

Week 2 Lecture Notes

ML: Linear Regression with Multiple Variables

Linear regression with multiple variables is also known as "multivariate linear regression".

We now introduce notation for equations where we can have any number of input variables.

$$\begin{aligned} x^{(i)} &= \text{value of feature } j \text{ in the } i^{\text{th}} \text{ training example} \\ x^{(i)} &= \text{the column vector of all the feature inputs of the } i^{\text{th}} \text{ training example} \\ m &= \text{the number of training examples} \\ n &= \lfloor x^{(1)} \rfloor, \text{ (the number of features)} \end{aligned}$$

Now define the multivariable form of the hypothesis function as follows, accommodating these multiple features:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

In order to develop intuition about this function, we can think about θ_0 as the basic price of a house, θ_1 as the price per square meter, θ_2 as the price per floor, etc. x_1 will be the number of square meters in the house, x_2 the number of floors, etc.

Using the definition of matrix multiplication, our multivariable hypothesis function can be concisely represented as:

$$h_{\theta}(x) = \begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

This is a vectorization of our hypothesis function for one training example; see the lessons on vectorization to learn more.

Remark: Note that for convenience reasons in this course Mr. Ng assumes $x_0^{(i)} = 1$ for $(i \in 1, \dots, m)$

[Note: So that we can do matrix operations with theta and x, we will set $x_0^{(i)} = 1$, for all values of i. This makes the two vectors 'theta' and $x_{(i)}$ match each other element-wise (that is, have the same number of elements: $n+1$).]

The training examples are stored in X row-wise, like such:

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} \\ x_0^{(2)} & x_1^{(2)} \\ x_0^{(3)} & x_1^{(3)} \\ \vdots & \vdots \end{bmatrix}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

You can calculate the hypothesis as a column vector of size $(m \times 1)$ with:

$$h_{\theta}(X) = X\theta$$

For the rest of these notes, and other lecture notes, X will represent a matrix of training examples $x_{(i)}$ stored row-wise.

Cost function

For the parameter vector θ (of type \mathbb{R}^{n+1} or in $\mathbb{R}^{(n+1) \times 1}$), the cost function is:

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

The vectorized version is:

$$J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y)$$

Where \vec{y} denotes the vector of all y values.

Gradient Descent for Multiple Variables

The gradient descent equation itself is generally the same form; we just have to repeat it for our 'n' features:

$$\begin{aligned} \text{repeat until convergence: } \{ \\ \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)} \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)} \\ \theta_2 &:= \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)} \\ &\dots \\ \} \end{aligned}$$

In other words:

$$\begin{aligned} \text{repeat until convergence: } \{ \\ \theta_j &:= \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad \text{for } j := 0..n \\ \} \end{aligned}$$

Matrix Notation

The Gradient Descent rule can be expressed as:

$$\theta := \theta - \alpha \nabla J(\theta)$$

Where $\nabla J(\theta)$ is a column vector of the form:

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

The j-th component of the gradient is the summation of the product of two terms:

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \\ &= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \cdot (h_{\theta}(x^{(i)}) - y^{(i)}) \end{aligned}$$

Sometimes, the summation of the product of two terms can be expressed as the product of two vectors.

Here, $x_j^{(i)}$, for $i = 1, \dots, m$, represents the m elements of the j-th column, \vec{x}_j , of the training set X.

The other term $(h_{\theta}(x^{(i)}) - y^{(i)})$ is the vector of the deviations between the predictions $h_{\theta}(x^{(i)})$ and the true values $y^{(i)}$. Re-writing $\frac{\partial J(\theta)}{\partial \theta_j}$, we have:

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{1}{m} \vec{x}_j^T (X\theta - \vec{y}) \\ \nabla J(\theta) &= \frac{1}{m} X^T (X\theta - \vec{y}) \end{aligned}$$

Finally, the matrix notation (vectorized) of the Gradient Descent rule is:

$$\theta := \theta - \frac{\alpha}{m} X^T (X\theta - \vec{y})$$

Feature Normalization

We can speed up gradient descent by having each of our input values in roughly the same range. This is because θ will descend quickly on small ranges and slowly on large ranges, and so will oscillate inefficiently down to the optimum when the variables are very uneven.

The way to prevent this is to modify the ranges of our input variables so that they are all roughly the same. Ideally:

$$-1 \leq x_{(i)} \leq 1$$

or

$$-0.5 \leq x_{(i)} \leq 0.5$$

These aren't exact requirements; we are only trying to speed things up. The goal is to get all input variables into roughly one of these ranges, give or take a few.

Two techniques to help with this are **feature scaling** and **mean normalization**. Feature scaling involves dividing the input values by the range (i.e. the maximum value minus the minimum value) of the input variable, resulting in a new range of just 1. Mean normalization involves subtracting the average value for an input variable from the values for that input variable, resulting in a new average value for the input variable of just zero. To implement both of these techniques, adjust your input values as shown in this formula:

$$x_i := \frac{x_i - \mu_i}{s_i}$$

Where μ_i is the **average** of all the values for feature (i) and s_i is the range of values (max - min), or s_i is the standard deviation.

Note that dividing by the range, or dividing by the standard deviation, give different results. The quizzes in this course use range - the programming exercises use standard deviation.

Example: x_i is housing prices with range of 100 to 2000, with a mean value of 1000. Then, $x_i := \frac{\text{price} - 1000}{1900}$.

Quiz question #1 on Feature Normalization (Week 2, Linear Regression with Multiple Variables)

Your answer should be rounded to exactly two decimal places. Use a '.' for the decimal point, not a ','. The tricky part of this question is figuring out which feature of which training example you are asked to normalize. Note that the mobile app doesn't allow entering a negative number (an 0.0), so you will need to use a browser to submit this quiz if your solution requires a negative number.

Gradient Descent Tips

Debugging gradient descent. Make a plot with *number of iterations* on the x-axis. Now plot the cost function, $J(\theta)$ over the number of iterations of gradient descent. If $J(\theta)$ ever increases, then you probably need to decrease α .

Automatic convergence test. Declare convergence if $J(\theta)$ decreases by less than ϵ in one iteration, where ϵ is some small value such as 10^{-3} . However in practice it's difficult to choose this threshold value.

It has been proven that if learning rate α is sufficiently small, then $J(\theta)$ will decrease on every iteration. Andrew Ng recommends decreasing α by multiples of 3.

Features and Polynomial Regression

We can improve our features and the form of our hypothesis function in a couple different ways.

We can **combine** multiple features into one. For example, we can combine x_1 and x_2 into a new feature x_3 by taking $x_1 \cdot x_2$.

Polynomial Regression

Our hypothesis function **will not** be linear (a straight line) if that does not fit the data well.

We can **change the behavior or curve** of our hypothesis function by making it a quadratic, cubic or square root function (or any other form).

For example, if our hypothesis function is $h_{\theta}(x) = \theta_0 + \theta_1 x_1$ then we can create additional features based on x_1 , to get the quadratic function $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$ or the cubic function $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

In the cubic version, we have created new features x_2 and x_3 where $x_2 = x_1^2$ and $x_3 = x_1^3$.

To make it a square root function, we could do: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

Note that at 2:52 and through 6:22 in the "Features and Polynomial Regression" video, the curve that Prof Ng discusses about "doesn't ever come back down" is in reference to the hypothesis function that uses the sqrt() function (shown by the solid purple line), not the one that uses size^2 (shown with the dotted blue line). The quadratic form of the hypothesis function would have the shape shown by the blue dotted line if θ_2 was negative.

One important thing to keep in mind is, if you choose your features this way then feature scaling becomes very important.

eg. if x_1 has range 1 - 1000 then range of x_1^2 becomes 1 - 1000000 and that of x_1^3 becomes 1 - 1000000000.

Normal Equation

The "Normal Equation" is a method of finding the optimum theta **without iteration**.

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

There is **no need** to do feature scaling with the normal equation.

Mathematical proof of the Normal equation requires knowledge of linear algebra and is fairly involved, so you do not need to worry about the details.

Proofs are available at these links for those who are interested:

[https://en.wikipedia.org/wiki/Linear_least_squares_\(mathematics\)](https://en.wikipedia.org/wiki/Linear_least_squares_(mathematics))

<http://eli.thegreenplace.net/2014/derivation-of-the-normal-equation-for-linear-regression>

The following is a comparison of gradient descent and the normal equation:

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$, need to calculate inverse of $X^T X$
Works well when n is large	Slow if n is very large

With the normal equation, computing the inversion has complexity $O(n^3)$. So if we have a very large number of features, the normal equation will be slow. In practice, when n exceeds 10,000 it might be a good time to go from a normal solution to an iterative process.

Normal Equation Noninvertibility

When implementing the normal equation in octave we want to use the 'pinv' function rather than 'inv'.

$X^T X$ may be **noninvertible**. The common causes are:

- Redundant features, where two features are very closely related (i.e. they are linearly dependent)
- Too many features (e.g. $m \leq n$). In this case, delete some features or use "regularization" (to be explained in a later lesson).

Solutions to the above problems include deleting a feature that is linearly dependent with another or deleting one or more features when there are too many features.

ML: Octave Tutorial

Basic Operations

```

1 %%% Change Octave prompt
2 PS1('>>> ');
3 %%% Change working directory in windows example:
4 cd 'c:/path/to/desired/directory name'
5 %%% Note that it uses normal slashes and does not use escape characters for the
   empty spaces.
6
7 %%% elementary operations
8 5+6
9 3-2
10 5*8
11 1/2
12 2^6
13 1 == 2 % False
14 1 == 2 % true, note, not "!="
15 1 && 0
16 1 || 0
17 xor(1,0)
18
19 %%% variable assignment
20 a = 3; % semicolon suppresses output
21 b = 'hi';
22 c = 3+1;
23
24 %%% Displaying then:
25 a = pi
26 disp(a)
27 disp(sprintf('2 decimals: %0.2f', a))
28 disp(sprintf('6 decimals: %0.6f', a))
29 format long
30 format short
31 a
32 a
33
34 %%% vectors and matrices
35 A = [1 2 3 4; 5 6]
36
37 v = [1 2 3]
38
39 v = 1:0.1:2 % from 1 to 2, with stepsize of 0.1. Useful for plot axes
40 v = 1:6 % from 1 to 6, assumes stepsize of 1 (row vector)
41
42 C = 2*ones(2,3) % same as C = [2 2 2; 2 2 2]
43 w = ones(1,3) % 1x3 vector of ones
44
45 w = rand(1,3) % drawn from a uniform distribution
46 w = randn(1,3) % drawn from a normal distribution (mean=0, var=1)
47 w = -5 + sqrt(10)*(randn(1,10000)); % (mean = -5, var = 10) - note: add the
   semicolon
48 hist(v,50) % plot histogram using 50 bins (default)
49 hist(v,50) % plot histogram using 50 bins
50 % Note: if hist() crashes, try "graphics_toolkit('gnu_plot')"
51
52 I = eye(4) % 4x4 Identity matrix
53
54 % help function
55 help eye
56 help rand
57 help help
58
59
60
```

Moving Data Around

Data files used in this section: [featuresX.dat](#), [priceY.dat](#)

```

1 %%% dimensions
2 sz = size(A) % 1x2 matrix: [(number of rows) (number of columns)]
3 size(A,1) % number of rows
4 size(A,2) % number of cols
5 length(v) % size of longest dimension
6
7 %%% loading data
8 pwd % show current directory (current path)
9 cd 'c:/Users/ang/Octave Files' % change directory
10 ls % list files in current directory
11 load gly.dat % alternatively, load('gly.dat')
12 load gly.dat
13
14 %%% list variables in workspace (detailed view)
15 whos % clear command without any args clears all vars
16 v = q1x(1:10); % first 10 elements of q1x (counts down the columns)
17 save hello.mat v % save variable v into file hello.mat
18 save hello.txt v -ascii; % save as ascii
19 % fopen, fread, fprintf, fscanf also work [[not needed in class]]
20
21 %%% indexing
22 A(1,2) % indexing ls (row,col)
23 A(2,:) % get the 2nd row.
24 A(:,2) % get the 2nd col.
25 A(1:3,:) % print all the elements of rows 1 and 3
26
27 A(:,2) = [10; 11; 12] % change second column
28 A = [A; 100; 101; 102]; % append column vec
29 A(:,3) % Select all elements as a column vector.
30
31 %%% Putting data together
32 A = [1 2 3 4; 5 6]
33 B = [11 12; 13 14; 15 16] % same dims as A
34 C = [A; B] % concatenating A and B matrices side by side
35 C = [A; B] % concatenating A and B matrices side by side
36 C = [A; B] % Concatenating A and B top and bottom
37
38
39
40
```

Computing on Data

```

1 %%% initialize variables
2 A = [1 2 3 4 5 6]
3 B = [11 12 13 14 15 16]
4 C = [1 1 2 2]
5 v = [1 2 3]
6
7 %%% matrix operations
8 A * C % matrix multiplication
9 A .* B % element-wise multiplication
10 % A ./ C or A ./ B gives error - wrong dimensions
11 A ./ 2 % element-wise square of each element in A
12 1./v % element-wise reciprocal
13 length(v) % functions like this operate element-wise on vecs or matrices
14 exp(v)
15 abs(v)
16
17 -v % -1*v
18
19 v + ones(length(v), 1)
20 % v + 1 % same
21
22 A' % matrix transpose
23
24 %%% misc useful functions
25 % max (or min)
26 a = [1 15 2 0.5]
27 val = max(a)
28 [val,idx] = max(a) % val - maximum element of the vector a and index - index
29 % value where maximum occur
30 val = max(A) % if A is matrix, returns max from each column
31
32 % compare values in a matrix & find
33 a = 3 % checks which values in a are less than 3
34 find(a < 3) % gives location of elements less than 3
35 A = magic(3) % generates a magic matrix - not much used in ML algorithms
36 [r,c] = find(A==7) % row, column indices for values matching comparison
37
38 % sum, prod
39 sum(a)
40 prod(a) % or cell()
41 floor(a)
42 max(rand(3),rand(3))
43 max(A,[],1) % maximum along columns (defaults to columns - max(A,[],))
44 max(A,[],2) % maximum along rows
45 A = magic(9)
46 sum(A,1)
47 sum(A,2)
48 sum(sum(A .* eye(9)))
49 sum(sum(A .* fliplr(eye(9))))
50
51 % Matrix inverse (pseudo-inverse)
52 pinv(A) % inv(A'*A)*A'
53
54
```

Plotting Data

```

1 %%% plotting
2 t = [0:0.01:0.98];
3 y1 = sin(2*pi*4*t);
4 plot(t,y1);
5 hold on; % "hold off" to turn off
6 hold on; % "hold off" to turn off
7 plot(t,y2,'r');
8 xlabel('time');
9 ylabel('value');
10 legend('sin','cos');
11 title('my plot');
12 print -dpng 'myPlot.png'
13 close; % or, "close all" to close all figs
14 figure(1); plot(t, y1);
15 figure(2); plot(t, y2);
16 figure(2); clf; % can specify the figure number
17 subplot(1,2,1); % Divide plot into 1x2 grid, access 1st element
18 plot(t,y1);
19 subplot(1,2,2); % Divide plot into 1x2 grid, access 2nd element
20 plot(t,y2);
21 axis([0.5 1 -1 1]); % change axis scale
22
23 %%% display a matrix (or image)
24 figure;
25 imagesc(magic(15)), colormap gray;
26 % comma-chaining function calls.
27 a=1;b=2;c=3;
28 a=1;b=2;c=3;
29
```

Control statements: for, while, if statements

```

1 v = zeros(10,1);
2 for i=1:10;
3     v(i) = 2^i;
4 end;
5 % Can also use "break" and "continue" inside for and while loops to control
   execution.
6
7 i = 1;
8 while i <= 5;
9     v(i) = 100;
10    i = i+1;
11 end
12
13 i = 1;
14 while true;
15     v(i) = 999;
16     i = i+1;
17     if i == 6;
18         break;
19     end;
20 end
21
22 if v(1)==1;
23     disp('The value is one!');
24 elseif v(1)==2;
25     disp('The value is two!');
26 else
27     disp('The value is not one or two!');
28 end
29
```

Functions

To create a function, type the function code in a text editor (e.g. gedit or notepad), and save the file as "functionName.m"

Example function:

```

1 function y = squareThisNumber(x)
2
3     y = x^2;
4
```

To call the function in Octave, do either:

1) Navigate to the directory of the functionName.m file and call the function:

```

1 %%% Navigate to directory:
2 cd /path/to/function
3
4 %%% Call the function:
5 functionName(args)
6
```

2) Add the directory of the function to the load path and save it: You should not use addpath/savepath for any of the assignments in this course. Instead use 'cd' to change the current working directory. Watch the video on submitting assignments in week 2 for instructions.

```

1 % To add the path for the current session of Octave:
2 addpath('/path/to/function')
3
4 % To remember the path for future sessions of Octave, after executing
   addpath above, also do:
5 savepath
6
```

Octave's functions can return more than one value:

```

1 function [y1, y2] = squareandCubeThisNo(x)
2
3     y1 = x^2
4     y2 = x^3
5
```

Call the above function this way:

```

1 [a,b] = squareandCubeThisNo(x)
2
```

Vectorization

Vectorization is the process of taking code that relies on **loops** and converting it into **matrix operations**. It is more efficient, more elegant, and more concise.

As an example, let's compute our prediction from a hypothesis. Theta is the vector of fields for the hypothesis and x is a vector of variables.

With loops:

```

1 prediction = 0.8;
2 for j = 1:n+1;
3     prediction = theta(j) * x(j);
4 end;
5
```

With vectorization:

```

1 prediction = 0.8;
2 for j = 1:n+1;
3     prediction = theta(j) * x(j);
4 end;
5
```