

# Machine Learning

---

## Machine Learning

### Week1

What is Machine Learning?

Parameter Learning

1. Algorithm to calculate  $J_{min}$  : Gradient descent(梯度下降)
2. Gradient Descent Intuition
3. Linear Algebra

### Week 2

1. Multivariate Linear Regression

- (1) Multiple Features
  - (2) Feature Scaling
  - (3) Learning Rate  $\alpha$
  - (4) Features and Polynomial Regression
2. Computing Parameters Analytically
- (1) Normal Equation
  - (2) Normal Equation Noninvertibility
3. Octave/Matlab Tutorial
- (1) Basic operations
  - (2) Moving Data Around
  - (3) Computing on Data
  - (4) Plotting Data
  - (5) Control Statement : `for` `while` `if`
  - (6) Vectorization
  - (7) Method learnt in the assignment

### Week 3

1. Classification -> logistic regression
2. Hypothesis Representation (function used)
  - (1) Logistic Regression Model
  - (2) Interpretation of Hypothesis Output
  - (3) Decision Boundary
3. Logistic Regression Model
  - (1) Cost function
  - (2) Simplified cost function and gradient descent
  - (3) Advanced optimization
4. Multi-class classification: one-vs-all
5. Solving the problem of overfitting
  - (1) The problem of overfitting
  - (2) Cost function
  - (3) regularized linear regression
  - (3) regularized logistic regression

### Week4

1. Motivations
2. Neural Networks
  - (1) Model Representation I
  - (2) Model representation II
3. Applications
  - (1) Example and Intuitions I
  - (2) Example and Intuitions II

### Week 5

1. Cost Function and Backpropagation
  - (1) Cost function
  - (2) Backpropagation Algorithm
  - (3) Backpropagation Intuition

## Week1

### What is Machine Learning?

- Field of study that gives computer the ability to learn without explicitly programmed
- A computer program learns from experience E with respect to some task T and some performance measure P, if its Performance of T, measured by P improves with E.
  - T: play check
  - P: probability to win
  - E: play thousands of games
- Machine learning algorithms
  - Supervised learning : teach computer how to do something
    - give the computer data (right answer) use the data to do something
    - regression problem: fit a line to predict
    - Classification: discrete valued output(0 or 1)
  - Unsupervised learning : Data has no label (clustering algorithm)
  - Reinforcement learning
  - recommender systems

### Model and Cost Function

#### 1. Model representation

- Some definitions
  - Data set -> Training set
  - m = training set . size()
  - x's "input" variables
  - y's "output" variables
  - (x,y) -> a training example
  - (x(i),y(i)) ->  $i_{th}$  training example
- Process
  - training set -> Learning algorithm -> h (function as output of learning algorithms) hypothesis
  - x -> h -> y
  -

$$h_\theta(x) = \theta_0 + \theta_1 x \quad (1)$$

- $\theta$  is parameter and the function represents a line: named as **linear regression** with one variable -> **Univariate(one var) linear regression**

## 2. Method

- Idea : choose  $\theta_0, \theta_1$  that  $h_\theta(x)$  is close to  $y$  for our training examples  $(x, y)$

$$\theta_0, \theta_1 - > \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2 \quad (2)$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2 \quad (3)$$

- $J$  is the cost function with input  $\theta$ , we need find  $J_{min}$ 
  - Suppose  $\theta_0$  is 0, then draw the graph of  $J \rightarrow$  a bow shape graph
  - for two parameters we use contour plots(figures)
    - X-axis :  $\theta_0$  Y-axis :  $\theta_1$
    - Every circle represent same value of  $J(\theta_0, \theta_1)$  (等高线, 3-D 图的投影)

## Parameter Learning

### 1. Algorithm to calculate $J_{min}$ : Gradient descent(梯度下降)

- Outline
  - Start with some  $\theta_0, \theta_1$  (say the two parameters are both zero)
  - Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until end up at a minimum

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1) \quad (for j = 0 and j = 1) \quad (4)$$

}

:= assignment

$\alpha$  is learning rate : size of steps down the hill

update  $\theta_0, \theta_1$  at the same time (simultaneous update)

```

1 | temp0 = theta0 - something(theta0, theta1);
2 | temp1 = theta1 - something(theta0, theta1);
3 | // update by old value first
4 | theta0 = temp0;
5 | theta1 = temp1;
```

## 2. Gradient Descent Intuition

- first suppose  $\theta_0 = 0$ 
  - Use  $\alpha \frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1)$  if delta  $< (>) 0$  then go right(left), when =0 stop
  - $\alpha$ : small->slow, large->fail to converge
- Used in Linear Regression

$$\frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1) = \frac{\delta}{\delta \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2 = \frac{\delta}{\delta \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^i - y^i)^2 \quad (5)$$

$$j = 0 : \frac{\delta}{\delta\theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2 \quad (6)$$

$$j = 1 : \frac{\delta}{\delta\theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2 * x^i \quad (7)$$

- $J$  is Convex function
- "Batch" Gradient Descent: Batch is each step of gradient descent uses **all** the training examples

### 3. Linear Algebra

- Matrix in Matlab/Octave

```

1 % The ; denotes we are going back to a new row.
2 A = [1, 2, 3; 4, 5, 6; 7, 8, 9; 10, 11, 12]
3
4 % Initialize a vector
5 v = [1;2;3]
6
7 % Get the dimension of the matrix A where m = rows and n = columns
8 [m,n] = size(A)
9
10 % You could also store it this way
11 dim_A = size(A)
12
13 % Get the dimension of the vector v
14 dim_v = size(v)
15
16 % Now let's index into the 2nd row 3rd column of matrix A
17 A_23 = A(2,3)
18 % Initialize matrix A and B
19 A = [1, 2, 4; 5, 3, 2]
20 B = [1, 3, 4; 1, 1, 1]
21
22 % Initialize constant s
23 s = 2
24
25 % See how element-wise addition works
26 add_AB = A + B
27
28 % See how element-wise subtraction works
29 sub_AB = A - B
30
31 % See how scalar multiplication works
32 mult_As = A * s
33
34 % Divide A by s
35 div_As = A / s
36
37 % What happens if we have a Matrix + scalar?
38 add_As = A + s
39
40 % The ; denotes we are going back to a new row.
41 A = [1, 2, 3; 4, 5, 6; 7, 8, 9; 10, 11, 12]
42

```

```

43 % Initialize a vector
44 v = [1;2;3]
45
46 % Get the dimension of the matrix A where m = rows and n = columns
47 [m,n] = size(A)
48
49 % You could also store it this way
50 dim_A = size(A)
51
52 % Get the dimension of the vector v
53 dim_v = size(v)
54
55 % Now let's index into the 2nd row 3rd column of matrix A
56 A_23 = A(2,3)
57
58 B=A*v
59 % Initialize a 3 by 2 matrix
60 A = [1, 2; 3, 4;5, 6]
61
62 % Initialize a 2 by 1 matrix
63 B = [1; 2]
64
65 % We expect a resulting matrix of (3 by 2)*(2 by 1) = (3 by 1)
66 mult_AB = A*B
67
68 % Make sure you understand why we got that result
69
70 % Initialize random matrices A and B
71 A = [1,2;4,5]
72 B = [1,1;0,2]
73
74 % Initialize a 2 by 2 identity matrix
75 I = eye(2)
76
77 % The above notation is the same as I = [1,0;0,1]
78
79 % What happens when we multiply I*A ?
80 IA = I*A
81
82 % How about A*I ?
83 AI = A*I
84
85 % Compute A*B
86 AB = A*B
87
88 % Is it equal to B*A?
89 BA = B*A
90
91 % Note that IA = AI but AB != BA
92
93 % transpose and invert
94 A_trans = A'
95 A_invA = inv(A)

```

- for  $\vec{x} = [x_1; x_2; \dots; x_n]$ , we use matrix to find  $\theta_0, \theta_1$ 
  - m represents number of sample

$$\begin{bmatrix} 1 & w_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad (8)$$

## Week 2

---

### 1. Multivariate Linear Regression

---

#### (1) Multiple Features

- Features  $\rightarrow \{x_1, x_2, \dots, x_n\} \rightarrow y \rightarrow$  Price
- Notation:
  - $n$  = number of features
  - $x^{(i)}$  = input (features) of  $i^{th}$  training example ( $n$  dimensional vector) : 第*i*组sample
  - $x_j^{(i)}$  = value of feature  $j$  in  $i^{th}$  training example: 第*j*个feature
- Hypothesis:

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T X \quad (9)$$

For convenience of notation, define  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T X \quad (10)$$

$$J(\theta) = J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (11)$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} = \text{ (for } j \in [0, n]) \quad (12)$$


---

#### (2) Feature Scaling

- When the two parameter have a great difference (say size(0-2000  $feet^2$ ), number of bedrooms(1-5)) The contour plots may be **tall and slim** makes it difficult to find the local minimum
- Feature Scaling
  - Use
 
$$x_1 = \frac{\text{size}(0 - 2000 \text{feet}^2)}{2000}, x_2 = \frac{\text{number of bedrooms}}{5} \quad (x_1, x_2 \in [0, 1]) \quad (13)$$
    - Only needed when the scale is slightly different
  - Mean Normalization
    - Replace  $x_i$  with  $x_i - \mu_i$  to make features have approximately zero mean (do not apply to  $x_0 = 1$ )
 
$$x_1 = \frac{\text{size} - 1000}{2000}, x_2 = \frac{\text{no. bedrooms} - 2}{5} \quad (14)$$
    - $x_i = \frac{x_i - \mu_i}{s_i}$  ( $\mu$  is average,  $s$  is range =  $\max - \min$ ) 
$$(15)$$

### (3) Learning Rate $\alpha$

- two point
    - **Debugging**: How to make sure gradient descent is working correctly
    - How to choose  $\alpha$
  - Debugging use a new function
    - Y-axis :  $\min J(\theta)$  X-axis: No. of iterations
    - $\min J(\theta)$  should decrease after every iteration
- 

### (4) Features and Polynomial Regression

- 2 Points
  - choice of features
  - Polynomial Regression
- Defining new features

$$h_{\theta}(x) = \theta_0 + \theta_1 * frontage + \theta_2 * depth \quad (16)$$

$$frontage * depth = Area = x \quad (17)$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad (18)$$

- Polynomial regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \quad (19)$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 \sqrt{x} \quad (20)$$

use cubic function because house price will not go down as size increase

We can simply use

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad (21)$$

$$x_1 = size, x_2 = size^2, x_3 = size^3 \quad (22)$$

Don't forget Feature scaling

## 2. Computing Parameters Analytically

---

### (1) Normal Equation

- Solve for  $\min J(\theta)$  at one step analytically
- Intuition: if 1D ( $\theta \in R$ )

$$J(\theta) = a\theta^2 + b\theta + c \quad (23)$$

$$\frac{\delta}{\delta\theta} J(\theta) = \dots - > 0 \quad (24)$$

But it is too difficult to implement

$$\theta = (X^T X)^{-1} X^T \vec{y} \quad (25)$$

- Explain

$X$  is the training set  $m \times n$  ( $m$  sets of data,  $n$  features)  $\vec{y}$  is the right answer

$$X\vec{\theta} = \vec{y} \text{ is not always solvable} \quad (26)$$

We find a similar answer  $\hat{\theta}$  such that  $X\hat{\theta} = \vec{p}$  (27)

$\vec{p}$  is a projection of  $\vec{y}$  on the column space of  $X$ , to make it solvable (28)

$$\hat{\theta} = X^{-1}\vec{p} \quad (29)$$

$\vec{p}$  and **column space of  $X$  are orthogonal to  $\vec{y} - \vec{p}$**   $\implies X^T \cdot (\vec{y} - \vec{p}) = \vec{0} \implies X^T \cdot (\vec{y} - X\hat{\theta}) = \vec{0}$  (30)

$$\hat{\theta} = (X^T X)^{-1} X^T \vec{y} \quad (31)$$

$$\vec{p} = X\hat{\theta} = X(X^T X)^{-1} X^T \vec{y} = P\vec{y} \quad (P \text{ is the projection matrix}) \quad (32)$$

$$P = X(X^T X)^{-1} X^T \quad (33)$$

- o  $P^T = P, P^2 = P$
- o  $P * \text{Anyone} \in \text{The column space of } X$
- Implementation (no need feature scaling)

```
1 | pinv(x'*x)*x'*y
```

- Pros
  - o No need  $\alpha$
  - o No need iterate
- Cons : if n is large(features), it is slow [ $O(n^3)$  while gradient descent is  $O(kn^2)$  fit for( $n < 10000$ )]

## (2) Normal Equation Noninvertibility

- In octave
  - o `pinv()` : pseudo invert (has value **even not invertable**)
  - o `inv()` : real invert
- $XX^T$  is not invertable
  - o Redundant features (dependent)
  - o To many featrues(e.g.  $m < n$ )
    - | Delete some features, or use **regularization**

## 3. Octave/Matlab Tutorial

### (1) Basic operations

```
1 | 1~=2    %not equal sign
2 | xor(1,0)  %yihuo
3 | 1&&0    %AND
4 | 1 || 0   %or
5 | PS1("somestring"); %change the string in front of command line
6 | x = 3;    %semicolon supressing output(not print the answer)
7 | a=pi
8 | disp(a); %display
```

```

9 disp(sprintf('2 decimals: %0.2f',a)) %display in two decimal place
10 v=1:0.1:2 %quick way to get row vector from 1 to 2 every column increase 0.1
11 v=1:6 %from 1 to 6 increase 1 by default
12 w=ones(2,3) %all entries is 1
13 w=zeros(2,3)
14 w=rand(2,3)
15 w=randn(1,3)%Gaussian distribution
16 w = -6+sqrt(10)*(randn(1,10000));
17 hist(w) %draw histogram
18 hist(w,50) %histogram with 50 bins
19 eye(4) %4by4 identity matrix
20 help eye
21 help rand %show the detail of method
22

```

## (2) Moving Data Around

```

1 sz=size(A) % show size of matrix sz is also a matrix
2 size(A,1) % row of A
3 size(A,2) % column of A
4 length(A) % longer dimention max(m,n)
5 pwd % current path you are in
6 cd % change the pass
7 ls % list the file
8 load featureX.dat %load the file
9 load('featuresX.dat')
10 who % show all the variables
11 whos % show in detail
12 clear featuresX %delete the variable
13 v = priceY(1:10) % store first 10 columns of priceY in the v
14 save hello.mat v %save v in the new file
15 save hello.txt v-ascii %save as txt
16 A(3,2) %i=3,j=2 elements of A
17 A(2,:) %: means the whole row or column
18 A([1 3],:) %all 1 rows and 3 rows
19 A(:,2)=[10;11;12] %assign
20 A=[A,[100;101;102]] %add a column
21 A(:) %put all elements of A into a single vector
22

```

## (3) Computing on Data

```

1 A.*B %get matrix c that c(i,j)=a(i,j)*b(i,j)
2 v = [1;2;3]
3 1./v %get matrix c that c(i,j) = 1/a(i,j)
4 log(v)
5 exp(v)
6 abs(v) % similar
7 % add 1 to every entries
8 v+1
9 v+ones(length(v),1)
10 A' % A transpose
11 A = [1 15 2 0.5]
12 val = max(A) % biggest number
13 max(A,[],1)

```

```

14 max(A,[],2)    % row and col max respectively
15 [val,ind]=max(a)% assign the position to ind
16 A < 3        % return true or false as a matrix
17 find(A<3)    % return all true entrys
18 A = magic(3)  % all the rows and columns are diagonals sum to a same number
19 [r,c]=find(A>=7)% A(ri,ci)>=7
20 sum(A)        % add all column or row for m = 1 matrix
21 sum(A,1)      % sum all the row
22 sum(A,2)      % sum all the col
23 sum(sum(A.*eyes(size(A),1))) % sum the diagonal of one side
24 sum(sum(A.*flipud(eyes(size(A),1)))) % sum the diagonal of other side
25 prod(A)       % multiple
26 floor(A)
27 ceil(A)       %for all entries
28
29

```

## (4) Plotting Data

```

1 plot(t,y1)      % t-x axis y-y axis
2 plot(t,y2,'r')  % print y2 in a different color
3 xlabel('time')
4 ylabel('value')
5 legend('sin','cos')
6 title('my plot')
7 cd 'your path'; print -dpng 'myPlot.png'
8
9 % plot in different pictures
10 figure(1); plot(t,y1);
11 figure(2); plot(t,y2);
12
13 % Divides plot a 1*2 grid, access first element
14 subplot(1,2,1);
15 plot(t,y1);
16 subplot(1,2,2);
17 plot(t,y2);
18 axis([0.5 1 -1 1]) %change range of x and y
19 clf;               % clear figure
20 A = magic(5);
21 imagesc(A);        %represent in different colors
22 % present the color bar and use black white style
23 imagesc(A),colorvar,colormap gray;
24

```

## (5) Control Statement: `for` `while` `if`

```

1 % for loop
2 v = zeros(1,10)
3 for i = 1:10,
4     v(i) = 2*i;
5 end;
6 % or
7 indices = 1:10;
8 for i = indices,
9     disp(i);

```

```

10 end;
11 % while loop
12 i = 1;
13 while i<=5,
14     v(i++) = 100;
15 end;
16 % or
17 i = 1;
18 while true,
19     v(i++) = 99;
20     if i == 6,
21         break;
22     end;
23 end;
24 % if - else
25 if v(1) == 1,
26     disp('The value is one');
27 elseif v(1) == 2,
28     disp('The value is two');
29 else
30     disp('The value is not one or two');
31 end;
32
33

```

- Create Function

- Create a new file named as 'the name of function.m'

```

1 function y = squareThisNumber(x)
2 y = x^2;
3 % use
4 addpath('path of the function')
5 squareThisNumber(5)
6
7 function [y1,y2] = squareAndCubeThisNumber(x)
8 y1 = x^2;
9 y2 = x^3;
10

```

- $J(\theta)$  in the example

```

1 % use the function
2 >> x = [1 1;1 2;1 3]
3 >> y = [1;2;3]
4 >> theta = [0;1];
5

```

```

1 % the function file 'costFunctionJ'
2 function J = costFunctionJ(X,y,theta)
3 m = size(X,1)
4 predictions = X*theta;           % m*2 . 2*1
5 sqrErrors = (predictions-y).^2;   % turn every entries into its square
6
7 J = 1/(2*m)*sum(sqrErrors);
8

```

## (6) Vectorization

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j = \theta^T X \quad (34)$$

- silly version

```

1 prediction = 0.0;
2 for j=1:n+1,
3     prediction = prediction + theta(j)*X(j);
4 end;

```

- Vectorized version

```
1 | prediction = theta'*X;
```

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \cdot x_j^{(i)} \quad (\text{for } j \in [0, n]) \quad (35)$$

- Implementation

```

1 prediction = X*theta;    % (m,1)
2 Error = prediction - Y; % (m,1)
3 delta = (1/m)*X'*Error; % (n+1,m)*(m,1)=(n+1,1)
4 theta=theta-alpha*delta;
5
6 theta = theta-alpha*X'*(X*theta-Y)

```

## (7) Method learnt in the assignment

```

1 data = load('filename'); % open the file and store it as a matrix
2 linspace - 生成线性间距向量
3
4 此 MATLAB 函数 返回包含 x1 和 x2 之间的 100 个等间距点的行向量。
5
6 y = linspace(x1,x2)
7 y = linspace(x1,x2,n)
8
9
10 logspace - 生成对数间距向量
11
12 此 MATLAB 函数 生成一个由在 10^a 和 10^b (10 的 N 次幂) 之间的 50 个对数间距点组成的行向量 y。logspace
13 函数对于创建频率向量特别有用。该函数是 linspace 和 : 运算符的对数等价函数。
14
15 y = logspace(a,b)
16 y = logspace(a,b,n)

```

```
17 |     y = logspace(a,pi)
```

- plot

```
1 | % two axis of 100 grids from (-10,10) (-1,4) respectively
2 | theta0_vals = linspace(-10,10,100);
3 | theta1_vals = linspace(-1,4,100);
4 | % draw the graph
5 | J_vals = zeros(length(theta0_vals),length(theta1_vals));
6 | % Fill
7 | for i = 1:length(theta0_vals)
8 |     for j = 1:length(theta1_vals)
9 |         t = [theta0_vals(i);theta1_vals(j)]; %a pair of theta
10 |        J_vals(i,j) = computeCost(x,y,t);
11 |    end
12 | end
13 | % Surface plot
14 | figure;
15 | surf(theta0_vals, theta1_vals, J_vals)
16 | xlabel('\theta_0'); ylabel('\theta_1');
17 |
18 | % Contour plot
19 | figure;
20 | % Plot J_vals as 15 contours spaced logarithmically between 0.01 and 100
21 | contour(theta0_vals, theta1_vals, J_vals, logspace(-2, 3, 20))
22 | xlabel('\theta_0'); ylabel('\theta_1');
23 | hold on;
24 | plot(theta(1), theta(2), 'rx', 'MarkerSize', 10, 'LineWidth', 2);
```

- Normalization for multiple features

```
1 | % normalization : after normalize you need to store the mean and standard(when predict new
2 | | value we need to change its scale again)
3 | mean(X); % return average of each column
4 | std(X); % return standard deviation of each column
5 | repmat(A,m,n) % treat A as an element create a matrix m by n of A as a entry
6 | % the code
7 | X_norm = X;
8 | mu = zeros(1, size(X, 2));
9 | sigma = zeros(1, size(X, 2));
10 | mu = mean(X)
11 | X_norm = X_norm-repmat(mu,size(X,1),1)
12 | sigma = std(X)
13 | X_norm = X_norm./repmat(sigma,size(X,1),1)
```

- Gradient Descent for multiple features

- cost function

```
1 | function J = costFunctionJ(X,y,theta);
2 | J = 1/(2*size(X,1))*((X*theta-y)'*(X*theta-y));
```

- Gradient iteration

```

1 | function [theta, J_history] = gradientDescentMulti(X, y, theta, alpha, num_iters);
2 | % Create a column vector of J
3 | J_history = zeros(num_iters, 1);
4 | % do gradient descent
5 | theta = theta - (1/length(y))*alpha*X'* (X*theta-y);
6 | % compute new J for debug
7 | J_history(iter) = computeCostMulti(X, y, theta);
8 | end

```

- Normal Equation

```

1 | function [theta] = normalEqn(X, y)
2 | theta = zeros(size(X, 2), 1);
3 | theta = pinv(X'*X)*X'*y;
4 | end

```

## Week 3

### 1. Classification -> logistic regression

- Email: spam/not spam
  - Online Transactions: Fraudulent ( Yes/ No )
  - Tumor: Malignant / Benign
- $y \in \{0, 1\}$  0: "Negative Class"; 1: "Positive Class"
- Why not use linear regression? : It is easily affected by noisy points

### 2. Hypothesis Representation (function used)

#### (1) Logistic Regression Model

Want  $0 \leq h_\theta(x) \leq 1$

$$h_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad (36)$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad (37)$$

$g$  is Sigmoid (Logistic) function

#### (2) Interpretation of Hypothesis Output

$h_\theta(x)$  = estimated probability that  $y=1$  on input  $x = P(y = 1 | x; \theta)$

$h_\theta(x) = 0.7 \Rightarrow$  Tell patient that 70% chance of tumor being malignant

$$P(y = 0 | x; \theta) = 1 - P(y = 1 | x; \theta) \quad (38)$$

### (3) Decision Boundary

- property of the g
  - Suppose predict "y=1" if  $h_\theta(x) \geq 0.5$
  - Predict "y=0" if  $h_\theta(x) < 0.5$
  - $g(z) \geq 0.5$  when  $z \geq 0$  ( $\theta^T x \geq 0$ )
- Decision Boundary

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \quad (39)$$

$$\theta = [-3, 1, 1] \quad (40)$$

Predict "y=1" if  $-3 + x_1 + x_2 \geq 0$   $x_1 + x_2 = 3$  is a line (Decision boundary)

Cut the  $x_1 x_2$  plane into two pieces

- What if the training set is more complex (cannot use a line) : Add polynomial parameter

## 3. Logistic Regression Model

### (1) Cost function

- Problem

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

m examples:  $x \in [x_0, x_1, \dots, x_n]$   $x_0 = 1, y \in \{0, 1\}$

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

How to choose  $\theta$

- Logistic regression cost function

$$Cost(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases} \quad (41)$$

- If  $y = 1$ ,  $h_\theta(x) = 1$  then Cost = 0
- But as  $h_\theta(z) \rightarrow 0$ , Cost  $\rightarrow \infty$
- 如果猜对了， 奖励这个算法， 猜错了则增大cost

### (2) Simplified cost function and gradient descent

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(h_\theta(x^{(i)}), y^{(i)}) \quad (42)$$

$$Cost(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases} \quad (43)$$

Note y = 0 or 1 always

$$Cost(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x)) \quad (44)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))] \quad (45)$$

- Gradient Descent

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta)$$

Repeat{

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

$$\frac{\delta}{\delta \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (46)$$

Which is similar to linear regression

### (3) Advanced optimization

Given  $\theta$  we have code that can compute  $J(\theta)$  and  $\frac{\delta}{\delta \theta_j} J(\theta)$

- Optimization algorithm
  - Conjugate gradient
  - BFGS
  - L-BFGS
  - Advantages
    - No need to manually pick  $\alpha$
    - Often faster than gradient descent
  - Disadvantages
    - More complex
- Example (Use derivative)

$$\theta = [\theta_1, \theta_2]$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\delta}{\delta \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\delta}{\delta \theta_2} J(\theta) = 2(\theta_2 - 5)$$

```

1 function [jVal,gradient] = costFunction(theta)
2 jVal = (theta(1)-5)^2+(theta(2)-5)^2;
3 gradient = zero(2,1);
4 gradient(1) = 2*(theta(1)-5);
5 gradient(2) = 2*(theta(2)-5);

```

```

1 % option is a data structure which can store your options
2 options = optimset('GradObj','on','MaxIter','100');
3 % set gradient objective parameter to on, and set the maximum number of iteration
4 initialTheta = zero(2,1);
5 % @represent a pointer to costFunction
6 [optTheta,functionVal,exitFlag] = fminunc(@costFunction,initialTheta,options);

```

## 4. Multi-class classification: one-vs-all

- Multiclass classification
  - Email folding/tagging: Work, friend, family, hobby
  - Medical: Not ill, cold, flu
  - Weather: Sunny,Cloudy,Rain,Snow
- One-vs-all (one-vs-rest) : use more binary classification

$$h_\theta^{(i)}(x) = P(y = i|x; \theta) \quad (i = 1, 2, 3) \text{ for three cluster} \quad (47)$$

- On a new  $x$ , to make a prediction, pick the class  $i$  that maximizes  $\max h_\theta^{(i)}(x)$

## 5. Solving the problem of overfitting

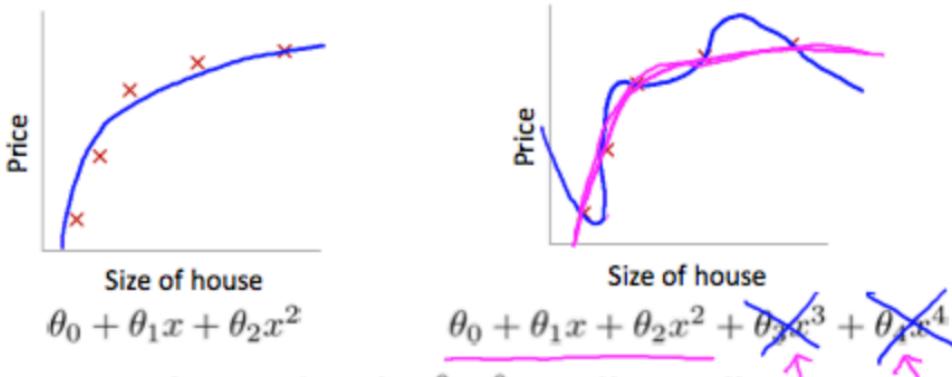
### (1) The problem of overfitting

- Three kinds of fitting
  - Cannot fit very well => "Underfit" or "High bias" (偏见, 偏差)
  - Fit well => "Just right"
  - 100% fitting (Although fitting, but the curve is meaningless) => "Overfit" "High variance"
- Overfitting : If we have too many features, the learned hypothesis may fit the training set very well ( $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$ ), but fail to generalize to new examples (predict prices on new examples)
- Addressing overfitting:
  - draw the graph
  - Options
    - (1) Reduce number of features
      - Manually select which features to keep
      - Model selection algorithm
    - (2) Regularization
      - Keep all the features, but reduce magnitude/values of parameters  $\theta_j$
      - Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

### (2) Cost function

- modify the cost function by add quadratic of overfitting  $\theta$ , and recalculate the minimum, to make them contribute less than before

#### Intuition



Suppose we penalize and make  $\theta_3, \theta_4$  really small.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000 \underline{\theta_3^2} + 1000 \underline{\theta_4^2}$$

$\underline{\theta_3 \approx 0} \quad \underline{\theta_4 \approx 0}$

- Regularization : small value for parameters  $\theta_0, \theta_1, \dots, \theta_n$ 
  - "Simpler" hypothesis
  - Less prone to overfitting
  - In real use, we don't know to shrink which parameter => shrink all parameters ( $\lambda$  is regularization parameter)

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right] \quad (48)$$

- If  $\lambda$  is too large => cause underfitting ( i.e.  $h_\theta(x) = \theta_0$ )

## (2) regularized linear regression

- Gradient descent (with regularization) => treat  $\theta_0$  differently

Repeat{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} [\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \lambda \theta_j] (j \neq 0)$$

Or  $\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} (j \neq 0) ((1 - \alpha \frac{\lambda}{m}) \text{ is always less than } 1)$

}

- Normal equation

$$\theta = (X^T X + \lambda L)^{-1} X^T y \quad (49)$$

Where L is a matrix

$$\begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad (50)$$

- Use when the matrix  $X^T X$  is singular

## (3) regularized logistic regression

- Cost function

$$J(\theta) = -\left[ \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))] \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad (j = 1, 2, \dots, n) \quad (51)$$

- Advanced optimization

# Week4

## 1. Motivations

- There may be a lot of features => it is not efficient to use all feature
- Neural network => build brain

## 2. Neural Networks

## (1) Model Representation I

- Neuron model : logistic unit

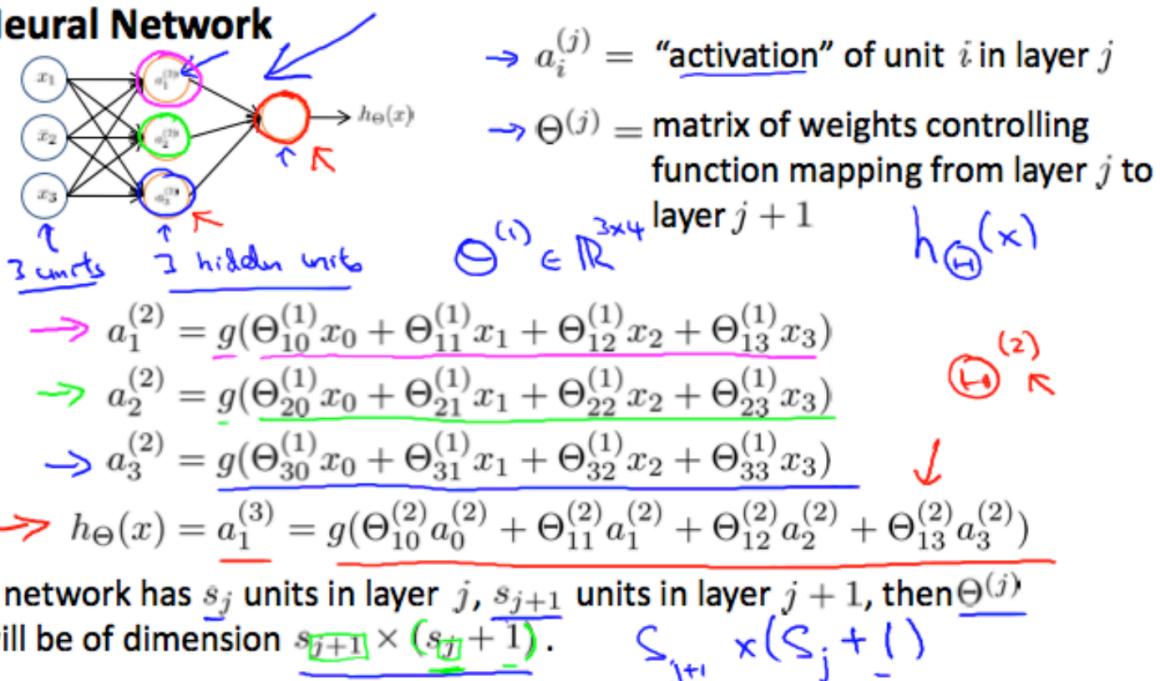
$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix} \Rightarrow [ ] \Rightarrow h_{\theta}(x) \quad (52)$$

- Sigmoid function => activation function  $g(z)$
- $\theta$  => weights or parameters
- We can have intermediate layers of nodes between the input and output layers called the "hidden layers."

$$\begin{aligned} a_1^{(2)} &= g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3) \\ a_2^{(2)} &= g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3) \\ a_3^{(2)} &= g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3) \\ h_{\Theta}(x) = a_1^{(3)} &= g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)}) \end{aligned}$$

- Neuron network

### Neural Network



Andrew N

- dimension of  $\Theta$ :  $s_{j+1} * (s_j + 1)$  where  $s_n$  is the number of variables in nth layer, 1 is the constant column
- $A^{(j)} * \Theta^{(j)T} = A^{(j+1)}$
- More specific description

- Input layer:

- W matrix: Each row of weight represent a neural, column represent features accept
- X matrix : Each row represent a feature, each column represent a sample
- $WX+B = Z$ : Each row of Z represent a neural, each column of z represent a sample(Neural\*feature.feature\*sample), where B is bias
- $A = g(Z)$ , activate output

- Hidden layer:
  - W matrix: Each row represent a neural of this layer, each column represent a neural of last layer, the first column represent the bias of last layer
  - $WA(\text{last})+B = Z$ : Each row of Z represent a neural in this layer, each column of z represent a sample(Neural\*(last Neural+1).(last Neural+1)\*sample)
- Output layer:
  - similar to Hidden layer, instead each row represent an output for this sample

## (2) Model representation II

- Forward propagation: Vectorized implementation

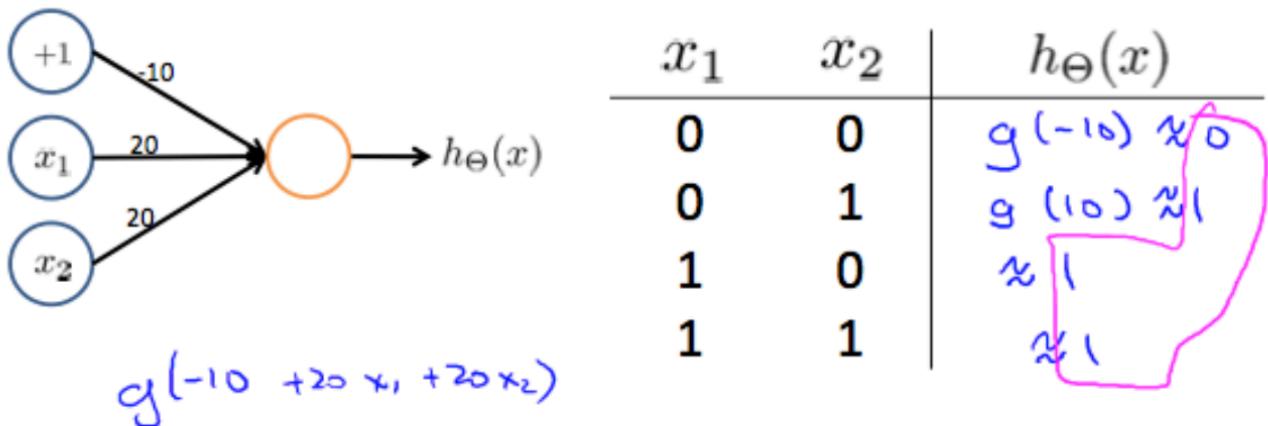
$$z^{(j+1)} = \Theta^{(j)} a^{(j)} \quad (53)$$

$$h_{\theta}(x) = a^{(j+1)} = g(z^{(j+1)}) \quad (54)$$

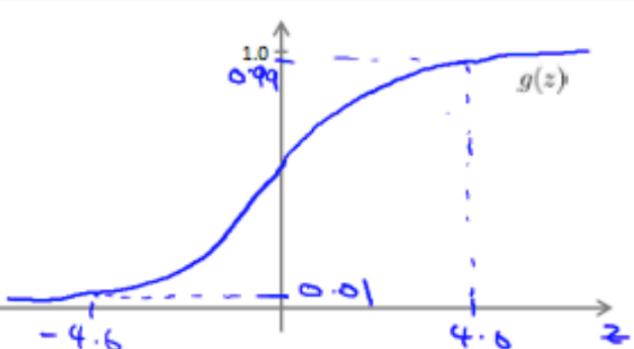
## 3. Applications

### (1) Example and Intuitions I

#### Example: OR function

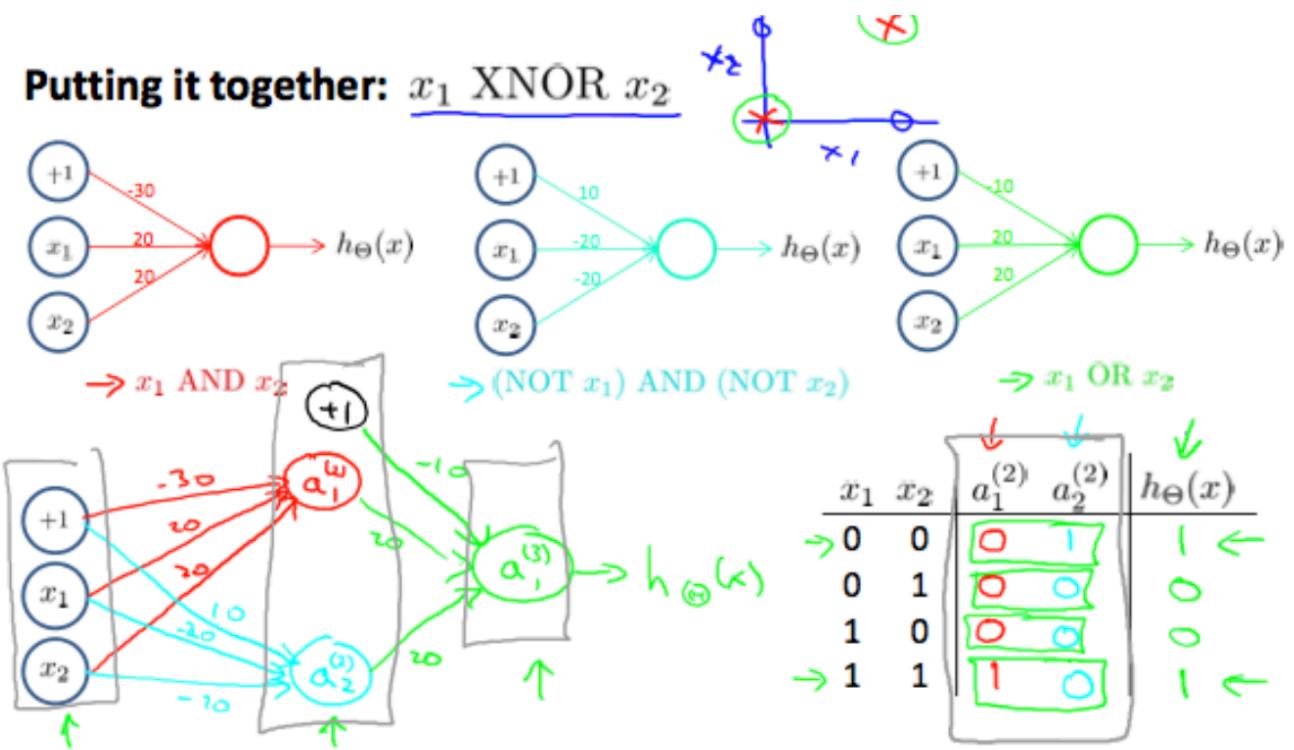


Where  $g(z)$  is the following:



## (2) Example and Intuitions II

**Putting it together:**  $x_1 \text{ XNOR } x_2$



## Week 5

### 1. Cost Function and Backpropagation

#### (1) Cost function

- Cost functions
  - Denotes
    - $L$  = total number of layers in the network
    - $S_l$  = number of units (not counting bias unit) in layer  $l$
    - $K$  = number of output units/classes (channels)
  - Logistic regression

$$J(\theta) = -\left[\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))]\right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad (55)$$

- Neural network (for multi-class problem there are  $K$  clusters)  $h_\Theta(x) \in R^K$  ( $h_\Theta(x)$ ) <sub>$i$</sub>  =  $i^{th}$  output
  - The first part add the cost of all the m-by-K results
  - The second part add all the theta used (every column, row, layer ( $L$ ))
  - regularization doesn't include the  $0^{th}$  weight

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - h_\Theta(x^{(i)}))_k \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\Theta_{ji}^l)^2 \quad (56)$$

## (2) Backpropagation Algorithm

- Proof the backpropagation

$$\text{权重和偏执项 } W^{[3]} = \begin{bmatrix} w_{11}^{[3]} & w_{12}^{[3]} & w_{13}^{[3]} \\ w_{21}^{[3]} & w_{22}^{[3]} & w_{23}^{[3]} \end{bmatrix}_{2 \times 3}, \quad B^{[3]} = \begin{bmatrix} b_1^{[3]} \\ b_2^{[3]} \end{bmatrix}_{2 \times 1}$$

该层的线性计算

$$Z^{[3]} = W^{[3]} A^{[2]} + B^{[3]} = \begin{bmatrix} w_{11}^{[3]} a_1^{[2](1)} + w_{12}^{[3]} a_2^{[2](1)} + w_{13}^{[3]} a_3^{[2](1)} + b_1^{[3]} & w_{11}^{[3]} a_1^{[2](2)} + w_{12}^{[3]} a_2^{[2](2)} + w_{13}^{[3]} a_3^{[2](2)} + b_1^{[3]} \\ w_{21}^{[3]} a_1^{[2](1)} + w_{22}^{[3]} a_2^{[2](1)} + w_{23}^{[3]} a_3^{[2](1)} + b_1^{[3]} & w_{21}^{[3]} a_1^{[2](2)} + w_{22}^{[3]} a_2^{[2](2)} + w_{23}^{[3]} a_3^{[2](2)} + b_1^{[3]} \end{bmatrix}_{2 \times 2}$$

$$\text{记作 } Z^{[3]} = \begin{bmatrix} z_1^{[3](1)} & z_1^{[3](2)} \\ z_2^{[3](1)} & z_2^{[3](2)} \end{bmatrix}_{2 \times 2}$$

$$\text{激活输出 } A^{[3]} = \sigma(Z^{[3]}) = \begin{bmatrix} a_1^{[3](1)} & a_1^{[3](2)} \\ a_2^{[3](1)} & a_2^{[3](2)} \end{bmatrix}_{2 \times 2} = \begin{bmatrix} \sigma(z_1^{[3](1)}) & \sigma(z_1^{[3](2)}) \\ \sigma(z_2^{[3](1)}) & \sigma(z_2^{[3](2)}) \end{bmatrix}_{2 \times 2}$$

$$\text{输出结果是 } Y = \begin{bmatrix} y_1^{(1)} & y_1^{(2)} \\ y_2^{(1)} & y_2^{(2)} \end{bmatrix}_{2 \times 2} = A^{[3]},$$

$$\text{对应的真实标签值记是 } \tilde{Y} = \begin{bmatrix} \tilde{y}_1^{(1)} & \tilde{y}_1^{(2)} \\ \tilde{y}_2^{(1)} & \tilde{y}_2^{(2)} \end{bmatrix}_{2 \times 2}$$

$$\text{设每个神经元的激活函数为最常用的Sigmoid函数: } \sigma(z) = \frac{1}{1 + e^{-z}}$$

- ■ each number in a "()" represent the sample id
-

$$L = -(\tilde{Y} \log(Y) + (1 - \tilde{Y}) \log(1 - Y)) = \begin{bmatrix} -(\tilde{y}_1^{(1)} \log(y_1^{(1)}) + (1 - \tilde{y}_1^{(1)}) \log(1 - y_1^{(1)})) & -(\tilde{y}_1^{(2)} \log(y_1^{(2)}) + (1 - \tilde{y}_1^{(2)}) \log(1 - y_1^{(2)})) \\ -(\tilde{y}_2^{(1)} \log(y_2^{(1)}) + (1 - \tilde{y}_2^{(1)}) \log(1 - y_2^{(1)})) & -(\tilde{y}_2^{(2)} \log(y_2^{(2)}) + (1 - \tilde{y}_2^{(2)}) \log(1 - y_2^{(2)})) \end{bmatrix}_{2 \times 2}$$

简记为：

$$L = \begin{bmatrix} l_1^{(1)} & l_1^{(2)} \\ l_2^{(1)} & l_2^{(2)} \end{bmatrix}$$

则由简单的链导法则可有：

$$\frac{dL}{dZ^{[3]}} = \frac{dL}{dA^{[3]}} \frac{dA^{[3]}}{dZ^{[3]}} = \begin{bmatrix} \frac{dl_1^{(1)}}{da_1^{[3](1)}} \frac{da_1^{[3](1)}}{dz_1^{[3](1)}} & \frac{dl_1^{(2)}}{da_1^{[3](2)}} \frac{da_1^{[3](2)}}{dz_1^{[3](2)}} \\ \frac{dl_2^{(1)}}{da_2^{[3](1)}} \frac{da_2^{[3](1)}}{dz_2^{[3](1)}} & \frac{dl_2^{(2)}}{da_2^{[3](2)}} \frac{da_2^{[3](2)}}{dz_2^{[3](2)}} \end{bmatrix}_{2 \times 2} = \begin{bmatrix} dz_1^{[3](1)} & dz_1^{[3](2)} \\ dz_2^{[3](1)} & dz_2^{[3](2)} \end{bmatrix}_{2 \times 2}$$

记  $\frac{dL}{dZ^{[3]}} = dZ^{[3]}$  ,  $\frac{dL}{dA^{[3]}} = dA^{[3]}$  ,  $\frac{dA}{dZ^{[3]}} = \sigma'(Z^{[3]})$  , 则, 显然有:

$dZ^{[3]} = dA^{[3]} * \sigma'(Z^{[3]})$  , 即:

$$\begin{bmatrix} dz_1^{[3](1)} & dz_1^{[3](2)} \\ dz_2^{[3](1)} & dz_2^{[3](2)} \end{bmatrix}_{2 \times 2} = \begin{bmatrix} da_1^{[3](1)} & da_1^{[3](2)} \\ da_2^{[3](1)} & da_2^{[3](2)} \end{bmatrix}_{2 \times 2} * \begin{bmatrix} d\sigma(dz_1^{[3](1)}) & d\sigma(dz_1^{[3](2)}) \\ d\sigma(dz_2^{[3](1)}) & d\sigma(dz_2^{[3](2)}) \end{bmatrix}_{2 \times 2}$$

其中的 \* 表示逐元素相乘 (比如MatLAB里的A.\*B, Python里的矩阵 A\*B)。这一块涉及的函数求导就不赘述了, 然后就可很容易计算出来下面的结果:

$$dZ^{[3]} = A^{[3]} - \tilde{Y} ,$$

$$\frac{dL}{dW^{[3]}} = \frac{dL}{dZ^{[3]}} \frac{dZ^{[3]}}{dW^{[3]}} = \begin{bmatrix} \frac{dl_1^{(1)}}{dw_{11}^{[3]}} + \frac{dl_1^{(2)}}{dw_{11}^{[3]}} & \frac{dl_1^{(1)}}{dw_{12}^{[3]}} + \frac{dl_1^{(2)}}{dw_{12}^{[3]}} & \frac{dl_1^{(1)}}{dw_{13}^{[3]}} + \frac{dl_1^{(2)}}{dw_{13}^{[3]}} \\ \frac{dl_1^{(1)}}{dw_{21}^{[3]}} + \frac{dl_1^{(2)}}{dw_{21}^{[3]}} & \frac{dl_1^{(1)}}{dw_{22}^{[3]}} + \frac{dl_1^{(2)}}{dw_{22}^{[3]}} & \frac{dl_1^{(1)}}{dw_{23}^{[3]}} + \frac{dl_1^{(2)}}{dw_{23}^{[3]}} \end{bmatrix}_{2 \times 3}$$

$$dW^{[3]} = \begin{bmatrix} \frac{dl_1^{(1)}}{dz_1^{[3](1)}} \frac{dz_1^{[3](1)}}{dw_{11}^{[3]}} + \frac{dl_1^{(2)}}{dz_1^{[3](2)}} \frac{dz_1^{[3](2)}}{dw_{11}^{[3]}} & \frac{dl_1^{(1)}}{dz_1^{[3](1)}} \frac{dz_1^{[3](1)}}{dw_{12}^{[3]}} + \frac{dl_1^{(2)}}{dz_1^{[3](2)}} \frac{dz_1^{[3](2)}}{dw_{12}^{[3]}} & \frac{dl_1^{(1)}}{dz_1^{[3](1)}} \frac{dz_1^{[3](1)}}{dw_{13}^{[3]}} + \frac{dl_1^{(2)}}{dz_1^{[3](2)}} \frac{dz_1^{[3](2)}}{dw_{13}^{[3]}} \\ \frac{dl_1^{(1)}}{dz_2^{[3](1)}} \frac{dz_2^{[3](1)}}{dw_{21}^{[3]}} + \frac{dl_1^{(2)}}{dz_2^{[3](2)}} \frac{dz_2^{[3](2)}}{dw_{21}^{[3]}} & \frac{dl_1^{(1)}}{dz_2^{[3](1)}} \frac{dz_2^{[3](1)}}{dw_{22}^{[3]}} + \frac{dl_1^{(2)}}{dz_2^{[3](2)}} \frac{dz_2^{[3](2)}}{dw_{22}^{[3]}} & \frac{dl_1^{(1)}}{dz_2^{[3](1)}} \frac{dz_2^{[3](1)}}{dw_{23}^{[3]}} + \frac{dl_1^{(2)}}{dz_2^{[3](2)}} \frac{dz_2^{[3](2)}}{dw_{23}^{[3]}} \end{bmatrix}_{2 \times 3}$$

上式计算后，结果为：

$$dW^{[3]} = \begin{bmatrix} dz_1^{[3](1)} a_1^{[2](1)} + dz_1^{[3](2)} a_1^{[2](2)} & dz_1^{[3](1)} a_2^{[2](1)} + dz_1^{[3](2)} a_2^{[2](2)} & dz_1^{[3](1)} a_3^{[2](1)} + dz_1^{[3](2)} a_3^{[2](2)} \\ dz_2^{[3](1)} a_1^{[2](1)} + dz_2^{[3](2)} a_1^{[2](2)} & dz_2^{[3](1)} a_2^{[2](1)} + dz_2^{[3](2)} a_2^{[2](2)} & dz_2^{[3](1)} a_3^{[2](1)} + dz_2^{[3](2)} a_3^{[2](2)} \end{bmatrix}_{2 \times 3}$$

从上面的式子可以看出，每个权重的梯度是每个样本得到的梯度之和，因此，这里都除以样本个数，求出平均梯度。整理一下，我们就得到：

$$dW^{[3]} = \frac{1}{2} \begin{bmatrix} dz_1^{[3](1)} & dz_1^{[3](2)} \\ dz_2^{[3](1)} & dz_2^{[3](2)} \end{bmatrix}_{2 \times 2} \begin{bmatrix} a_1^{[2](1)} & a_2^{[2](1)} & a_3^{[2](1)} \\ a_1^{[2](2)} & a_2^{[2](2)} & a_3^{[2](2)} \end{bmatrix}_{2 \times 3}$$

$$\text{即 } dW^{[3]} = \frac{1}{2} dZ^{[3]} A^{[2]^T}$$

- ■ m = 2 thus use  $\frac{1}{2}$
- 

## 二、隐含层 (Hidden)

$$\text{权重和偏执项 } W^{[2]} = \begin{bmatrix} w_{11}^{[2]} & w_{12}^{[2]} \\ w_{21}^{[2]} & w_{22}^{[2]} \\ w_{31}^{[2]} & w_{32}^{[2]} \end{bmatrix}_{3 \times 2}, B^{[2]} = \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix}_{3 \times 1}$$

该层的线性计算

$$Z^{[2]} = W^{[2]} A^{[1]} + B^{[2]} = \begin{bmatrix} w_{11}^{[2]} a_1^{[1](1)} + w_{12}^{[2]} a_2^{[1](1)} + b_1^{[2]} & w_{11}^{[2]} a_1^{[1](2)} + w_{12}^{[2]} a_2^{[1](2)} + b_1^{[2]} \\ w_{21}^{[2]} a_1^{[1](1)} + w_{22}^{[2]} a_2^{[1](1)} + b_2^{[2]} & w_{21}^{[2]} a_1^{[1](2)} + w_{22}^{[2]} a_2^{[1](2)} + b_2^{[2]} \\ w_{31}^{[2]} a_1^{[1](1)} + w_{32}^{[2]} a_2^{[1](1)} + b_3^{[2]} & w_{31}^{[2]} a_1^{[1](2)} + w_{32}^{[2]} a_2^{[1](2)} + b_3^{[2]} \end{bmatrix}_{3 \times 2}$$

$$\text{记作 } Z^{[2]} = \begin{bmatrix} z_1^{[2](1)} & z_1^{[2](2)} \\ z_2^{[2](1)} & z_2^{[2](2)} \\ z_3^{[2](1)} & z_3^{[2](2)} \end{bmatrix}_{3 \times 2}$$

$$\text{激活输出 } A^{[2]} = \sigma(Z^{[2]}) = \begin{bmatrix} a_1^{[2](1)} & a_1^{[2](2)} \\ a_2^{[2](1)} & a_2^{[2](2)} \\ a_3^{[2](1)} & a_3^{[2](2)} \end{bmatrix}_{3 \times 2} = \begin{bmatrix} \sigma(z_1^{[2](1)}) & \sigma(z_1^{[2](2)}) \\ \sigma(z_2^{[2](1)}) & \sigma(z_2^{[2](2)}) \\ \sigma(z_3^{[2](1)}) & \sigma(z_3^{[2](2)}) \end{bmatrix}_{3 \times 2}$$

◦

现在输出层的都求出来，然后就再往回一层，求隐含层的梯度，因此，中间链导需要经过  $A^{[2]}$  :

$$dA^{[2]} = \begin{bmatrix} dz_1^{[3](1)} \frac{dz_1^{[3](1)}}{da_1^{[2](1)}} + dz_2^{[3](1)} \frac{dz_2^{[3](1)}}{da_1^{[2](1)}} & dz_1^{[3](2)} \frac{dz_1^{[3](2)}}{da_1^{[2](2)}} + dz_2^{[3](2)} \frac{dz_2^{[3](2)}}{da_1^{[2](2)}} \\ dz_2^{[3](1)} \frac{dz_1^{[3](1)}}{da_2^{[2](1)}} + dz_2^{[3](1)} \frac{dz_2^{[3](1)}}{da_2^{[2](1)}} & dz_1^{[3](2)} \frac{dz_1^{[3](2)}}{da_2^{[2](2)}} + dz_2^{[3](2)} \frac{dz_2^{[3](2)}}{da_2^{[2](2)}} \\ dz_1^{[3](1)} \frac{dz_1^{[3](1)}}{da_3^{[2](1)}} + dz_2^{[3](1)} \frac{dz_2^{[3](1)}}{da_3^{[2](1)}} & dz_1^{[3](2)} \frac{dz_1^{[3](2)}}{da_3^{[2](2)}} + dz_2^{[3](2)} \frac{dz_2^{[3](2)}}{da_3^{[2](2)}} \end{bmatrix}_{3 \times 2}$$

$$dA^{[2]} = \begin{bmatrix} dz_1^{[3](1)} w_{11}^{[3]} + dz_2^{[3](1)} w_{21}^{[3]} & dz_1^{[3](2)} w_{11}^{[3]} + dz_2^{[3](2)} w_{21}^{[3]} \\ dz_2^{[3](1)} w_{12}^{[3]} + dz_2^{[3](1)} w_{22}^{[3]} & dz_1^{[3](2)} w_{12}^{[3]} + dz_2^{[3](2)} w_{22}^{[3]} \\ dz_1^{[3](1)} w_{13}^{[3]} + dz_2^{[3](1)} w_{23}^{[3]} & dz_1^{[3](2)} w_{13}^{[3]} + dz_2^{[3](2)} w_{23}^{[3]} \end{bmatrix}_{3 \times 2}$$

即：

$$dA^{[2]} = W^{[3]^T} dA^{[3]}$$

然后就可以推导隐含层的梯度了，过程和上面是一样的，所以直接给出结果：

$$dW^{[2]} = \frac{1}{2} dZ^{[2]} A^{[1]^T}$$

$$dB^{[2]} = \frac{1}{2} \text{sum}(dZ^{[2]}, axis=1)$$

◦

然后更新权重即可：

$$\begin{aligned} W^{[i]} &= W^{[i]} + \eta dW^{[i]} \\ B^{[i]} &= B^{[i]} + \eta dB^{[i]} \\ i &= 1, 2, 3 \end{aligned}$$

最后把上面的推导总结一下，对于第  $l$  层：

*Input* :  $dA^{[l]}, A^{[l-1]}, W^{[l]}, Z^{[l]}$

..... 若该层是输出层，则  $dZ^{[l]} = A^{[l]} - \tilde{Y}$ ，否则  $dZ^{[l]} = dA^{[l]} * \sigma'(Z^{[l]})$

.....  $dW^{[l]} = \frac{1}{N} dZ^{[l]} A^{[l-1]^T}$

.....  $dB^{[l]} = \frac{1}{N} \text{sum}(dZ^{[l]}, axis=1)$

.....  $dA^{[l-1]} = W^{[l]^T} dZ^{[l]}$

*Output* :  $dA^{[l-1]}, dW^{[l]}, dB^{[l]}$

- Gradient computation

- we need:  $\min_{\Theta} J(\Theta)$  and  $\frac{\delta}{\delta \Theta_{ij}^{(l)}} J(\Theta)$
- Given one training example  $(x, y)$
- Intuition:  $\delta_j^{(l)}$  = "errors" of node  $j$  in layer  $l$  =  $\frac{\delta_j^{(l)} J}{\delta Z^{(l)}}$

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad (57)$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot * (a^{(3)} \cdot * (1 - a^{(3)})) \quad (58)$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot * (a^{(2)} \cdot * (1 - a^{(2)})) \quad (59)$$

$$\frac{\delta}{\delta \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)} \quad (60)$$

- Backpropagation algorithm

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

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

For  $i=1$  to  $m$

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

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

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

Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  using  $\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) \cdot * a^{(l)} \cdot * (1 - a^{(l)})$

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

### 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  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

For  $i = 1$  to  $m$  ←  $(x^{(i)}, y^{(i)})$ .

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

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

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

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

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

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

→  $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$  if  $j = 0$

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

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

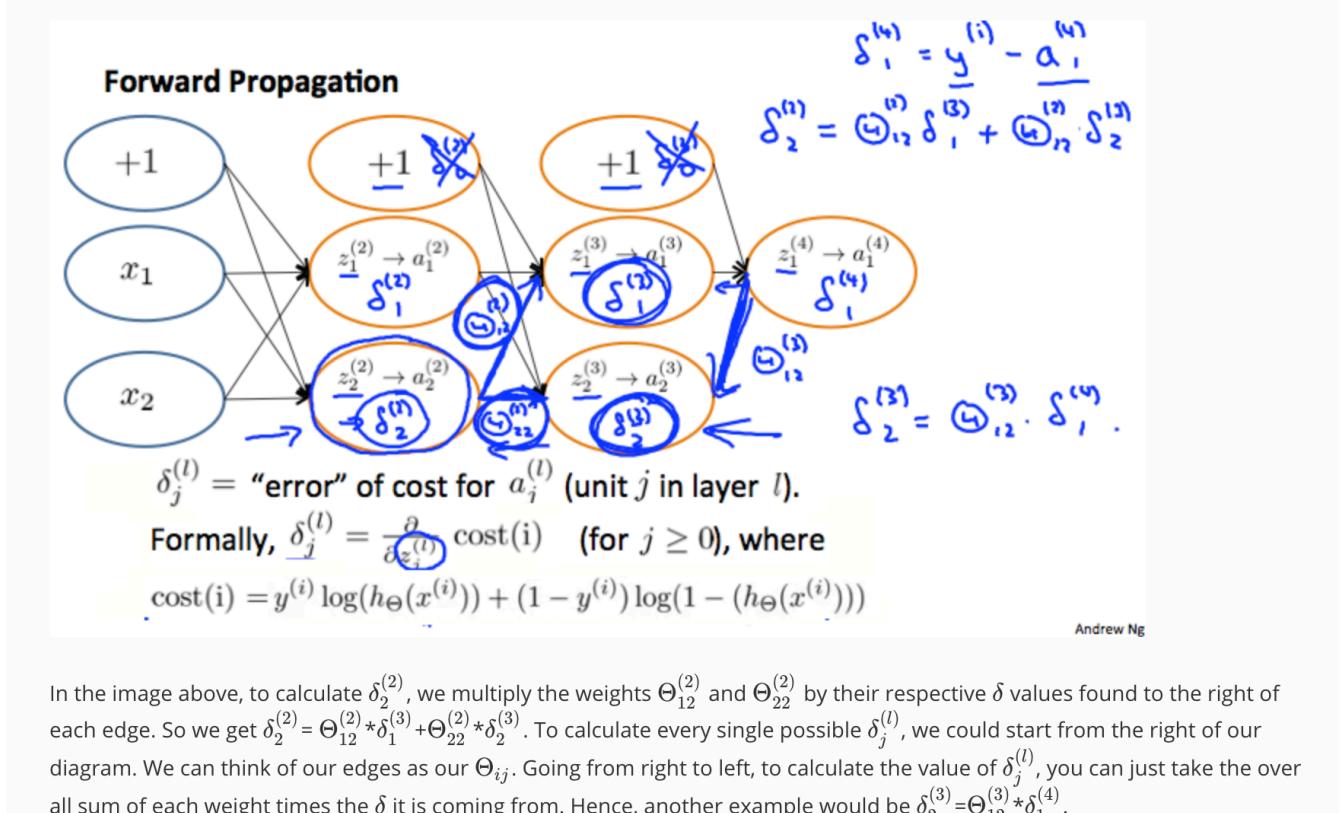
### (3) Backpropagation Intuition

- Forward Propagation

$$z_1^{(l)} = \Theta_{10}^{(l-1)} * 1 + \Theta_{11}^{(l-1)} * a_1^{(l-1)} + \Theta_{12}^{(l-1)} * a_2^{(l-1)} \quad (61)$$

Focusing on a single example  $x^{(i)}$   $y^{(i)}$  ignore regularization

$$cost = -[y^{(i)} \log h_{\Theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\Theta}(x^{(i)}))] \quad (62)$$



In the image above, to calculate  $\delta_2^{(2)}$ , we multiply the weights  $\Theta_{12}^{(2)}$  and  $\Theta_{22}^{(2)}$  by their respective  $\delta$  values found to the right of each edge. So we get  $\delta_2^{(2)} = \Theta_{12}^{(2)} * \delta_1^{(1)} + \Theta_{22}^{(2)} * \delta_2^{(3)}$ . To calculate every single possible  $\delta_j^{(l)}$ , we could start from the right of our diagram. We can think of our edges as our  $\Theta_{ij}$ . Going from right to left, to calculate the value of  $\delta_j^{(l)}$ , you can just take the overall sum of each weight times the  $\delta$  it is coming from. Hence, another example would be  $\delta_2^{(3)} = \Theta_{12}^{(3)} * \delta_1^{(2)}$ .

## (4) Backpropagation in Practice

- Add all matrices of each layer into one big vector

```

1 thetaVec = [Theta1(:); Theta2(:); Theta3(:)];
2 % reshape row 1 to 110 to a 10*11 matrix
3 Theta1 = reshape(thetaVec(1:110), 10, 11)

```

### Learning Algorithm

- Have initial parameters  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ .
- Unroll 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 prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ .
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.

```

- Gradient checking

- use an approximate value  $\frac{J(\theta+\epsilon) - J(\theta-\epsilon)}{2\epsilon}$

$\rightarrow \theta \in \mathbb{R}^n$  (E.g.  $\theta$  is “unrolled” version of  $\underline{\Theta^{(1)}}, \underline{\Theta^{(2)}}, \underline{\Theta^{(3)}}$ )

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

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

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

$\vdots$

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

```
1 for i=1:n,
2     thetaplus = theta;
3     thetaplus(i) = thetaplus(i) + EPSILON;
4     thetaMinus = theta;
5     thetaMinus(i) = thetaMinus(i) - EPSILON;
6     gradApprox(i) = (J(thetaPlus)-J(thetaMinus))/(2*EPSILON);
7 end;
```

- Check that  $\text{gradApprox} \approx \text{DVec}$  (but it is slow)