

1 Linear algorithms

1.1 Definitions

def. Regression problem - our goal is to predict a **continuous** valued output.

def. Classification problem - our goal is to predict a **discrete** valued output.

def. Supervised learning - we have data set, which matches some response to the input data. (Having the idea that there is a relationship between the input and the output.)

def. Unsupervised learning - we have a data set and we don't know what to do with it. Our unsupervised learning should highlight some structure, **clusters**. So this is called a **clustering algorithm**.

As an example of the application of this algorithm, we can consider the separation of conversation from music and music from conversation, this is called the **cocktail party problem algorithm**.

It can be written in one line:

$$[W, s, v] = \text{svd}((\text{ repmat}(\text{sum}(x .* x, 1), \text{size}(x, 1), 1) .* x) * x');$$

1.2 What is machine learning ?

Example: playing checkers.

E = the **experience** of playing many games of checkers

T = the **task** of playing checkers

P = the **probability** that the program will win the next game

1.3 What about development environment?

You need to use Octave or Matlab, because you can write code in one line in them, as shown in paragraph 1.1.

1.4 Some notations

1. m - number of training examples

2. x's - input features

3. y's - output target variable

(x, y) - one training example

$(x^{(i)}, y^{(i)})$ - i^{th} training example

$h_{\Theta}(x) = h(x) = \Theta_0 + \Theta_1 \cdot x$ - **linear regression** or **univariate linear regression**

h is called a hypothesis, $h : X \rightarrow Y$

1.5 Cost function (linear regression)

Idea: Choose Θ_0, Θ_1 so that $h_{\Theta}(x)$ is close to y for our training examples (x, y)

Squared error (function):

$$\frac{1}{2m} \cdot \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right)^2 = \mathbb{J}(\Theta_0, \Theta_1) \xrightarrow{\Theta_0, \Theta_1} \min, h_{\Theta}(x^{(i)}) = \widehat{y^{(i)}} \quad (1)$$

Why $\frac{1}{2}$? The mean is halved as a convenience for the computation of the gradient descent, as the derivative term of the square function will cancel out the 0.5 term.

We can simplify it and set $\Theta_0 = 0$, then \mathbb{J} is a function of one variable.

Then let's take a data set of the following pairs $(1,1);(2,2);(3,3)$. We will take different values $\Theta_1 = 0, 0.5, 1$ in our simplified formula. We will substitute them into the function for h , and then for \mathbb{J} . As a result, we will know the values of \mathbb{J} and will be able to plot \mathbb{J} . From the constructed graph, the minimum value is visible, and therefore the answer.

def. A **contour plot** is a graph that contains many contour lines. A contour line of a two variable function has a constant value at all points of the same line.

If we have both parameters involved, then this graph can be drawn in **3D**:

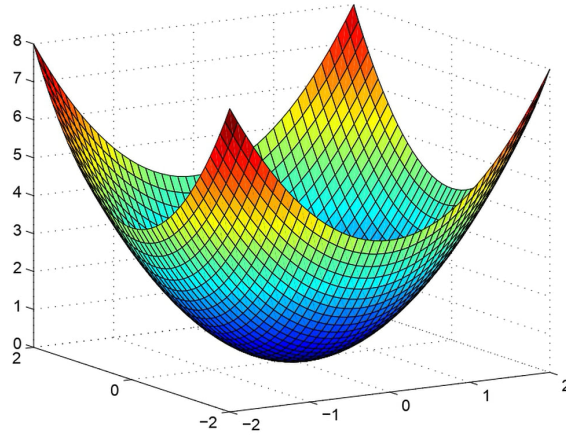


Figure 1: Cost function in 3D

1.6 Gradient Descent (algorithm)

In the previous paragraph, we drew a lot of graphs and built a lot of functions. Now we want to understand how to quickly find these parameters in which the value of the function \mathbb{J} reaches its minimum. Then we can build a perfect straight line to find the answer.

A brief idea of the algorithm:

- Start with some Θ_0, Θ_1 (maybe = 0).
- Keep changing Θ_0, Θ_1 to reduce $\mathbb{J}(\Theta_0, \Theta_1)$ until we hopefully end up at a minimum. (or local minimum)

Now maths, repeat until convergence:

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} \cdot \mathbb{J}(\theta_0, \theta_1), j = 0..1 \quad (2)$$

For example (simultaneous update):

$$\begin{aligned} \text{temp0} &:= \theta_0 - \alpha \cdot \frac{\partial}{\partial \theta_0} \cdot \mathbb{J}(\theta_0, \theta_1) \\ \text{temp1} &:= \theta_1 - \alpha \cdot \frac{\partial}{\partial \theta_1} \cdot \mathbb{J}(\theta_0, \theta_1) \\ \theta_0 &:= \text{temp0} \text{ and } \theta_1 := \text{temp1} \end{aligned}$$

Minor clarification:

α - this is **learning rate**, determines the step size of our algorithm, if α is very large, then gradient descent is very aggressive.

For example, you can again consider a function from one variable and draw a parabola. If we look at the right branch, then our derivative is > 0 and we move to the left with step α , otherwise vice versa.

The α parameter is very important, because if it is very small, then the algorithm will be very slow, if it is large, then we will jump over the minimum or even begin to **diverge**.

And as gradient descent runs, you will automatically take smaller and smaller steps, therefore, it is not necessary to reduce the α .

1.7 Gradient Descent for linear regression

Let's combine gradient descent and regression together:

$$\frac{\partial}{\partial \theta_j} \mathbb{J}(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right)^2 = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m \left(\Theta_0 + \Theta_1 x^{(i)} - y^{(i)} \right)^2 \quad (3)$$

Take partial derivatives:

$$j = 0 : \frac{\partial}{\partial \theta_0} \mathbb{J}(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right)$$
$$j = 1 : \frac{\partial}{\partial \theta_1} \mathbb{J}(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$

In total, we have the transition of our first algorithm:

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right) \quad (4)$$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)} \quad (5)$$

Function for linear regression is always **convex** function. Another name for this algorithm is "Batch" Gradient Descent. **"Batch"**: each step of gradient descent uses all the training examples. Now we know how to implement gradient descent for linear regression.

1.8 Linear Algebra Review

How to use linear algebra in Matlab:

```
1 A = [1, 2, 3; 4, 5, 6; 7, 8, 9; 10, 11, 12]
2 v = [1; 2; 3]
3 [m,n] = size(A)
4 dim_A = size(A)
5 dim_v = size(v)
6 A_23 = A(2,3)
```

```
1 s = 2
2 add_AB = A + B
3 sub_AB = A - B
4 mult_As = A * s
5 div_As = A / s
6 add_As = A + s
```

Matrices in machine learning: $A = \begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$, $h_{\theta}(x) = -40 + 0.25 \cdot x \Rightarrow$ result is $A \cdot \begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$

Identity matrix in Matlab, transposed and inversed matrix:

```

1 A = [2, 4; 5, 6]
2 % Initialize a 2 by 2 identity matrix
3 I = eye(2)
4 % The above notation is the same as I = [1,0;0,1]
5 % Inversed matrix
6 pinv(A)
7 inv(A)
8 % Transposed matrix
9 transpose(A)
10 A'
```

2 Linear Regression with Multiple Variables, Octave, Matlab

The basic skills of using Matlab can be found here: <https://www.coursera.org/learn/machine-learning/supplement/Mlf3e/more-octave-matlab-resources>

Here you can find Matlab online: <https://matlab.mathworks.com/>

2.1 Multivariate Linear Regression

Notations:

n = number of features

$x^{(i)}$ = input (features) of i^{th} training example

$x_j^{(i)}$ = value of feature j in i^{th} training example

Hypothesis:

$$h_{\Theta}(x) = \Theta_0 + \Theta_1 x_1 + \dots + \Theta_n x_n$$

For **convenience** of notation, define $x_0 = 1$. This is an additional zero feature for each training example: $x_0^{(i)} = 1$.

Matrix notation:

$$\mathbb{X} = \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}, \quad \Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \dots \\ \Theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

Now you can write like this:

$$h_{\Theta}(x) = \Theta_0 x_0 + \Theta_1 x_1 + \dots + \Theta_n x_n = \Theta^T \mathbb{X} \quad (6)$$

2.2 Gradient Descent For Multiple Variables

The algorithm is very similar to the one we have already passed, so I'll just write down the transition formula:

$$\text{repeat until convergence: } \theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m \left(h_{\Theta}(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)} \quad \text{for } j := 0 \dots n$$

2.3 What can be done to make the algorithm work faster?

Idea: Let's make sure that all the features belong to the same range, for example $-1 \leq x_i \leq 1$, then the algorithm will converge to the global minimum much faster. Our elongated contour ellipses will become circles.

Mean normalization: Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean (Do not apply to $x_0 = 1$). For example, all features: $-0.5 \leq x_i \leq 0.5$. General formula: $x_i = \frac{x_i - \mu_i}{s_i}$, where μ_i - average value of x_i in training set, s_i - range, max - min (or standard deviation).

Debugging: How to make sure that the descent works correctly? Let's plot $J(\theta)$ from the number of iterations (100-400). It should decrease like hyperbole. From the graphics, you can also understand when it has become close to the horizontal line, so you can correctly choose the number of iterations to convergence.

Automatic convergence test: Declare convergence if $J(\theta)$ decreases by less than $\epsilon = 10^{-3}$ in one iteration.

2.4 Polynomial regression

Perhaps a quadratic or cubic function is better suited than just a straight line, so you can modify the multivariate linear regression to a polynomial.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 = \theta_0 + \theta_1(var) + \theta_2(var)^2 + \theta_3(var)^3$$

And just in this algorithm, it is very important to adjust the values to the same range. Perhaps it is better not to switch to the cubic function, but to replace the quadratic term with the square root: \sqrt{var} .

2.5 Computing Parameters Analytically

Normal equation: Method to solve for θ analytically.

$$\begin{aligned} \theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_n) &= \frac{1}{2m} \cdot \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \\ \frac{\partial}{\partial \theta_j} J(\theta) &= \dots = 0 \quad (\text{for every } j) \\ \text{Solve for } \theta_0, \theta_1, \dots, \theta_n \end{aligned}$$

Let's solve:

$$\mathbb{X}_{m \times (n+1)} = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad y_{m \times 1} = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

Solution: $\theta = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T y$, \mathbb{X} -design matrix, $\mathbb{X}_i = (x^{(i)})^T$.

In this method, there is **no need** to use features scaling.

Problem: what if the inverse matrix doesn't exist? There are two functions in Matlab: *inv* and *pinv*. The latter will work even if the matrix is not reversible. There are two possible reasons: **redundant features** (linear dependent) or too **many features** ($m \leq n$) \Rightarrow delete some features or use regularization (later in this course).

3 Implementation

3.1 Implementation of the gradient descent algorithm with one variable

```
1 % load data (comma separated)
2 data = load('ex1data1.txt');
3 X = data(:, 1);
4 y = data(:, 2);
5
6 % plot
7 plotData(X,y);
8
9 % initialization
10 m = length(X);
11 X = [ones(m,1), data(:,1)];
12 theta = zeros(2, 1);
13 iterations = 1500;
14 alpha = 0.01;
15
16 % Compute theta, run gradient descent
17 theta = gradientDescent(X, y, theta, alpha, iterations);
18
19 % Print theta to screen
20 fprintf('Theta computed from gradient descent:\n%f,\n%f\n', theta(1), theta(2));
21
22 fprintf('J-function value: %f\n', computeCost(X, y, theta));
23
24 hold on; % keep previous plot visible
25 plot(X(:,2), X*theta, '-')
26 legend('Training data', 'Linear regression')
27 hold off % don't overlay any more plots on this figure
28
29 % Predict value for population size of 70,000
30 predict = [1, 7] * theta;
31 fprintf('For population = 70,000, we predict a profit of %f\n', predict*10000);
32
33 % drawing a graph of our data
34 function plotData(x, y)
35 figure;
36 plot(x, y, 'rx', 'MarkerSize', 5);
37 ylabel('Profit in $10,000s');
38 xlabel('Population of City in 10,000s');
39 end
40
41 % Gradient Descent function
42 function [theta, J_history] = gradientDescent(X, y, theta, alpha, num_iters)
43 m = length(y);
44 J_history = zeros(num_iters, 1);
45
46 for iter = 1 : num_iters
47     sum_theta1 = 0;
48     sum_theta2 = 0;
49     for i = 1 : m
```

```

50         sum_theta1 = sum_theta1 + (theta' * X(i, :) - y(i)) * X(i, 1);
51         sum_theta2 = sum_theta2 + (theta' * X(i, :) - y(i)) * X(i, 2);
52     end
53     theta1 = theta(1) - alpha / m * sum_theta1;
54     theta2 = theta(2) - alpha / m * sum_theta2;
55     theta = [theta1; theta2];
56 end
57 end
58
59 % J-function
60 function J = computeCost(X, y, theta)
61 m = length(y);
62 J = 0;
63 for i = 1 : m
64     J = J + (theta' * X(i, :) - y(i)) ^ 2;
65 end
66 J = J / (2 * m);
67 end

```

3.2 Visualizing $\mathbb{J}(\theta)$

```

1 % Grid over which we will calculate J
2 theta0_vals = linspace(-10, 10, 100);
3 theta1_vals = linspace(-1, 4, 100);
4
5 % Initialize J_vals to a matrix of 0's
6 J_vals = zeros(length(theta0_vals), length(theta1_vals));
7
8 % Fill out J_vals
9 for i = 1 : length(theta0_vals)
10     for j = 1 : length(theta1_vals)
11         t = [theta0_vals(i); theta1_vals(j)];
12         J_vals(i, j) = computeCost(X, y, t);
13     end
14 end
15
16 % Surface plot
17 J_vals = J_vals';
18 figure;
19 surf(theta0_vals, theta1_vals, J_vals)
20 xlabel('\theta_0'); ylabel('\theta_1');
21
22 % Contour plot
23 figure;
24 % Plot J_vals as 15 contours spaced logarithmically between 0.01 and 100
25 contour(theta0_vals, theta1_vals, J_vals, logspace(-2, 3, 20))
26 xlabel('\theta_0'); ylabel('\theta_1');
27 hold on;
28 plot(theta(1), theta(2), 'rx', 'MarkerSize', 10, 'LineWidth', 2);
29 hold off;

```

3.3 Implementation of the gradient descent algorithm with two variable

```

1 % Load Data
2 data = load('ex1data2.txt');
3 X = data(:, 1:2);
4 y = data(:, 3);
5 m = length(y);

```

```

6
7 % First 10 examples from the dataset
8 fprintf(' x = [%0.0f %0.0f], y = %0.0f \n', [X(1:10,:) y(1:10,:)] ');
9 % Scale features and set them to zero mean
10 [X, mu, sigma] = featureNormalize(X);
11
12 % Add intercept term to X
13 X = [ones(m, 1) X];
14
15 % Feature Normalize
16 function [X_norm, mu, sigma] = featureNormalize(X)
17 X_norm = X;
18 mu = zeros(1, size(X, 2));
19 sigma = zeros(1, size(X, 2));
20 cols = size(X, 2);
21 for col = 1:cols
22     cur_mean = mean(X(:, col));
23     cur_std = std(X(:, col));
24     X_norm(:, col) = (X(:, col) - cur_mean) / cur_std;
25     mu(1, col) = cur_mean;
26     sigma(1, col) = cur_std;
27 end
28 end
29
30 % J-function multi
31 function J = computeCostMulti(X, y, theta)
32 m = length(y);
33 sqr_term = (X * theta - y);
34 J = sqr_term' * sqr_term / (2 * m);
35 end

```

$$\mathbb{J}(\theta) = \frac{1}{2m} \cdot (X\theta - \vec{y})^T \cdot (X\theta - \vec{y}), \text{ vectorized form of cost function} \quad (7)$$

```

1 % Gradient Descent
2 function [theta, J_history] = gradientDescentMulti(X, y, theta, alpha, num_iters)
3 m = length(y);
4 J_history = zeros(num_iters, 1);
5 cols = size(X, 2);
6
7 for iter = 1 : num_iters
8     sum_theta = zeros(cols, 1);
9     for i = 1 : m
10         for col = 1 : cols
11             sum_theta(col, 1) = sum_theta(col, 1) +
12                 (theta' * X(i, :) - y(i)) * X(i, col);
13         end
14     end
15     for col = 1 : cols
16         theta(col, 1) = theta(col, 1) - alpha / m * sum_theta(col, 1);
17     end
18     J_history(iter) = computeCostMulti(X, y, theta);
19 end
20 end
21
22 % Run gradient descent
23 alpha = 0.1;
24 num_iters = 400;
25
26 % Init Theta and Run Gradient Descent
27 theta = zeros(3, 1);
28 [theta, ~] = gradientDescentMulti(X, y, theta, alpha, num_iters);
29

```



```

30 % Display gradient descent's result
31 fprintf('Theta computed from gradient descent:\n%f\n%f\n%f',
32 theta(1),theta(2),theta(3))
33 % Estimate the price of a 1650 sq-ft, 3 br house
34 size_of_flat = (1650 - mu(1)) / sigma(1);
35 amount_of_rooms = (3 - mu(2)) / sigma(2);
36 price = theta' * [1; size_of_flat; amount_of_rooms];
37 fprintf('Predicted price of a 1650 sq-ft, 3 br house (using gradient descent):
38 \n $%f', price);

```

3.4 Implementation of Normal Equations

```

1 % Solve with normal equations:
2 % Load Data
3 data = csvread('ex1data2.txt');
4 X = data(:, 1:2);
5 y = data(:, 3);
6 m = length(y);
7
8 % Add intercept term to X
9 X = [ones(m, 1) X];
10
11 % Calculate the parameters from the normal equation
12 theta = normalEqn(X, y);
13
14 % Display normal equation's result
15 fprintf('Theta computed from the normal equations:\n%f\n%f\n%f',
16 theta(1),theta(2),theta(3));
17
18 % Estimate the price of a 1650 sq-ft, 3 br house
19 price = theta' * [1; 1650; 3];
20 fprintf('Predicted price of a 1650 sq-ft, 3 br house (using gradient descent):
21 \n $%f', price);
22
23 % Normal Equations function
24 function [theta] = normalEqn(X, y)
25 theta = pinv(X' * X) * X' * y;
26 end

```

3.5 Basic Matlab features

```

1 xor(1, 0) % xor
2 1 ~= 2 % not equals
3 PS1('>> ') % Octave function to change prompt style
4 a = 3; % semicolon supressing output
5
6 a = pi;
7 disp(a) % just print
8 disp(sprintf('%0.8f', a)) % C-style output = 3.14159265
9
10 % set 3.141592653589793 or 3.1416 by default output
11 format long;
12 format short;
13
14 % start from 1, end with 10, step is 0.5
15 v = 1:0.5:10
16 % generates vector with step 1
17 v_step_one = 1:10

```

```

18 m = ones(4, 5) * 4 % matrix of 4 (zeros - of 0)
19
20 % random matrix of random numbers: 0 <= x <= 1
21 m = rand(3, 3)
22
23 % Gaussian distribution with mean zero and std equal to one
24 g = randn(2, 2)
25
26 % Example of Gaussian distribution
27 w = -6 + sqrt(10) * (randn(1, 10000));
28 hist(w, 100)

```

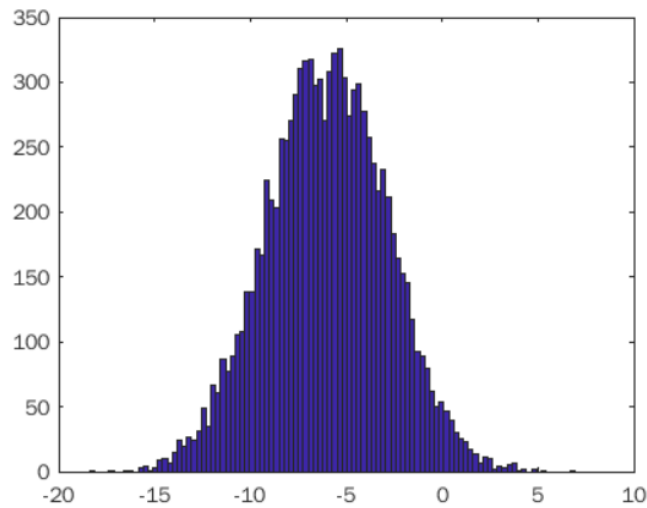


Figure 2: Gaussian distribution

3.6 Working with data

```

1 % size of matrix
2 rows = size(x, 1)
3 cols = size(x, 2)
4
5 length(x) % the biggest dimension of x
6
7 % loading data
8 load ex1data1.txt
9 % or <=>
10 load('ex1data1.txt')
11
12 who % prints all variables in our workspace
13 whos % detail info
14
15 clear var % delete variable var
16
17 % saving data
18 v = rand(1, 5)
19 % compressed binary format
20 save res.mat v;
21

```

```

22 % format for human, you can read it
23 save res.txt v -ascii;
24
25 % select some rows and some cols
26 A([1 3], [4 5])
27
28 % append column [1; 2; 3] to the end of matrix A
29 A = [A, [1; 2; 3]];
30
31 A(:) % convert matrix into vector of matrix's elements
32
33 % pow corresponding elements
34 % C(1, 1) = A(1, 1) ^ B(1, 1)
35 C = A .^ B
36
37 % similarly we can do
38 X = 1 ./ A
39 % A .* X = ones
40 % similarly log(A), exp(A), abs(A), -A
41
42 % max, min
43 v = rand(1, 5) * 4
44 [val, ind] = max(v)
45 % or just max(v)
46
47 b = v < 10 % element-wise operation
48 find(b) % returns indexes which satisfy condition
49
50 res = magic(4)
51 % returns a matrix where sums over all columns, rows, diagonals are equal
52 % to one number
53 [r, c] = find(res >= 7)
54
55 a = [1, 2, 3]
56 % sum(a), prod(a), floor(a), ceil(a)
57
58 % element-wise maximum in two random matrices
59 max(rand(3, 3), rand(3, 3))
60
61 max(A, [], 1) % columns-wise max
62 max(A, [], 2) % rows-wise max
63
64 max(max(A)) % max(matrix) or
65 max(A(:))
66
67 sum(A, 1) % columns sums
68 sum(A, 2) % rows sums
69
70 sum(sum(A .* eye(3))) % main diag sum
71 sum(sum(A .* flipud(eye(3)))) % secondary diag sum
72 flipud(A) % turns the matrix upside down

```

3.7 Plotting data

```

1 t = [0:0.01:0.98];
2 y1 = sin(2*pi*4*t);
3 plot(t, y1);
4
5 y2 = cos(2*pi*4*t);
6 plot(t, y1);
7 hold on;
8 plot(t, y2, 'r'); % 'r' - red color

```

```

9  xlabel('time');
10 ylabel('value');
11 legend('sin', 'cos');
12 title('sin and cos');
13 print -dpng 'sin_and_cos.png'; % save your plot

```

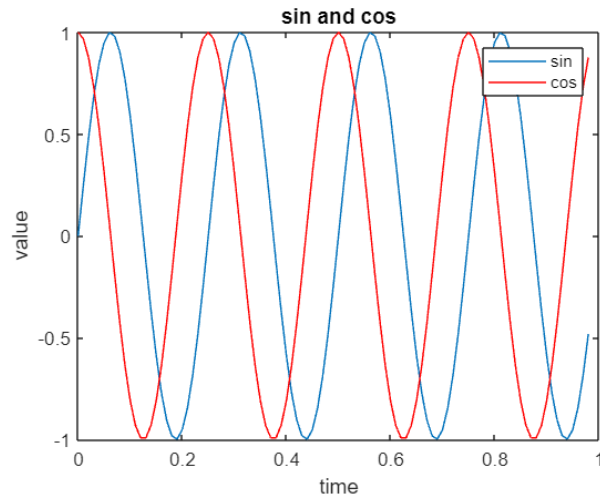


Figure 3: Plotting data

```

1  close; % closes all opened figures
2
3  % If you want two different figures:
4  figure(1); plot(t, y1);
5  figure(2); plot(t, y2);
6
7  % What if I want many plots on one image?
8  % I can divide my screen into grid!
9  subplot(1, 3, 1); % divide 1x3 grid and access first
10 plot(t, y1);
11 subplot(1, 3, 2); % access second
12 plot(t, y2);
13 subplot(1, 3, 3); % access third
14 plot(t, y1);
15
16 % How to change axis?
17 axis([-10 10 -1 1]);
18 % set x from -10 to 10
19 % set y from -1 to 1
20
21 clf; % clear figure, not delete
22
23 % I want to print my matrix
24 % Where for different elements different colors
25 A = [1 2 3; 3 2 1; 4 5 1];
26 % imagesc(A);
27 imagesc(A), colorbar, colormap autumn;
28 % second is more complicated command
29 % it plots matrix, sets color map to autumn
30 % and prints color bar on the right of our plot
31 % here is an example:

```

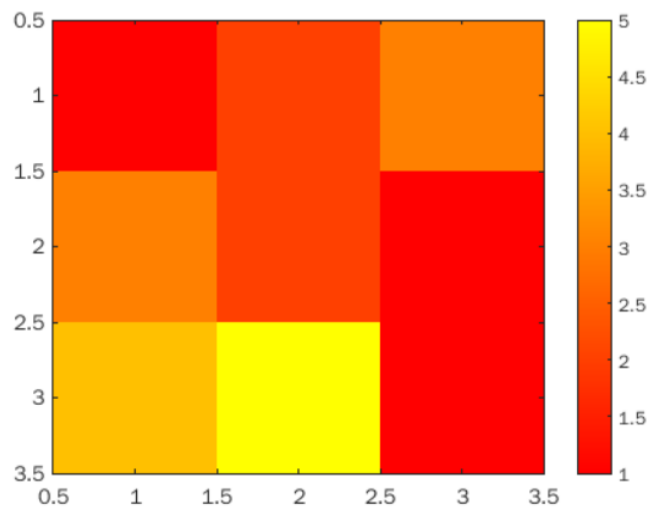


Figure 4: Plotting matrix

3.8 Control Statements

```

1  % `for` loop
2  vals=1:2:7;
3  for i=vals
4      disp(i);
5  end
6
7  % `while` loop
8  % `break`, `continue` statements
9  i = 1;
10 while i <= 10
11     disp(i);
12     i = i + 1;
13     if i == 6
14         disp("end")
15         break
16     end
17 end
18
19 % `if`-`else` statement
20 x=1;
21 if x==1, disp("one")
22 elseif x==2, disp("two")
23 else , disp("another")
24 end
25
26 pwd, cd, addpath % Octave functions
27 % current dir, move to another dir, add path to PATH
28
29 % definition of functions
30 % all data, that returns function
31 [ans1, ans2] = SquareAndCubeFunction(5);
32 fprintf('%d, %d\n', ans1, ans2);
33 % only first variable, that returns function
34 fprintf('%d\n', SquareAndCubeFunction(5));
35

```

```

36 % function that returns square and cube
37 function [square, cube] = SquareAndCubeFunction(x)
38 square = x ^ 2;
39 cube = x ^ 3;
40 end

```

3.9 Vectorization

Vectorization is the idea of using functions from standard libraries. This reduces the code, increases its readability. These functions are highly optimized and can use parallelism, so the code will run much faster and more efficiently.

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

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

$$\Theta := \Theta - \alpha \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) * x^{(i)}$$

4 Logistic Regression

4.1 Classification and Representation

4.1.1 Classification

Email: Spam / Not Spam ?

Online transactions: Fraudulent (Yes / No) ?

Tumor: Malignant / Benign ?

$y \in \{0, 1\}$

0: "Negative Class" (benign tumor)

1: "Positive Class" (malignant tumor)

Threshold classifier output $h_{\theta}(x)$ at 0.5: (linear regression)

if $h_{\theta}(x) \geq 0.5$, predict "y = 1"

if $h_{\theta}(x) < 0.5$, predict "y = 0"

This method can work on a specific example, but in most cases it will work extremely poorly.

It may even happen that $h_{\theta}(x)$ can be > 1 or < 0 .

Therefore, we will develop a **logistic regression** algorithm: $0 \leq h_{\theta}(x) \leq 1$.

4.1.2 Hypothesis Representation

$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}} \text{ - sigmoid (logistic) function}$$

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

$h_{\theta}(x)$ = estimated probability that $y = 1$ on input x , or in the language of mathematics: $h_{\theta}(x) = \mathbb{P}(y = 1 | x; \theta)$.

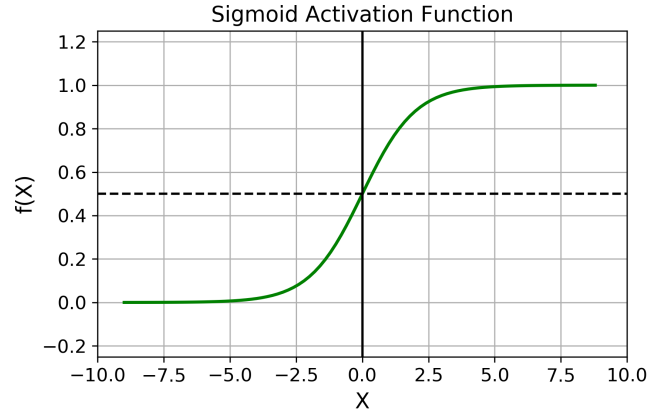


Figure 5: Sigmoid function

The relationship of probabilities:

$$\begin{aligned}\mathbb{P}(y = 0|x; \theta) + \mathbb{P}(y = 1|x; \theta) &= 1 \Leftrightarrow \\ \mathbb{P}(y = 0|x; \theta) &= 1 - \mathbb{P}(y = 1|x; \theta)\end{aligned}$$

4.1.3 Decision Boundary

How to understand what we predicted 0 or 1: $g(z) \geq 0.5 \Leftrightarrow z \geq 0 \Leftrightarrow \theta^T x \geq 0$. Accordingly, if we reverse all the signs, then we predict 0. By writing $\theta^T x$, we can build a figure on the graph that separates our sets. (line, parabola, circle) This equation is called **decision boundary**.

4.2 Logistic Regression Model

4.2.1 Cost Function

If we leave our cost function as it was in linear regression, then given that our hypothesis is very complex, we get a non-convex function. This function will have many local optima.

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \quad (8)$$

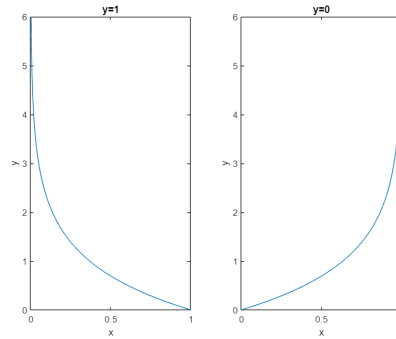


Figure 6: Cost Function for Logistic Regression

4.2.2 Simplified Cost Function and Gradient Descent

$$\mathbb{J}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost} \left(h_{\theta}(x^{(i)}), y^{(i)} \right) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta} \left(x^{(i)} \right) + (1 - y^{(i)}) \log \left(1 - h_{\theta} \left(x^{(i)} \right) \right) \right]$$
$$\theta_j := \theta_j - \frac{1}{m} \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Vector form of writing our formulas:

$$h = g(\mathbb{X}\theta)$$
$$\mathbb{J}(\theta) = -\frac{1}{m} \cdot (y^T \log(h) + (1 - y)^T \log(1 - h))$$
$$\theta := \theta - \frac{\alpha}{m} \cdot \mathbb{X}^T (g(\mathbb{X}\theta) - \vec{y})$$

To simplify the cost function, we merged this system into one formula, so it's easier to count and calculate the derivative. It turned out to be a very interesting coincidence: the gradient descent formulas coincided for linear regression and logistic regression. Vector formulas are obvious and do not require explanation.

4.2.3 Advanced Optimization

Optimization algorithms:

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:

- No need to manually pick α . In fact, they have a clever inter-loop called a **line search algorithm** that automatically tries out different values for the learning rate α and automatically picks a good learning rate α so that it can even pick a different learning rate for every iteration. And so then you don't need to choose it yourself.
- Often faster than gradient descent

Disadvantages:

- More complex

To use these advanced algorithms, you need to implement a special function:

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

Now we can call this function: (exitFlag = did the algorithm converge?!)

```
1 options = optimset('GradObj', 'on', 'MaxIter', 100);
2 initialTheta = zeros(2, 1);
3 [optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta, options);
```


4.3 Multiclass Classification

4.3.1 Multiclass Classification: One-vs-all

Email foldering/tagging: Work, Friends, Family, Hobby

Medical diagrams: Not ill, Cold, Flu

Weather: Sunny ($y=1$), Cloudy ($y=2$), Rain ($y=3$), Snow ($y=4$)

How will the algorithm work? We sort through our sets in order. We take the current set and denote it as a positive class, and everything else as a negative class. We run logistic regression on this data and repeat the same with another set. Each logistic regression classifier we denote as $h_{\theta}^{(i)}(x)$. We fit n classifiers: $h_{\theta}^{(i)}(x) = \mathbb{P}(y = i|x; \theta)$.

On a new input x , to make a prediction, pick the class i that maximizes $\max_i \left(h_{\theta}^{(i)}(x) \right)$.

5 Regularization

5.1 Solving the Problem of Overfitting

5.1.1 The Problem of Overfitting

Let's say we have data that can be depicted on a plane as part of a function x^4 . We can train a linear model $\theta_0 + \theta_1 x$, but then our predictions will not correspond to reality at all accurately. This problem is called **underfitting**, sometimes it is said that the algorithm has **high bias**. If we train a quadratic model, then it will perfectly fit our data. If we take more parameters, for example, a function of the fourth degree, then even if our curve passes through all the given ones, it will still be **wavy**. They say about it like this: algorithm has **high variance** or problem **overfitting** (too many features). Even if the value of the \mathbb{J} function is approximately 0, then everything will be bad in the new examples.

Options:

- Reduce number of features.
 - Manually select which features to keep.
 - Model selection algorithm (later in course).
- Regularization.
 - Keep all the features, but reduce magnitude / values of parameters θ_j .
 - Works well when we have a lot of features, each of which contributes a bit to predicting y .

5.1.2 Cost Function

Suppose we penalize and make θ_3, θ_4 really small:

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

There may be a situation when it is not clear which parameters we want to compress. Then we have to compress all the parameters, for this we will add another term to our cost function.

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

This is regularization term, and λ is regularization parameter. The function is smoother.

5.1.3 Regularized Linear Regression

It is important to remember that we compress all θ values **except** θ_0 .

$$\begin{aligned}\theta_0 &:= \theta_0 - \alpha \frac{1}{m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)} \\ \theta_j &:= \theta_j - \alpha \cdot \left[\frac{1}{m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \iff \\ \theta_j &:= \theta_j \cdot \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}\end{aligned}$$

The term $1 - \alpha \frac{\lambda}{m} \approx 0.99$, so it shrinks θ_j on each iteration.

Normal equation:

$$\Theta = \left(\mathbb{X}^T \mathbb{X} + \lambda \cdot \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & 0 & 1 & \dots \\ 0 & \dots & 0 & 1 \end{bmatrix}_{(n+1) \times (n+1)} \right)^{-1} \cdot \mathbb{X}^T y$$

Interesting fact: if $\lambda > 0$, then matrix is **invertible** (this can be proven). What is written in square brackets is the corresponding partial derivative.

5.1.4 Regularized Logistic Regression

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

The transition formulas for gradient descent are the same as for linear regression.

6 Implementation

6.1 Required Functions

6.1.1 Cost Function, Gradient

```
1 function [J, grad] = costFunction(theta, X, y)
2 m = length(y);
3 h = sigmoid(X * theta);
4 J = -1 / m * (y' * log(h) + (1 - y)' * log(1 - h));
5 grad = X' * (h - y) / m;
6 end
```

6.1.2 Plotting Data

```
1 function plotData(X, y)
2 figure;
3 hold on;
4 pos = find(y == 1);
5 neg = find(y == 0);
6 plot(X(pos, 1), X(pos, 2), 'k+', 'LineWidth', 2, 'MarkerSize', 7);
```

```

7 plot(X(neg, 1), X(neg, 2), 'ko', 'MarkerFaceColor', 'y','MarkerSize', 7);
8 hold off;
9 end

```

6.1.3 Prediction

```

1 function p = predict(theta, X)
2 p = X * theta >= 0;
3 end

```

6.1.4 Sigmoid

```

1 function g = sigmoid(z)
2 g = 1 ./ (1 + exp(-z));
3 end

```

6.1.5 Map Feature

```

1 function out = mapFeature(X1, X2)
2 degree = 6;
3 out = ones(size(X1(:,1)));
4 for i = 1:degree
5     for j = 0:i
6         out(:, end + 1) = (X1 .^ (i - j)) .* (X2 .^ j);
7     end
8 end
9 end

```

6.1.6 Plotting Decision Boundary

```

1 function plotDecisionBoundary(theta, X, y)
2 plotData(X(:, 2:3), y);
3 hold on
4 if size(X, 2) <= 3
5     plot_x = [min(X(:, 2)), max(X(:, 2))];
6     plot_y = -(theta(2) * plot_x + theta(1)) / theta(3);
7     plot(plot_x, plot_y)
8     legend('Admitted', 'Not admitted', 'Decision Boundary')
9     axis([30, 100, 30, 100])
10 else
11     u = linspace(-1, 1.5, 50);
12     v = linspace(-1, 1.5, 50);
13     z = zeros(length(u), length(v));
14     for i = 1:length(u)
15         for j = 1:length(v)
16             z(i,j) = mapFeature(u(i), v(j)) * theta;
17         end
18     end
19     z = z';
20     contour(u, v, z, [0, 0], 'LineWidth', 2)
21 end
22 hold off

```

6.1.7 Regularized Cost Function, Gradient

```

1 function [J, grad] = costFunctionReg(theta, X, y, lambda)
2 m = length(y);
3 n = size(theta, 1);
4 h = sigmoid(X * theta);
5 s = 0;
6 for i=2:n
7     s = s + theta(i) ^ 2;
8 end
9 add = lambda / (2 * m) * s;
10 J = -1 / m * (y' * log(h) + (1 - y)' * log(1 - h)) + add;
11 grad = X' * (h - y) / m;
12 shift = lambda / m * theta;
13 grad(2:end) = grad(2:end) + shift(2:end);
14 end

```

6.2 Launching

6.2.1 First part

```

1 data = load('ex2data1.txt');
2 X = data(:, [1, 2]);
3 y = data(:, 3);
4
5 plotData(X, y);
6 xlabel('Exam 1 score')
7 ylabel('Exam 2 score')
8 legend('Admitted', 'Not admitted')
9
10 [m, n] = size(X);
11 X = [ones(m, 1) X];
12 initial_theta = zeros(n + 1, 1);
13
14 options = optimoptions(@fminunc, 'Algorithm', 'Quasi-Newton', 'GradObj', 'on',
15 'MaxIter', 400);
16 [theta, cost] = fminunc(@(t)(costFunction(t, X, y)), initial_theta, options);
17
18 plotDecisionBoundary(theta, X, y);
19 hold on;
20 xlabel('Exam 1 score')
21 ylabel('Exam 2 score')
22 legend('Admitted', 'Not admitted')
23 hold off;
24
25 prob = sigmoid([1 45 85] * theta);
26 fprintf('For a student with scores 45 and 85, we predict an admission probability\n\n', prob);
27
28
29 p = predict(theta, X);
30 fprintf('Train Accuracy: %f\n', mean(double(p == y)) * 100);

```

6.2.2 Second part

```

1 data = load('ex2data2.txt');
2 X = data(:, [1, 2]); y = data(:, 3);
3
4 plotData(X, y);
5 hold on;
6 xlabel('Microchip Test 1')
7 ylabel('Microchip Test 2')
8 legend('y = 1', 'y = 0')
9 hold off;
10
11 X = mapFeature(X(:, 1), X(:, 2));
12
13 initial_theta = zeros(size(X, 2), 1);
14 lambda = 1;
15 [cost, grad] = costFunctionReg(initial_theta, X, y, lambda);
16
17 initial_theta = zeros(size(X, 2), 1);
18 lambda = 0; % try: 0, 1, 10, 100
19 options = optimoptions(@fminunc, 'Algorithm', 'Quasi-Newton', 'GradObj', 'on',
20 'MaxIter', 1000);
21 [theta, J, exit_flag] = fminunc(@(t)(costFunctionReg(t, X, y, lambda)),
22 initial_theta, options);
23
24 plotDecisionBoundary(theta, X, y);
25 hold on;
26 title(sprintf('lambda = %g', lambda))
27 xlabel('Microchip Test 1')
28 ylabel('Microchip Test 2')
29 legend('y = 1', 'y = 0', 'Decision boundary')
30 hold off;
31
32 p = predict(theta, X);
33 fprintf('Train Accuracy: %f\n', mean(double(p == y)) * 100);

```

7 Neural Networks: Representation

7.1 Motivations

7.1.1 Non-linear Hypotheses

Why do we need neural networks, because we have linear and logistic regression? Let's say we have 100 parameters, then in order to include all the quadratic terms in the hypothesis, we will need about ≈ 5000 features. This is bad for two reasons. Firstly, it is expensive in memory and time, and secondly, most likely we will retrain. As a solution to the problem, we can leave only a subset, but then we reduce our capabilities.

In order to train the classifier, we need to show him a lot of photos with our object and without our object. To recognize photos, we need a nonlinear model, because the dependencies are quite complex. If we have a 50×50 image, that we need 2500 pixels, if we have RGB, then we have 7500 features, which means in total we have 3 million parameters. That's a lot.

7.1.2 Neurons and the Brain

One part of the brain can hear and see and feel, which means that there is some kind of specialized algorithm that can do everything. This is the motivation for neural networks. You can add various sensors and devices to a person that perform some function, then the brain will learn to perform it very quickly.

7.2 Neural Networks

7.2.1 Model Representation 1

Simplified neuron model (logistic unit): $(x_0 = 1), x_1, x_2, \dots, x_n \rightarrow \text{yellow circle} \rightarrow h_\theta(x)$. $x_0 = 1$ - "bias" unit. Hypotheses: $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$. θ - weights. $h_\theta(x)$ - sigmoid (logistic) activation function.

Neural Network

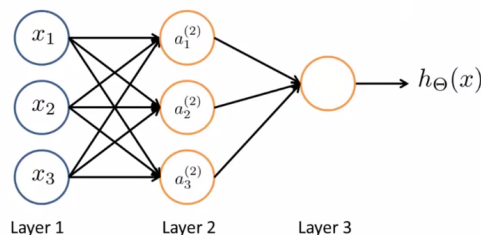


Figure 7: Neural Network Model

Layer 1 - input layer, Layer 3 - output layer, Layer 2 - hidden layer.

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

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

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

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

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

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

If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.

7.2.2 Model Representation 2

Forward propagation: Vectorized implementation

Designations: $a_1^{(2)} = g(z_1^{(2)})$, $a_2^{(2)} = g(z_2^{(2)})$, $a_3^{(2)} = g(z_3^{(2)})$.

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$z^{(2)} = \Theta^{(1)} x \quad \text{and} \quad a^{(2)} = g(z^{(2)})$$

If we put $a^{(1)} = x$, then: $z^{(2)} = \Theta^{(1)} a^{(1)}$.

Add $a_0^{(2)} = 1$, $z^{(3)} = \Theta^{(2)} a^{(2)}$, $h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$.

7.3 Applications

7.3.1 Multiclass Classification

$$h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad \text{etc.}$$

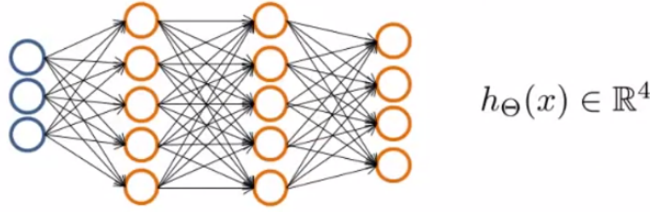


Figure 8: Multiclass Classification Problem

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

7.4 Implementation

7.4.1 All Functions

```

1 % Function #1
2 function [h, display_array] = displayData(X, example_width)
3 if ~exist('example_width', 'var') || isempty(example_width)
4     example_width = round(sqrt(size(X, 2)));
5 end

```

```

6 colormap(gray);
7 [m, n] = size(X);
8 example_height = (n / example_width);
9 display_rows = floor(sqrt(m));
10 display_cols = ceil(m / display_rows);
11 pad = 1;
12 display_array = - ones(pad + display_rows * (example_height + pad), ...
13     pad + display_cols * (example_width + pad));
14 curr_ex = 1;
15 for j = 1:display_rows
16     for i = 1:display_cols
17         if curr_ex > m, break;
18         end
19         max_val = max(abs(X(curr_ex, :)));
20         display_array(pad + (j - 1) * (example_height + pad) + (1:example_height), ...
21             pad + (i - 1) * (example_width + pad) + (1:example_width)) = ...
22             reshape(X(curr_ex, :), example_height, example_width) / max_val;
23         curr_ex = curr_ex + 1;
24     end
25     if curr_ex > m, break;
26     end
27 end
28 h = imagesc(display_array, [-1 1]);
29 axis image off
30 drawnow;
31 end
32
33 % Function #2
34 function [J, grad] = lrCostFunction(theta, X, y, lambda)
35 m = length(y);
36 h = sigmoid(X * theta);
37 J = -1 / m * (y' * log(h) + (1 - y)' * log(1 - h)) + lambda / (2 * m) * ...
38     (theta' * theta - theta(1) ^ 2);
39 grad = X' * (h - y) / m;
40 shift = lambda / m * theta;
41 grad(2:end) = grad(2:end) + shift(2:end);
42 end
43
44 % Function #3
45 function [all_theta] = oneVsAll(X, y, num_labels, lambda)
46 m = size(X, 1);
47 n = size(X, 2);
48 all_theta = zeros(num_labels, n + 1);
49 X = [ones(m, 1) X];
50 for c=1:num_labels
51     initial_theta = zeros(n + 1, 1);
52     options = optimset('GradObj', 'on', 'MaxIter', 50);
53     [theta] = fmincg(@(t)(lrCostFunction(t, X, (y == c), lambda)), ...
54         initial_theta, options);
55     all_theta(c, :) = theta';
56 end
57 end
58
59 % Function #4
60 function p = predict(Theta1, Theta2, X)
61 m = size(X, 1);
62 p = zeros(size(X, 1), 1);
63 X = [ones(m, 1) X];
64 for i=1:m
65     [~, pos] = max(sigmoid(Theta2 * [1; sigmoid(Theta1 * X(i, :))']));
66     p(i) = pos;
67 end
68 end
69
70 % Function #5

```



```

71 function p = predictOneVsAll(all_theta, X)
72 m = size(X, 1);
73 p = zeros(size(X, 1), 1);
74 X = [ones(m, 1) X];
75 for i=1:m
76     [~, pos] = max(all_theta * X(i, :)');
77     p(i) = pos;
78 end
79 end

```

7.4.2 Launch

```

1 load('ex3data1.mat');
2
3 m = size(X, 1);
4 rand_indices = randperm(m);
5 sel = X(rand_indices(1:100), :);
6 displayData(sel);
7
8 theta_t = [-2; -1; 1; 2];
9 X_t = [ones(5,1) reshape(1:15,5,3)/10];
10 y_t = ([1;0;1;0;1] >= 0.5);
11 lambda_t = 3;
12 [J, grad] = lrCostFunction(theta_t, X_t, y_t, lambda_t);
13 fprintf('Cost: %f | Expected cost: 2.534819\n', J);
14 fprintf('Gradients:\n'); fprintf('%f\n', grad);
15 fprintf('Expected gradients:\n 0.146561\n -0.548558\n 0.724722\n 1.398003');
16
17 num_labels = 10; % 10 labels, from 1 to 10
18 lambda = 0.1;
19 [all_theta] = oneVsAll(X, y, num_labels, lambda);
20
21 pred = predictOneVsAll(all_theta, X);
22 fprintf('\nTraining Set Accuracy: %f\n', mean(double(pred == y)) * 100);
23
24
25 load('ex3data1.mat');
26 m = size(X, 1);
27
28 sel = randperm(size(X, 1));
29 sel = sel(1:100);
30 displayData(X(sel, :));
31
32 load('ex3weights.mat');
33 pred = predict(Theta1, Theta2, X);
34 fprintf('\nTraining Set Accuracy: %f\n', mean(double(pred == y)) * 100);
35
36 rp = randi(m);
37 pred = predict(Theta1, Theta2, X(rp,:));
38 fprintf('\nNeural Network Prediction: %d (digit %d)\n', pred, mod(pred, 10));
39
40 displayData(X(rp, :));

```

8 Neural Networks: Learning

8.1 Cost Function and Backpropagation

8.1.1 Cost Function

Designations:

$$\begin{aligned} & \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\} \\ & L = \text{total num. of layers in network} \\ & s_l = \text{num. of units (not counting bias unit) in layer } l \\ & K = \text{number of classes} \end{aligned}$$

Cost function:

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

$$(h_{\Theta}(x))_i = i^{th} \text{ output}$$

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

8.1.2 Backpropagation Algorithm

Gradient computation:

We need to compute $\mathbb{J}(\Theta)$ and $\frac{\partial}{\partial \Theta_{ij}^{(l)}} \mathbb{J}(\Theta)$.

Given one training example (x, y) :

Forward propagation:

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

Gradient computation: Backpropagation algorithm

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

For each output unit (layer L): $\delta_j^{(L)} = a_j^{(L)} - y_j \Leftrightarrow \delta^{(L)} = a^{(L)} - y$.

$$\delta^{(l)} = (\Theta^{(l)})^T \delta^{(l+1)} \cdot * g'(z^{(l)}) = (\Theta^{(l)})^T \delta^{(l+1)} \cdot * a^{(l)} \cdot * (1 - a^{(l)})$$

For the first level, we do not consider the error function!

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} \mathbb{J}(\Theta) = a_j^{(l)} \delta_i^{(l+1)} \text{ (ignoring regularization)}$$

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 :

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 $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$
5. $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \iff \Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

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

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

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

8.1.3 Backpropagation Intuition

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

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

(Think of $\text{cost}(i) \approx (h_{\Theta}(x^{(i)}) - y^{(i)})^2$)

Formally (for $j \geq 0$):

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$$

Example:

$$\begin{aligned} \delta_2^{(2)} &= \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \delta_2^{(3)} \\ \delta_2^{(3)} &= \Theta_{12}^{(3)} \delta_1^{(4)} \end{aligned}$$

8.2 Backpropagation in Practice

8.2.1 Implementation Note: Unrolling Parameters

Advanced optimization

```
function [jVal, gradient] = costFunction(theta)
...
optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network (L=4):

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

"Unroll" into vectors

Example:

$$\begin{aligned} s_1 &= 10, s_2 = 10, s_3 = 1 \\ \Theta^{(1)} &\in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11} \\ D^{(1)} &\in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11} \end{aligned}$$

```
thetaVec = [Theta1(:); 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);
```

Learning Algorithm

Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
 Unroll to get initialTheta to pass to fminunc.

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

8.2.2 Gradient Checking

Numerical estimation of gradients

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

Parameter vector θ

$$\begin{aligned} \theta &\in \mathbb{R}^n \text{ ("unrolled")} \\ \theta &= \theta_1, \theta_2, \dots, \theta_n \\ \frac{\partial}{\partial \theta_1} \mathbb{J}(\theta) &\approx \frac{\mathbb{J}(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - \mathbb{J}(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon} \\ \frac{\partial}{\partial \theta_2} \mathbb{J}(\theta) &\approx \frac{\mathbb{J}(\theta_1, \theta_2 + \epsilon, \dots, \theta_n) - \mathbb{J}(\theta_1, \theta_2 - \epsilon, \dots, \theta_n)}{2\epsilon} \\ &\dots \\ \frac{\partial}{\partial \theta_n} \mathbb{J}(\theta) &\approx \frac{\mathbb{J}(\theta_1, \theta_2, \dots, \theta_n + \epsilon) - \mathbb{J}(\theta_1, \theta_2, \dots, \theta_n - \epsilon)}{2\epsilon} \end{aligned}$$

Check that $\text{gradApprox} \approx \text{DVec}$. Once you have verified once that your backpropagation algorithm is correct, you don't need to compute gradApprox again. The code to compute gradApprox can be very slow.

8.2.3 Random Initialization

If we initialize $\Theta_{ij}^{(l)} = 0$, then there will be a lot of identical values in our neural network that change by the same amount. After each update, parameters corresponding to inputs going into each of two hidden units are identical.

Random initialization: Symmetry breaking

Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$ (i.e $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)
 Theta1 = rand(10, 11) * (2 * initEpsilon) - initEpsilon
 Theta2 = rand(1, 11) * (2 * initEpsilon) - initEpsilon
 rand - between 0 and 1

8.2.4 Putting It Together

If there are several hidden layers in our architecture, then it is desirable that they are the same size. This size should be comparable to the number of input and output elements. $\mathbb{J}(\Theta)$ - non-convex function.

One effective strategy for choosing e_{init} is to base it on the number of units in the network. A good choice of e_{init} is $e_{\text{init}} = \frac{\sqrt{6}}{\sqrt{L_{\text{in}} + L_{\text{out}}}}$ where $L_{\text{in}} = s_l$ and $L_{\text{out}} = s_{l+1}$ are the number of units in the layers adjacent to $\Theta^{(l)}$. This range of values ensures that the parameters are kept small and makes the learning more efficient.

8.3 Implementation

8.3.1 Functions

```
1 % Function #1: Numerical Gradient
2 function numgrad = computeNumericalGradient(J, theta)
3 numgrad = zeros(size(theta));
4 perturb = zeros(size(theta));
5 e = 1e-4;
6 for p = 1:numel(theta)
7     perturb(p) = e;
8     loss1 = J(theta - perturb);
9     loss2 = J(theta + perturb);
10    numgrad(p) = (loss2 - loss1) / (2*e);
11    perturb(p) = 0;
12 end
13 end
14
15 % Function #2: Regularized Cost Function and Gradient
16 function [J, grad] = nnCostFunction(nn_params, ...
17                                     input_layer_size, ...
18                                     hidden_layer_size, ...
19                                     num_labels, ...
20                                     X, y, lambda)
21 Theta1 = reshape(nn_params(1:hidden_layer_size * (input_layer_size + 1)), ...
22                  hidden_layer_size, (input_layer_size + 1));
23
24 Theta2 = reshape(nn_params((1 + (hidden_layer_size * (input_layer_size + 1))) : end), ...
25                  num_labels, (hidden_layer_size + 1));
26 m = size(X, 1);
27 J = 0;
28 Theta1_grad = zeros(size(Theta1));
29 Theta2_grad = zeros(size(Theta2));
30 for i=1:m
31     a1 = [1 X(i, :)]';
32     a2 = [1; sigmoid(Theta1 * a1)];
33     h = sigmoid(Theta2 * a2);
34     vy = size(num_labels, 1);
35     for j=1:num_labels
36         vy(j, 1) = j;
37     end
38     vy = (vy == y(i, 1));
39     delta_last = h - vy;
40     delta_hidden = Theta2' * delta_last .* a2 .* (1 - a2);
41     delta_hidden = delta_hidden(2:end);
42     Theta1_grad = Theta1_grad + delta_hidden * (a1)';
43     Theta2_grad = Theta2_grad + delta_last * (a2)';
44     J = J + sum(-1 * vy .* log(h) - (1 - vy) .* log(1 - h));
45 end
46 Theta1_grad = Theta1_grad / m;
47 Theta2_grad = Theta2_grad / m;
```

```

48 for i=1:size(Theta1_grad, 1)
49     for j=2:size(Theta1_grad, 2)
50         Theta1_grad(i, j) = Theta1_grad(i, j) + Theta1(i, j) * lambda / m;
51     end
52 end
53 for i=1:size(Theta2_grad, 1)
54     for j=2:size(Theta2_grad, 2)
55         Theta2_grad(i, j) = Theta2_grad(i, j) + Theta2(i, j) * lambda / m;
56     end
57 end
58 J = J / m;
59 reg = 0;
60 for i=1:size(Theta1, 1)
61     for j=2:size(Theta1, 2)
62         reg = reg + Theta1(i, j) ^ 2;
63     end
64 end
65 for i=1:size(Theta2, 1)
66     for j=2:size(Theta2, 2)
67         reg = reg + Theta2(i, j) ^ 2;
68     end
69 end
70 J = J + lambda / (2 * m) * reg;
71 grad = [Theta1_grad(:) ; Theta2_grad(:)];
72 end
73
74 % Function #3: Random Initial Weights
75 function W = randInitializeWeights(L_in, L_out)
76     epsilon_init = 0.12;
77     W = rand(L_out, 1 + L_in) * 2 * epsilon_init - epsilon_init;
78 end
79
80 % Function #4: Sigmoid Derivate
81 function g = sigmoidGradient(z)
82     g = sigmoid(z) .* (1 - sigmoid(z));
83 end

```

8.3.2 Launching

```

1 % Part #1: Display Data
2 load('ex4data1.mat');
3 m = size(X, 1);
4 sel = randperm(size(X, 1));
5 sel = sel(1:100);
6 displayData(X(sel, :));
7
8 % Part #2: Load Weights
9 load('ex4weights.mat');
10
11 % Part #3: Unregularized Cost Function
12 input_layer_size = 400;
13 hidden_layer_size = 25;
14 num_labels = 10;
15 nn_params = [Theta1(:) ; Theta2(:)];
16 lambda = 0;
17 J = nnCostFunction(nn_params, input_layer_size, hidden_layer_size, num_labels,
18 X, y, lambda);
19 fprintf('Cost at parameters (loaded from ex4weights): %f', J);
20
21 % Part #4: Regularized Cost Function
22 lambda = 1;
23 J = nnCostFunction(nn_params, input_layer_size, hidden_layer_size, num_labels,

```

```

24 X, y, lambda);
25 fprintf('Cost at parameters (loaded from ex4weights): %f', J);
26
27 % Part #5: Random Initialization
28 initial_Theta1 = randInitializeWeights(input_layer_size, hidden_layer_size);
29 initial_Theta2 = randInitializeWeights(hidden_layer_size, num_labels);
30 initial_nn_params = [initial_Theta1(:) ; initial_Theta2(:)];
31
32 % Part #6: Learning
33 options = optimset('MaxIter', 500);
34 lambda = 1;
35 costFunction = @(p) nnCostFunction(p, input_layer_size, hidden_layer_size,
36 num_labels, X, y, lambda);
37 [nn_params, ~] = fmincg(costFunction, initial_nn_params, options);
38 Theta1 = reshape(nn_params(1:hidden_layer_size * (input_layer_size + 1)),
39 hidden_layer_size, (input_layer_size + 1));
40 Theta2 = reshape(nn_params((1 + (hidden_layer_size * (input_layer_size + 1))):end),
41 num_labels, (hidden_layer_size + 1));
42 pred = predict(Theta1, Theta2, X);
43 fprintf('\nTraining Set Accuracy: %f\n', mean(double(pred == y)) * 100);

```

9 Advice for Applying Machine Learning

9.1 Evaluating a Learning Algorithm

9.1.1 Deciding What to Try Next

Debugging a learning algorithm

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- Get more training examples
- Try smaller sets of features (avoid overfitting)
- Try getting additional features
- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc.)
- Try decreasing λ
- Try increasing λ

Machine learning diagnostic:

Diagnostic: A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance. Diagnostic can take time to implement, but doing so can be a very good use of your time.

9.1.2 Evaluating a Hypotheses

You can take our dataset and divide it into training and test parts, in a ratio of 70% to 30%.

- m_{test} = no. of test examples
- $(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$ - one test

Before separating the data, they need to be randomly mixed.

Training/testing procedure

- Learn parameter θ from training data
- Compute test set error $\mathbb{J}_{\text{test}}(\theta)$
- Misclassification error (0/1 misclassification error)

$$\text{err}(h_{\Theta}(x), y) = \begin{cases} 1 & \text{if } h_{\Theta}(x) \geq 0.5 \text{ and } y = 0 \text{ or } h_{\Theta}(x) < 0.5 \text{ and } y = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test Error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\Theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

9.1.3 Model Selection and Train/Validation/Test Sets

Model selection - this is the process of selecting features, the regularization parameter, the degree of the polynomial, etc. Let's denote d = degree of polynomial. Now let's train our models for different d and calculate our error function on them, choose the degree that has the least error. Despite the fact that we are looking at the error on the test set, the model can still generalize poorly.

Solution. Let's divide our dataset into 3 parts: Training set, **Cross Validation set (CV)**, Test set. Now we will choose our degree based on validation data, and evaluate the generalization error based on test data. Ratio: 60% : 20% : 20%.

9.2 Bias vs. Variance

9.2.1 Diagnosing Bias vs. Variance

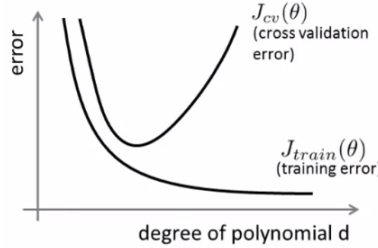


Figure 9: Bias vs. Variance

Is it a bias problem or a variance problem?

Bias (underfit, leftward): train and cv errors are high, and they are approximately equal.

Variance (overfit, on the right): train error is low, but in contrast, cv error is high.

The **optimal** value is in the **middle**.

9.2.2 Regularization and Bias/Variance

Large λ , high bias: $\lambda = 10000, \theta_1 \approx 0, \theta_2 \approx 0, \dots, h_\theta(x) \approx \theta_0$.

Intermediate λ : "Just right".

Small λ , high variance: $\lambda = 0$.

Choosing the **regularization** parameter λ **automatically**:

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$
$$\mathbb{J}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$
$$\mathbb{J}_{train,cv,test}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Try $\lambda = 0, 0.01, 0.02, 0.04, 0.08, \dots, 10$.

9.2.3 Learning Curves

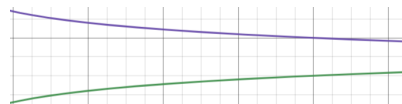


Figure 10: Learning Curves

Let's build such a graph - the dependence of the error on the number of examples. On top - an error on validation data, on the bottom - on test examples. If the curves converge, the difference between the straight lines is small, but the error is big - it's a lack of training. If the difference between the straight lines is large, there is a small error on the training data, and a large one on the validation data, then this is retraining. If you add a lot of test cases, then the lack of training will not be corrected, but retraining will become much better. Similarly with neural networks: few layers - under-training, many - retraining. Also, 2 is computationally more difficult.

9.3 Implementation

9.3.1 Functions

```
1 % #1
2 function [J, grad] = linearRegCostFunction(X, y, theta, lambda)
3 m = length(y);
4
5 J = (X * theta - y)' * (X * theta - y) / (2 * m) + ...
6     lambda / (2 * m) * (theta' * theta - theta(1, 1) ^ 2);
7
8 grad = zeros(size(theta));
9 for j=1:size(theta, 1)
10     sum = 0;
11     for i=1:m
12         sum = sum + (X(i, :) * theta - y(i, 1)) * X(i, j);
13     end
14     if j >= 2
15         sum = sum + lambda * theta(j, 1);
16     end
17     grad(j, 1) = sum / m;
18 end
19 grad = grad(:);
20 end
21
22 % #2
23 function [error_train, error_val] = ...
24     learningCurve(X, y, Xval, yval, lambda)
25 m = size(X, 1);
26 error_train = zeros(m, 1);
27 error_val = zeros(m, 1);
28 for i=1:m
29     [theta] = trainLinearReg(X(1:i,:), y(1:i), lambda);
30     error_train(i) = linearRegCostFunction(X(1:i,:), y(1:i), theta, 0);
31     error_val(i) = linearRegCostFunction(Xval, yval, theta, 0);
32 end
33 end
34
35 % #3
36 function [X_poly] = polyFeatures(X, p)
37 X_poly = zeros(numel(X), p);
38 for i=1:p
39     X_poly(:, i) = X(:, 1) .^ i;
40 end
41 end
42
43 % #4
44 function [lambda_vec, error_train, error_val] = ...
45     validationCurve(X, y, Xval, yval)
46 lambda_vec = [0 0.001 0.003 0.01 0.03 0.1 0.3 1 3 10]';
47 error_train = zeros(length(lambda_vec), 1);
48 error_val = zeros(length(lambda_vec), 1);
49 for i=1:length(lambda_vec)
50     [theta] = trainLinearReg(X, y, lambda_vec(i, 1));
51     error_train(i) = linearRegCostFunction(X, y, theta, 0);
52     error_val(i) = linearRegCostFunction(Xval, yval, theta, 0);
53 end
54 end
```

9.3.2 Code

```

1  load ('ex5data1.mat');
2  m = size(X, 1);
3  figure;
4  plot(X, y, 'rx', 'MarkerSize', 10, 'LineWidth', 1.5);
5  xlabel('Change in water level (x)');
6  ylabel('Water flowing out of the dam (y)');
7
8  theta = [1 ; 1];
9  J = linearRegCostFunction([ones(m, 1) X], y, theta, 1);
10 fprintf('Cost at theta = [1 ; 1]: %f', J);
11
12 [J, grad] = linearRegCostFunction([ones(m, 1) X], y, theta, 1);
13 fprintf('Gradient at theta = [1 ; 1]: [%f; %f] \n', grad(1), grad(2));
14
15 lambda = 0;
16 [theta] = trainLinearReg([ones(m, 1) X], y, lambda);
17
18 figure;
19 plot(X, y, 'rx', 'MarkerSize', 10, 'LineWidth', 1.5);
20 xlabel('Change in water level (x)');
21 ylabel('Water flowing out of the dam (y)');
22 hold on;
23 plot(X, [ones(m, 1) X]*theta, '--', 'LineWidth', 2)
24 hold off;
25
26 lambda = 0;
27 [error_train, error_val] = learningCurve([ones(m, 1) X], y,
28 [ones(size(Xval, 1), 1) Xval], yval, lambda);
29
30 plot(1:m, error_train, 1:m, error_val);
31 title('Learning curve for linear regression')
32 legend('Train', 'Cross Validation')
33 xlabel('Number of training examples')
34 ylabel('Error')
35 axis([0 13 0 150])
36
37 p = 8;
38 X_poly = polyFeatures(X, p);
39 [X_poly, mu, sigma] = featureNormalize(X_poly);
40 X_poly = [ones(m, 1), X_poly];
41
42 X_poly_test = polyFeatures(Xtest, p);
43 X_poly_test = X_poly_test - mu;
44 X_poly_test = X_poly_test ./ sigma;
45 X_poly_test = [ones(size(X_poly_test, 1), 1), X_poly_test];
46
47 X_poly_val = polyFeatures(Xval, p);
48 X_poly_val = X_poly_val - mu;
49 X_poly_val = X_poly_val ./ sigma;
50 X_poly_val = [ones(size(X_poly_val, 1), 1), X_poly_val];
51
52 lambda = 0;
53 [theta] = trainLinearReg(X_poly, y, lambda);
54
55 plot(X, y, 'rx', 'MarkerSize', 10, 'LineWidth', 1.5);
56 plotFit(min(X), max(X), mu, sigma, theta, p);
57 xlabel('Change in water level (x)');
58 ylabel('Water flowing out of the dam (y)');
59 title(sprintf('Polynomial Regression Fit (lambda = %f)', lambda));
60
61 [error_train, error_val] = learningCurve(X_poly, y, X_poly_val, yval, lambda);
62 plot(1:m, error_train, 1:m, error_val);
63 title(sprintf('Polynomial Regression Learning Curve (lambda = %f)', lambda));
64 xlabel('Number of training examples')

```

```

65  ylabel('Error')
66  axis([0 13 0 100])
67  legend('Train', 'Cross Validation')
68
69  [lambda_vec, error_train, error_val] = validationCurve(X_poly, y, X_poly_val,
70  yval);
71  plot(lambda_vec, error_train, lambda_vec, error_val);
72  legend('Train', 'Cross Validation');
73  xlabel('lambda');
74  ylabel('Error');

```

10 Machine Learning System Design

10.1 Building a Spam Classifier

Let's create a list of the most common words in our dataset. (10000-50000 words) We will put 1 if the word occurs and 0 if not. So we made up x and y. It is recommended to write a simple quick solution, and then analyze the result with your eyes and hands, this can help you decide what to do next. Do I need to distinguish between large and small letters in words, do I need to highlight the common root of words? It is very difficult to answer right away, so there is a need for one number that will rate the algorithm.

10.2 Handling Skewed Data

Let's say our test system detects 99 percent of cases of the disease. That is, it has 1% errors. At the same time, only 0.5% of all people actually have the disease. Then our system is actually very bad. Therefore, a different error metric is needed.

Precision/Recall

$y = 1$ in presence of rare class that we want to detect

Actual class | Predicted class

1 | 1 - True Positive

1 | 0 - False Negative

0 | 1 - False Positive

0 | 0 - True Negative

Precision

(Of all patients where we predicted $y = 1$, what fraction actually has cancer?)

$$\text{Precision} = \frac{\text{True pos.}}{\text{\#predicted pos.}} = \frac{\text{True pos.}}{\text{True pos.} + \text{False pos.}}$$

Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\text{Recall} = \frac{\text{True pos.}}{\text{\#actual pos.}} = \frac{\text{True pos.}}{\text{True pos.} + \text{False neg.}}$$

Good classifier: Recall and Precision $\rightarrow 1$.

How to estimate the algorithm?

Average: $\frac{P+R}{2}$.

F_1 Score(F score): $2 \cdot \frac{PR}{P+R}$.

Accuracy = $\frac{\text{True pos.} + \text{True neg.}}{\text{\#total number}}$.

11 Support Vector Machines

11.1 Large Margin Classification

11.1.1 Optimization Objective

Alternative view of logistic regression

In logistic regression, we had some kind of curve function for the cost function. Let's simplify it a bit and divide it into two parts: an inclined straight line and a horizontal straight line - a polyline. If $y = 1$, then we denote it $cost_1(z)$, otherwise $cost_0(z)$.

$$\min(\theta) : C \cdot \sum_{i=1}^m \left[y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$h_{\theta}(x) = \begin{cases} 1, & \text{if } \theta^T x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

11.1.2 Large Margin Intuition

If $y = 1$, we want $\theta^T x \geq 1$.

If $y = 0$, we want $\theta^T x \leq -1$.

SVM Decision Boundary: Linearly separable case

Other algorithms will separate the two clusters with a line so that the clusters are on different sides of the line. The SVM should divide so that the clusters are at the same distance from the line, that is, the straight line will be better circumscribed. Therefore, this algorithm is sometimes called "Large margin classifier".

Large margin classifier **in presence of outliers**

If we have two well-defined clusters and one outlier, then our algorithm can be very sensitive to this outlier. It doesn't seem right. If C is large, then the algorithm is sensitive, if C is small, then the border will most likely not change. $C \approx \frac{1}{\lambda}$.

11.1.3 Mathematics

Vector Inner Product

$$u^T v = v^T u = p \cdot \|u\| = u_1 v_1 + u_2 v_2$$

p = length of projection of v onto u

SVM Decision Boundary

Simplification: $\theta_0 = 0, n = 2$.

$$\min \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \left(\sqrt{\theta_1^2 + \theta_2^2} \right)^2 = \frac{1}{2} (\|\theta\|)^2$$

$$\theta^T x^{(i)} = p^{(i)} \cdot \|\theta\| = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Consider the **first case**. Let our boundary be drawn poorly and practically touches the clusters, then it can be shown that the θ vector will be perpendicular to our boundary, which means that the projections of the tangent clusters on the perpendicular will be very small, then in order for the equation to be fulfilled $p^{(i)} \cdot \|\theta\| \geq 1$, a very large theta is needed, which contradicts the minimized function.

Consider the **second case**. We have drawn a good boundary, and similarly we can show that the θ vector is perpendicular to the boundary. Then the projections are quite large, and therefore θ is quite small, so we will choose exactly such a boundary.

11.2 Kernels

11.2.1 Kernels 1

$$\begin{aligned}\Theta_0 + \Theta_1 f_1 + \Theta_2 f_2 + \Theta_3 f_3 + \dots \\ f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, \dots\end{aligned}$$

Given x , compute new feature depending on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$.
Given x :

$$\begin{aligned}f_1 &= \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) \\ f_2 &= \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right) \\ f_3 &= \text{similarity}(x, l^{(3)}) = \exp\left(-\frac{\|x - l^{(3)}\|^2}{2\sigma^2}\right) \\ \text{similarity} &\Leftrightarrow k \\ f_1 &= \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)\end{aligned}$$

k - Kernel (Gaussian Kernel) function

If $x \approx l^{(1)}$: $f_1 \approx \exp(0) \approx 1$.

If x is far from $l^{(1)}$: $f_1 \approx 0$. If the sigma value is small ($\sigma^2 = 0.5$), then our graph will be narrow,

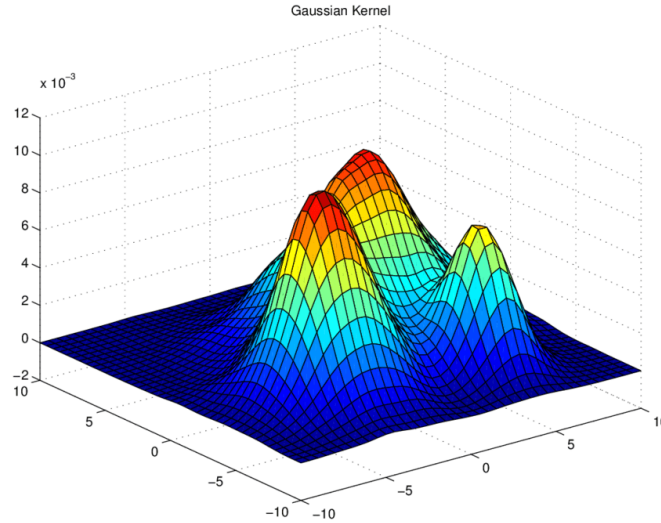


Figure 11: Kernel function

sharp descents. If it is bigger ($\sigma^2 = 3$), then the graph will become wider, gentle descents. At the same time, the maximum of 1 will remain. We will predict "1" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

11.2.2 Kernels 2

Given $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$, choose $l^{(1)} = x^{(1)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$f_i = \text{similarity}(x, l^{(i)})$$

$$f = \begin{bmatrix} f_0 = 1 \\ f_1 \\ \dots \\ f_m \end{bmatrix}$$

$$x^{(i)} \in \mathbb{R}^{n+1} \Rightarrow f^{(i)} \in \mathbb{R}^{m+1} \text{ and } \Theta \in \mathbb{R}^{m+1}$$

Predict "y=1" if $\theta^T f \geq 0$

Training:

$$\min(\theta) : C \cdot \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

If the number of examples is very large (≈ 10000), then in real conditions this formula is used for regularization $\Theta^T \cdot M \cdot \Theta$ instead of this one $\Theta^T \cdot \Theta$. Support vector machine and kernels - very fast and good idea which optimizes well. But logistic regression and kernels is bad idea because it will be computationally expensive. How to optimize this cost function? It's not important, because there are already a lot of perfect software with some kinds of tricks.

SVM parameters:

$C (= \frac{1}{\lambda})$: Large C - lower bias, high variance. Small C - higher bias, low variance.

σ^2 : Large σ^2 - features f_i vary more smoothly. Higher bias, lower variance. Small σ^2 - features f_i vary less smoothly. Lower bias, higher variance (derivatives).

11.3 SVMs in Practice

11.3.1 Using An SVM

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters θ .

Need to specify:

- Choice of parameter C.
- Choice of kernel (similarity function):
 - E.g. No kernel ("linear kernel"). Predict "y=1" if $\theta^T x \geq 0$. n - large, m - small.
 - Gaussian kernel. Need to choose σ^2 . n - small, m - large. Complex decision boundary.

Kernel (similarity) functions:

function $f = \text{kernel}(x1, x2)$

$$f = \exp \left(- \frac{\|x1 - x2\|^2}{2\sigma^2} \right)$$

return

Note: Do perform feature scaling before using the Gaussian kernel.

Other choices of kernel:

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels. (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages's optimizations run correctly, and do not diverge).

Many off-the-shelf kernels available:

- Polynomial kernel: $k(x, l) = (x^T l + \text{constant})^{\text{degree}}$.
- More esoteric: String, chi-square kernels, ...

If n is large (relative to m) - use logistic regression, or SVM without a kernel (linear).

If n is small, m is intermediate - use SVM with Gaussian kernel.

If n is small, m is large - create/add more features, then use logistic regression or SVM without a kernel.

Neural network likely to work well for most of these settings, but may be slower to train.