

-ws- Backpropagation Algorithm-

★ Aslında bunun için değil. Bunu sadece Gradient öğrenmesi için kullanacağız.

➤ An algorithm for minimizing the Cost Function
➔ Geçen dersle $J(\Theta)$ formunu yadık.

• Yapmamız gerekeni Find Parameters Θ minimizes $J(\Theta)$

• GDA veya Adv. Opt. Alg.s kullanabilmek için aşağıdaki iki kavramı hesaplayabiliyor olmalıyız:

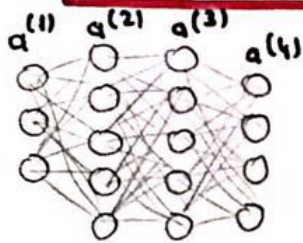
• $J(\Theta)$ ✓ Bunu biliyoruz.

• $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ where $\Theta_{ij}^{(l)} \in \mathbb{R}$
Bunu nasıl hesaplarız ona bakalım:

Let's assume we have only one training ex: (x, y) :

• Let's look at the sequence of calculations:

• **First: Forward Propagation** ile inputa konşilik gelen $h(x)$ 'i bulalım:



Forward Propagation

$$\begin{aligned} \rightarrow a^{(1)} &= x \rightarrow z^{(2)} = \Theta^{(1)} \cdot a^{(1)} \rightarrow a^{(2)} = g(z^{(2)}) \text{ (add } a_0^{(2)}) \\ \rightarrow z^{(3)} &= \Theta^{(2)} \cdot a^{(2)} \rightarrow a^{(3)} = g(z^{(3)}) \text{ (add } a_0^{(3)}) \\ \rightarrow z^{(4)} &= \Theta^{(3)} \cdot a^{(3)} \rightarrow a^{(4)} = h_{\Theta}(x) = g(z^{(4)}) \end{aligned}$$

Remember: $a^{(2)} = \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \\ a_4^{(2)} \\ a_5^{(2)} \end{bmatrix}, \dots$

• **Next: Backpropagation Algorithm** in order to compute the derivative terms or Gradients:

★ Yeni bir kavram tanımlıyoruz " $\delta_J^{(l)}$ " her node için tanımlı. 1. layerdaki J. unit'in "Activation Error'u" ($a_J^{(1)}$ 1. layer J. unit'in activation'u)

For each output unit (layer 4): $\delta_J^{(4)} = a_J^{(4)} - y_J = (h(x))_J - y_J$

vectorized form

$$\delta^{(4)} = a^{(4)} - y$$

Then for other layers:

$$\delta^{(3)} = (\Theta^{(3)})^T \cdot \delta^{(4)} \cdot g'(z^{(3)})$$

where

$$g'(z^{(3)}) = a^{(3)} \cdot (1 - a^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \cdot \delta^{(3)} \cdot g'(z^{(2)})$$

where

$$g'(z^{(2)}) = a^{(2)} \cdot (1 - a^{(2)})$$

No $\delta^{(1)}$ term! Input's Error is almost 0.

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \cdot \delta_i^{(l+1)}$$

➔ (ignoring reg. terms on assuming $\lambda = 0$)

WS - Backpropagation Algorithm II

Gözen sayfa da anlatılanlar biraraya getirilip, Backpropagation Algorithm derivative terms hesaplamak için nasıl kullanıldığını bakarsak:

- Training Set: $\{ (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}) \}$
- Set $\Delta_{ij}^{(1)} = 0$ (for all i, j) (Used to compute $\frac{\partial}{\partial \Theta_{ij}^{(1)}} J(\Theta)$)
- For $i = 1:m$

- Set $a^{(1)} = x^{(i)}$
- Perform Forw. Prop. to compute $a^{(l)}$ for $l = 2, 3, \dots, L$
- Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$
- Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ (there is no $\delta^{(1)}$)
- $\Delta_{ij}^{(1)} := \Delta_{ij}^{(1)} + a_j^{(1)} \cdot \delta_i^{(2)}$

end

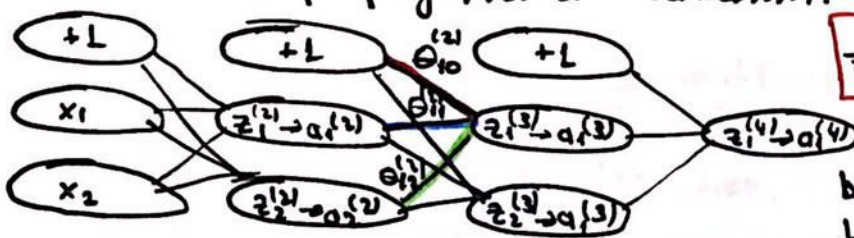
↳ Vectorized Form: $\Delta^{(L)} := \Delta^{(1)} + \delta^{(2+1)} \cdot (a^{(1)})^T$

- $D_{ij}^{(1)} := \frac{1}{m} \Delta_{ij}^{(1)} + 2 \Theta_{ij}^{(1)}$ if $J \neq 0$
- $D_{ij}^{(1)} := \frac{1}{m} \Delta_{ij}^{(1)}$ if $J = 0$

$$\frac{\partial}{\partial \Theta_{ij}^{(1)}} J(\Theta) = D_{ij}^{(1)}$$

- WS - Backpropagation Intuition -

- Backpropagation'ın mekanizmasını daha iyi kavrayabilmek için öncelikle forward-propagation'a bakalım.



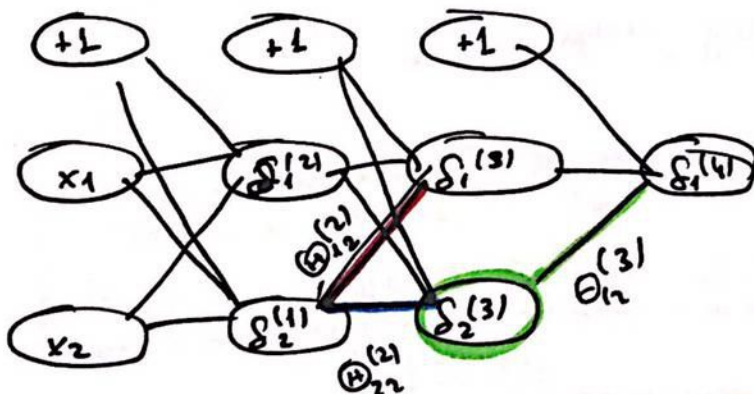
$$z_1^{(3)} = \theta_{10}^{(2)} \cdot 1 + \theta_{11}^{(2)} \cdot a_1^{(2)} + \theta_{12}^{(2)} \cdot a_2^{(2)}$$

Backpropagation da buna benzer çalışıyor fakat computation sağdan sola doğru oluyor.

- $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \text{cost}(i)$ gibi düşünebiliriz. $\text{cost}(i)$ burada single training ex.'in hatası! Normalde logaritmik formda fakat intüition dan $\text{cost}(i) \approx (h(x) - y(i))^2$ gibi düşünebiliriz.

- $\delta_j^{(l)}$ = "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 $j \geq 0$)



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

$$\delta_2^{(3)} = \theta_{12}^{(3)} \cdot \delta_1^{(4)} + \theta_{22}^{(3)} \cdot \delta_2^{(4)}$$

$$\delta_2^{(3)} = \theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

- Bazen bias units için de δ_0 'lar hesaplanır ancak kullanılmaz çünkü anlamsız!

W5 - Implementation Note: Unrolling Parameters from matrices to vectors.

- Bu unrolling olayına **advanced optimization algorithms** için gerek olacak

Advanced Optimization

function [JVal, gradient] = costFunction(theta)
 ↳ \mathbb{R}^{n+1} vector ↳ \mathbb{R}^{n+1} vectors!
 optTheta = fminunc(@costFunction, initialTheta, options)

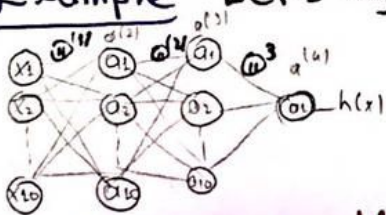
- Yukarıdaki rutinler theta'nın ve gradient'in vector olduğunu assume eden, log. reg. için zaten böyleydi ama NNs için bunlar matris's.

- $L=4$ olan bir NN için

Parameters • $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (Theta1, Theta2, Theta3)
Gradients • $D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)

- Bu videoda amaç NNs için Parameters ve Gradients'i nasıl matris formundan vector forma unroll edebileceğimizi öğrenmek:

Example Let's say we have NN with $s_1=10, s_2=10, s_3=10, s_4=1$



• $\Theta^{(1)} \in \mathbb{R}^{10 \times 10}, \Theta^{(2)} \in \mathbb{R}^{10 \times 10}, \Theta^{(3)} \in \mathbb{R}^{1 \times 10}$ Parameters (Matris)
 • $D^{(1)} \in \mathbb{R}^{10 \times 10}, D^{(2)} \in \mathbb{R}^{10 \times 10}, D^{(3)} \in \mathbb{R}^{1 \times 10}$ Gradients (Matris)

Unrolling

ThetaVec = [Theta1(:); Theta2(:); Theta3(:)];
 DVec = [D1(:); D2(:); D3(:)]

Bu şekilde matrisleri açılıp vucra ekler ve sütun vektörleri elde edilmiş olur

Roll Back

Theta1 = reshape(ThetaVec(1:10), 10, 10);
 Theta2 = reshape(ThetaVec(11:20), 10, 10);
 Theta3 = reshape(ThetaVec(21:21), 1, 10);

When to use Unrolling / Rolling Parameters

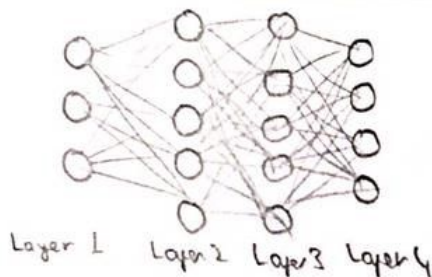
Neden? Çünkü Forward and Backprop. ve $J(\theta)$ calculation is easier by matrices but opt. alg. need vector forms.

- Let's say we have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$
- Unroll them to get "initialTheta" to pass to: fminunc(@costFunction, initialTheta, options)
- function [JVal, gradientVec] = costFunction(ThetaVec)
- From ThetaVec, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ use forward/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\theta)$
- Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec

-WS- NNs Cost Function

- NNs are one of the most powerful Learning Algorithms we have today.
- Bu derslerde NN parametreleri training set'e fit etmeyi göreceğiz
- We gonna focus on the application of NNs to classification problems. Yani başka uygulamalarda da kullanılabiliyor demek ki.

Neural Network (Classification)



- Training Set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- m: # of training examples
- L: # of layers in network ($L=4$)
- s_l: # of units (w/o bias) in layer l ($s_1=3, s_2=4, s_3=5, \dots$)

→ We are gonna consider 2 Types of Classification Problems:

Binary Classification:

- $y = 0$ or 1
- 1 output unit computes $h(x)$
- $h_\theta(x) \in \mathbb{R}$ • $s_L = 1$ • $K = 1$

Multi-class Classification (K classes):

- $y \in \mathbb{R}^K$ E.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ car, $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ truck, $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ bike
- $s_L = K$
- K output units (K classes)

Cost Function for NN: Generalization of the one we used for Log. Reg.

Cost Func. for Log. Reg: $J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \cdot \log(h(x^{(i)})) + (1 - y^{(i)}) \cdot \log(1 - h(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

For NN: $J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \cdot \log(h(x^{(i)})_k) + (1 - y_k^{(i)}) \cdot \log(1 - h(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$

- $h_\theta(x) \in \mathbb{R}^K$ ($K=1$ if binary classification)
- $(h_\theta(x))_i = i$ th output

⚠ Instead of having 1 log. reg. output unit we have K of them.

⚠ Basically we are summing log. reg. cost function over each of my four outputs. Her unit için y_k ile $h(x)_k$ karşılaştırılıyor.

Adv. Opt. Methods can

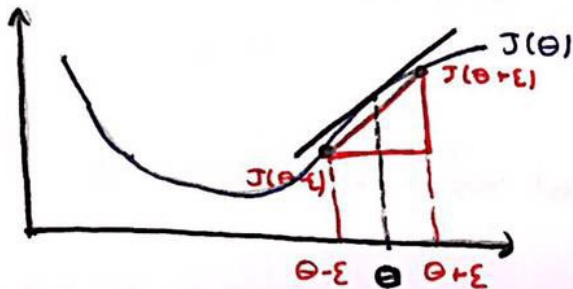
- $J(\theta)$

• Partial Derivative terms $\frac{\partial}{\partial \theta_{ij}^{(l)}}$ for every i, j, l
 → hence # of layers

WS - NNs: Gradient Checking -

- ⊗ Backpropagation is a little tricky. Bazen işe yaramıyş gibi görünür hatta $J(\theta)$ giderek azalar ama sonuçlar $J(\theta)$ 'yi minimize edemediğinden durabilir (backprop alg. dahî buglardan dolayı)
- ⊗ Gradient Checking ile Forw/Back prop. implementasyonımızın 100% doğru olduğundan emin olabiliriz.

- Numerical Estimation of Gradients -



• Assume I have a cost function $\theta \in \mathbb{R}$.

⊗ Amacım numeric olarak θ noktasında $J(\theta)$ 'nin eğimini bulabilecek bir yaklaşım elde etmek.

→ $\epsilon \approx 10^{-4}$, zaten ϵ küçüldükçe \approx giderek $=$ e yaklaşır. Türevin tanımı zaten bu.

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

Implement Code:

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

- General Case when θ is a vector Parameter -

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

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

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

⋮

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

WS - NNs: Gradient Checking II -

- Biraz ervel açılacak general case için aşağıdaki kod kullanılır:
for $i = 1:n$,

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + ϵ ; \rightarrow

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - ϵ ;

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

end

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

$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \text{thetaPlus bu oldu}$

- We implement this for loop to check that $\boxed{\text{gradApprox} \approx \text{DVec}}$
 \nearrow from for loop
 \nwarrow from backpropagation
- Backpropagation ile Gradyent's hesaplamak relatively efficient.
Numerik sonuç ile bunu karşılaştırarak backprop.'un doğru
çalıştığından emin olabiliriz.

- Putting Everything Together: Implementation -

- Implement backpropagation to compute $\boxed{\text{DVec}}$ (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$)
- Implement numerical gradient checking to compute $\boxed{\text{gradApprox}}$.
- Make sure $\boxed{\text{gradApprox} \approx \text{DVec}}$

⚠ Turn off gradient checking before seriously training the network.
Because numerical gradient checking code is computationally expensive.
Yani anladığımız kadarıyla sadece 1 tur için yani 1 theta ile bir
numeric ve backprop ile elde edilen Gradyent karşılaştırılır
doğruluğundan emin olduktan sonra backprop. kullanılmaya
devam edilir.

Important: Be sure to disable your gradient checking code
before training the classifier. If you run gradient computation
on every iteration of GDA (or in the inner loop of cost function),
your code will be very slow.

-WS - Random Initialization -

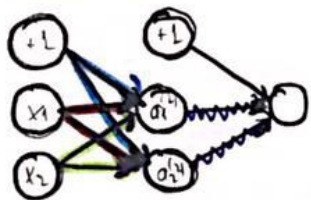
Initial Value of Θ : For GDA and Advanced Optimization Algorithms we need initial value for Θ .

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

Consider gradient descent algorithm,

Setting $\text{initialTheta} = \text{zeros}(n, 1)$ has worked for logistic regression but it does not work we are training a Neural Network!

Zero Initialization Example: Aşağıdaki NN'i eğitme çalışması düzelir.



$\Theta_{ij}^{(1)} = 0$ for all i, j , Θ olarak initialize ederse;

→ Tüm training examples için $a_1^{(2)} = a_2^{(2)}$
→ Ayrıca aynı a'dan çıkan $\Theta^{(2)}$ 'nin weight'lerinde yarı $\Theta_1^{(2)}$ ve $\Theta_2^{(2)}$ de 0 olduğundan, $f_1^{(2)}$ ve $f_2^{(2)}$ aynı!

→ Ayrıca $\frac{\partial}{\partial \Theta_{10}^{(1)}} J(\Theta) = \frac{\partial}{\partial \Theta_{20}^{(1)}} J(\Theta)$ bu da şu demek GDA uygularken sonra bile $\Theta_{20}^{(1)} = \Theta_{10}^{(1)}$ olacak. Aynı renkli weights eşit olacak!

Sonuçta bu da şu demek sonuçta $a_1^{(2)} = a_2^{(2)}$ olacak! Redundant representation. 50 tane hidden unit olsa dahi hepsi aynı sonuç verecek.

Solution: Random Initialization: Symmetry Breaking

Initialize $\Theta_{ij}^{(1)}$ to a random value in $[-A, A]$ (i.e. $-A \leq \Theta_{ij}^{(1)} \leq A$), random 10x11 matrix (between 0-1)

$$\Theta_{10}^{(1)} = \text{rand}(10, 11) * (2 * A) - A;$$

$$\Theta_{20}^{(1)} = \text{rand}(1, 11) * (2 * A) - A;$$

As usually small around zero

★ Peki zero initialization neden log. reg. sorun değildi? Çünkü cost function farklı olduğu için grad'da farklıydı ve her grad için farklı oluyordu bu yüzden farklı sonuçlar geliyordu ama NN için grads backprop ile Θ ile hesaplanır ve aynı geliyor.

WS - Pulling All Together -

How do we implement learning algorithm for NNs.

When training a NN first thing we should do is pick a network architecture:

- # of input units: Dimension of features $x^{(i)}$
- # of output units: Number of classes

Reminder: Don't forget to write y in vector form like $y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ 6 class vantage 3 dimensional

→ Reasonable default: Using a single hidden layer, or if ≥ 2 hidden layer have same # of hidden units in every layer (usually the more the better) → more computationally expensive but rather than that it's a good thing.

Ex: 3-5-4, 3-5-5-4, 3-5-5-5-4.

Training a NN:

- 1) Randomly initialize weights (small values near zero)
- 2) Implement forward propagation to get $h(x^{(i)})$ for any $x^{(i)}$
- 3) Implement code to compute $J(\theta)$
- 4) Implement backprop to compute partial derivatives $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

→ To implement backpropagation we usually use a for loop advanced method uses once, sum, query

for $i = 1:m$ {

Perform forward/back prop. using example $(x^{(i)}, y^{(i)})$

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

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

}

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

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

6) Using gradient descent algorithm or advanced optimization method with backpropagation to try to minimize $J(\theta)$ as a function of parameters $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

③.

⚠ $J(\theta)$ is NON-CONVEX for Neural Networks so GDA or other Adv. Opt. Alg. can get stuck in a local optima. But in practice this is not usually a huge problem. You generally get a global optima but maybe by just trying.