R}^K = 1$ (for binary classification)
 $(h_{\theta}(x))_{i=1}^K = i^{\text{th}}$ output

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

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

regularization term.

do not sum over the axes
 where $i=0$

say first hidden layer then

$$\theta_{10}^{(2)} x_0 + \theta_{11}^{(2)} x_1 +$$

\uparrow
 x_0
 (we don't sum over this term).

$$\sum \begin{pmatrix} 0 & k=1 \\ 0 & \vdots \\ 0 & k=K \\ 0 & \end{pmatrix} y_k = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$$

what we want

Notes: Week 5

Back propagation Algorithm

cost function $J(\theta) = \frac{1}{n} \left[\sum_{I=1}^m \sum_{K=1}^K \right]$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$$

Need code to compute:

$$J(\theta)$$

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

parameters are

$$\theta_{ij}^{(l)} \in \mathbb{R}$$

where all parameters are real numbers.

Gradient computation:

Given one training example (x, y) :

lets see what happens.

~~Backward~~ first forward propagation:

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

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

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

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

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

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

$$a^{(4)} = h_0(x) = g(z^{(4)})$$

activation values

	$a^{(1)}$	$a^{(2)}$	$a^{(3)}$	$a^{(4)}$
	1	0	0	0
	0	0	0	0
	0	0	0	0
	0	0	0	0
Layer	2	3	4	

Notes: Week 5

Gradient computation: Backpropagation algorithm

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

$a_j^{(l)}$ = activation of element " j " in layer l

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

$$\delta_j^{(4)} = \underbrace{a_j^{(4)}}_{\substack{\uparrow \\ \text{activation} \\ \text{of the} \\ \text{unit}}} - \underbrace{y_j}_{\substack{\text{the actual} \\ \text{value of the unit} \\ \text{observed in the training} \\ \text{example}}} = (h \theta^{(4)})_j$$

The δ term is just the difference between what our hypothesis $(h \theta^{(4)})_j$ output, and our original value in the training set.

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

Vector dimension \equiv number of output units in our network

derivatives calculated by calculus

$$\begin{aligned} & a^{(3)} * (1 - a^{(3)}) \\ & a^{(2)} * (1 - a^{(2)}) \end{aligned}$$

$$\text{so: } \delta_j^{(4)} = a_j^{(4)} - y_j$$

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

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

no $\delta^{(1)}$!!

ignoring regularization

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

Notes: Week 5

Gradient computation: backpropagation Algorithm

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

set $\Delta_{ij}^l = 0$ (for all l, i, j)

used to compute $\frac{\partial}{\partial \theta_{ij}^l} J(\theta)$

For $i = 1$ to m ^{mth element = 's} $(x^{(i)}, y^{(i)})$

① Set $a^{(1)} = x^{(i)}$

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

③ Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

④ Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ (backprop)

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

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

finally two last pieces outside the for loop:

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

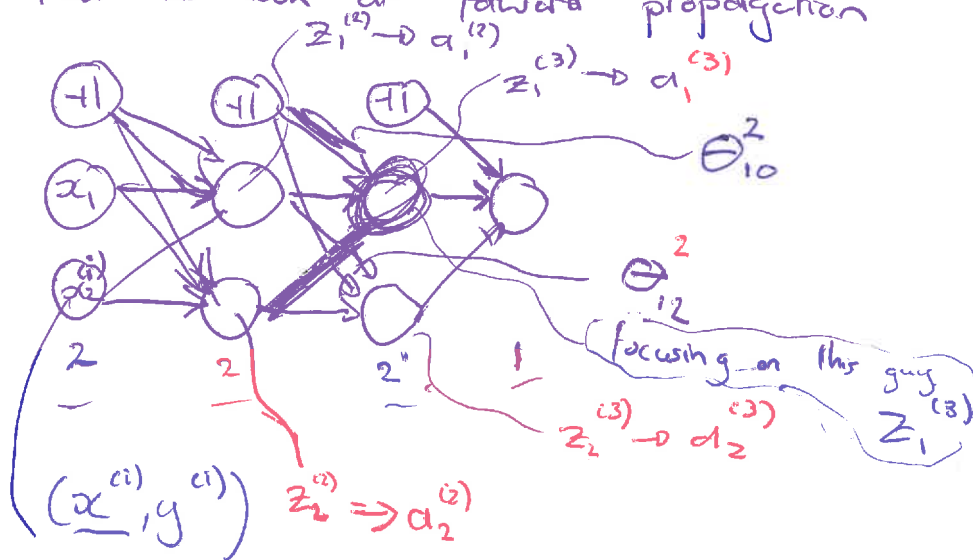
$$D_{ij}^l := \frac{1}{m} \Delta_{ij}^{(l)} \text{ if } j = 0. \text{ (bias term)}$$

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

Notes: Week 5

Backpropagation Intuition

First let's look at forward propagation



forward propagation.

$$z_1^{(3)} = \theta_{10}^{(2)} \times 1 + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)}$$

backpropagation is very similar to FP only that the computations are now flowing back through the network in the opposite direction.

What is backpropagation doing?

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

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

how well is the network doing on example i ?

think of $\text{cost}(i) \approx (h_\theta(x^{(i)}) - y^{(i)})^2$ squared error cost function.

$z_j^{(i)}$ if we can go inside the network and change the terms

$\delta_j^{(i)}$ \rightarrow partial derivative of the cost function with respect to the intermediate terms we are computing.

And so, they are a measure of how much we would like to change the neural networks weights, in order to affect the z -step intermediate values of the (FP) computation, so to affect the overall output $h_\theta(x)$ and \therefore the overall cost.

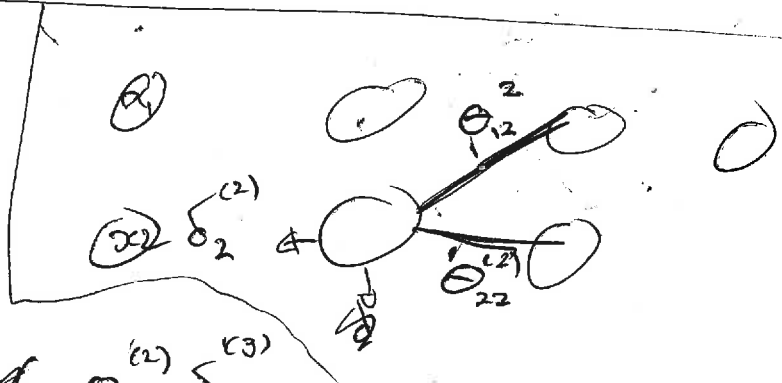
Formally $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ " (for $j \geq 0$) where
 $\text{cost}(i) = y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log h_\theta(x^{(i)})$

so first

$$\delta_1^{(4)} = y^{(i)} - a_1^{(4)}$$

propagate backwards:

$$\delta_2^{(2)} = \Theta_{12}^{(2)} \delta_1^{(3)} + \cancel{1} \Theta_{22}^{(2)} \delta_2^{(3)}$$



Notes - Week 5

Neural Networks: learning: Implementation notes: unrolling parameters

Advanced optimization

function [JVal, gradient] = costFunction(theta)

...
optTheta = fminunc(@costFunction, initialTheta, options)

\mathbb{R}^n (real numbers)
vectors theta
Assumes that both are vectors

taking Neural Network ($L=4$) as an example:

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (theta 1, theta 2, theta 3)

$D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)

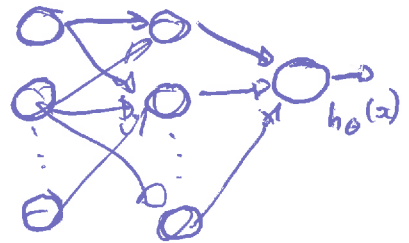
"Unroll" into vectors

example

two hidden units

$S_1 = 10, S_2 = 10, S_3 = 1$ - output unit

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



matlab command:

thetaVec = [Theta2(:); Theta2(:); Theta3(:)];

DVec = [D1(:); D2(:); D3(:)];

Theta1 = reshape(thetaVec(1:110), 10, 11);

Theta2 = reshape(thetaVec(111:220), 10, 11);

Theta3 = reshape(thetaVec(221:231), 1, 11);

elements that you want dimensions of matrix that you are after.

Learning Algorithm

Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

Unroll to get initialTheta to pass to (step 1)

fminunc(@costFunction, initialTheta, options)

Notes: - Week 5

Neural Networks: Learning: Implementation note: unrolling parameters

Learning algorithm

function [Jval, gradientVec] = costFunction(thetaVec)

From thetaVec, get $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ | reshape to matrices

Use forward prop / back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ & $J(\theta)$

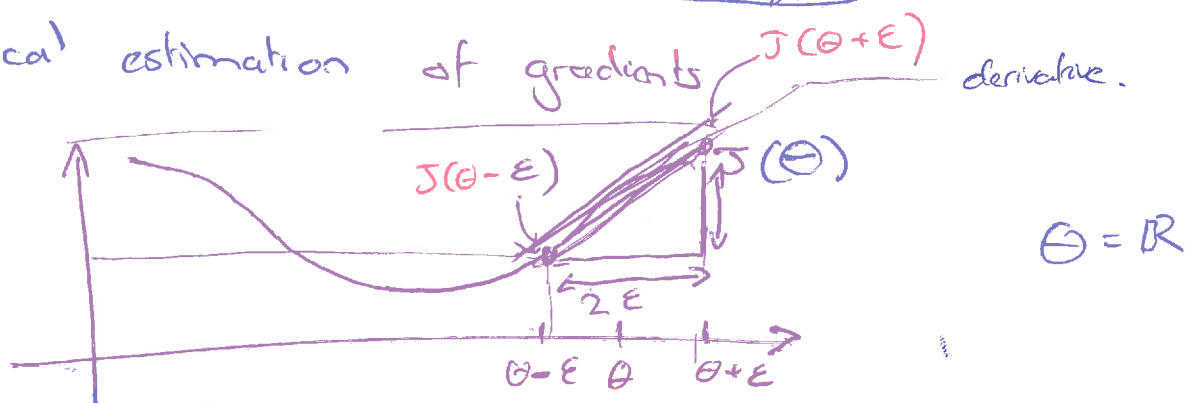
Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec

↗
forward
prop & BP
to compute
derivatives & cost function.

Notes - Week 5

Neural Networks: learning: gradient checking

Numerical estimation of gradients



Slope of the little line is the approximation of the derivative.

$\frac{d}{d\theta} J(\theta) \approx$ the approximation is \therefore given by:

"two sided difference" for estimating the derivative
"one sided" \rightarrow

$$\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

~~$$\frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}$$~~

$\epsilon = 10^{-4}$
with a small ϵ
this becomes very close to the derivative

Implementation: $\boxed{\text{gradApprox}} = \frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 \times \text{EPSILON}}$

Parameter vector θ

$\theta \in \mathbb{R}^n$ (eg θ is "unrolled" version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

$$\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$$

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

$$\vdots$$

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

Oldae code:

... usually with unrolled version of theta,

for $i = 1:n$

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + EPSILON; \rightarrow roughly equal to

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{i-1} + \epsilon \\ \vdots \\ \theta_n \end{bmatrix}$$

thetaMinus = theta;

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

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

[check that Gradapprox is roughly \sim to $\frac{D_{vec}}{\epsilon}$]

\uparrow
from backprop.

Implementation Note:

- Implement backprop to compute D_{vec} (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$)
- Implement numerical gradient check to compute gradApprox.
- Make sure they have similar values.
- Turn off gradient checking - using backprop code for learning.

Important

- disable gradient checking code before training classifier
↳ otherwise code will be very slow.

Notes - Week 5

Neural Networks learning: Random initialization

Initial value of Θ

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

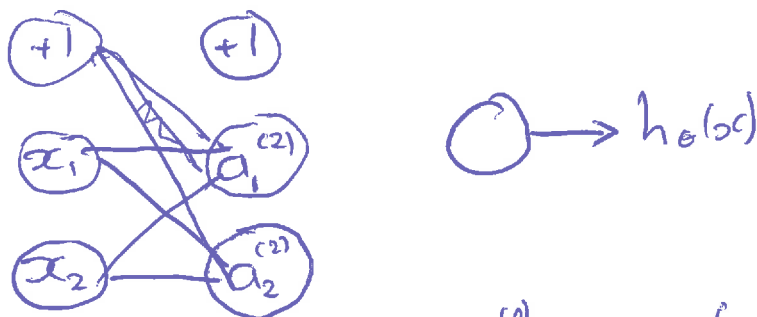
$\text{OptTheta} = \text{fminunc}(@\text{costFunction}, \text{initialTheta}, \text{options})$

Consider gradient descent:

Set $\text{initialTheta} = \text{zeros}(n, 1)$?

↑ \rightarrow this does not work when training a neural network.

Zero initialization



$$a_1^{(2)} = a_2^{(2)} \text{ if } \theta_{ij}^{(2)} = 0 \text{ for all } i, j, l. \quad \left\{ \begin{array}{l} \text{Also } \delta_1^{(2)} = \delta_2^{(2)} \end{array} \right.$$

$$\frac{\partial}{\partial \theta_{o_1}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{o_2}^{(1)}} J(\theta) \quad \left| \begin{array}{l} \text{even after gradient descent} \\ \theta_{o_1}^{(1)} = \theta_{o_2}^{(1)} \end{array} \right.$$

After each update, parameters corresponding to inputs going into each of two hidden units will be identical.

Notes - Week 5

Neural Networks Learning: Random initialization: Symmetry breaking

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

→ random 10x11 matrix (between 0 and 1) different ϵ
then gradient checking

E.g. / $\text{Theta1} = \boxed{\text{rand}(10, 11)} * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$ → $[-\epsilon, \epsilon]$

$$\text{Theta2} = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

Notes - Week 5:

Neural Networks Learning: Putting it All Together

Training a neural network?

- Pick a neural network architecture (connectivity pattern between the neurons)
- No. of input units: Dimension of features $x^{(i)}$ ✓
- No. of output units: number of classes ✓

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

$$y = 5$$

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

- Reasonable default: 1 hidden layer or if > 1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

Training a neural network

- 1) Randomly initialize weights
- 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 $j = 1:m$

perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$

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

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

compute $\frac{\partial}{\partial \theta_{jk}^l} J(\theta)$:

- 5) Use gradient checking to compare $\frac{\partial}{\partial \theta_{jk}^l} J(\theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\theta)$.

then disable gradient checking code ✓

- 6) Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\theta)$ as a function of parameters θ . $J(\theta)$ - non convex \therefore susceptible to local minima.

"backpropagation" predicts the direction of the gradient as gradient descent moves towards a local optimum.