

IV. LINEAR REGRESSION WITH MULTIPLE VARIABLES

Lecture 19: Multiple Features

In the housing price example, our goal was to predict the price of a house as a function of its size. Of course, there are many more factors which could affect house prices:

multiple features (variables)

| Size x_1 | Number of bedrooms x_2 | Number of floors x_3 | Age of home (years) x_4 | Price (\$1000) y |
|---------------|--------------------------------|------------------------------|---------------------------------|-----------------------|
| 2104 | 5 | 1 | 45 | 460 |
| 1916 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| ... | ... | ... | ... | ... |

Notation:

$n = \text{number of features}$

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

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

Example above: $n=4$; $x_3^{(2)}=2$; $x^{(3)} = \begin{bmatrix} 1534 \\ 3 \\ 2 \\ 30 \end{bmatrix}$

Our hypothesis now is:

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

Notation:

$$\text{if } x_0 = 1 \Leftrightarrow x_0^{(i)} = 1$$

$$* \quad \underline{x}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$* \quad \underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

so that

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

Lecture 20: Gradient Descent for Multiple Variables

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

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[\sum_{j=1}^n \theta_j x_j^{(i)} - y^{(i)} \right]^2$$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{2m} \sum_{i=1}^m 2 \left(\sum_{j=1}^n \theta_j x_j^{(i)} - y^{(i)} \right) x_j^{(i)} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Note that we can write the cost function as:

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

where $X_{ij} = x_j^{(i)}$ & $y_i = y^{(i)}$

Also, gradient descent can be written as

Repeat until converge {

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

Of course, these only hold for linear regression.

Gradient Descent: repeat until convergence {

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

} (simultaneously update θ_j for $j=0, 1, \dots, n$)

Lecture 21: Gradient Descent in Practice I - Feature Scaling

→ Different features may have very different values.

Example) → Size of house ~ 1000 (feet²)
 → Number of rooms $\sim 1-5$

This causes the cost function J to be a lot more stretched in certain directions than others. For instance, in the example above, taking $x_1 = \text{size of house}$ & $x_2 = \text{number of units}$ & the hypothesis function:

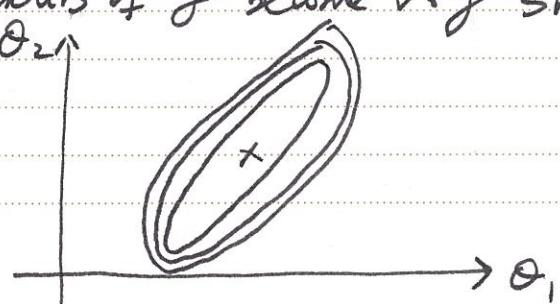
$$h_\theta(x) = \theta_1 x_1 + \theta_2 x_2$$

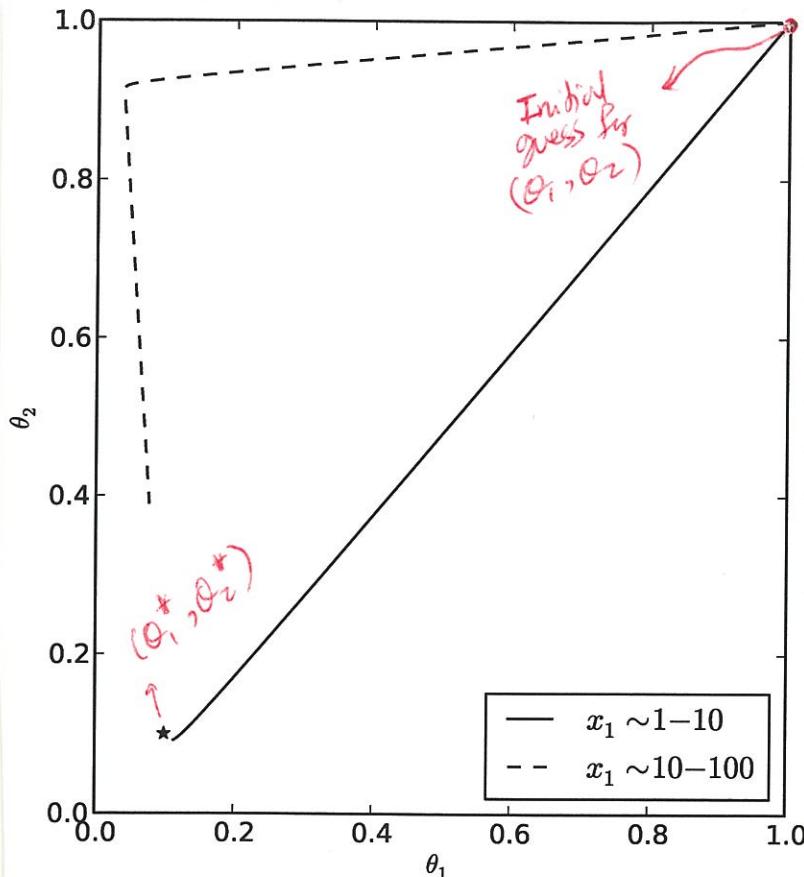
Then, an order of magnitude estimate of the cost function $J(\theta)$ would be:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} - y^{(i)})^2$$

$$\sim 10^6 \theta_1^2 + 10 \theta_2^2 + \dots$$

so that the contours of J become very skewed:





To see the effect this could have on gradient descent, I did the following exercise:

① Pick a target value θ_1^*, θ_2^* .

② Pick m random numbers uniformly for both $x_1 \sim x_2$:

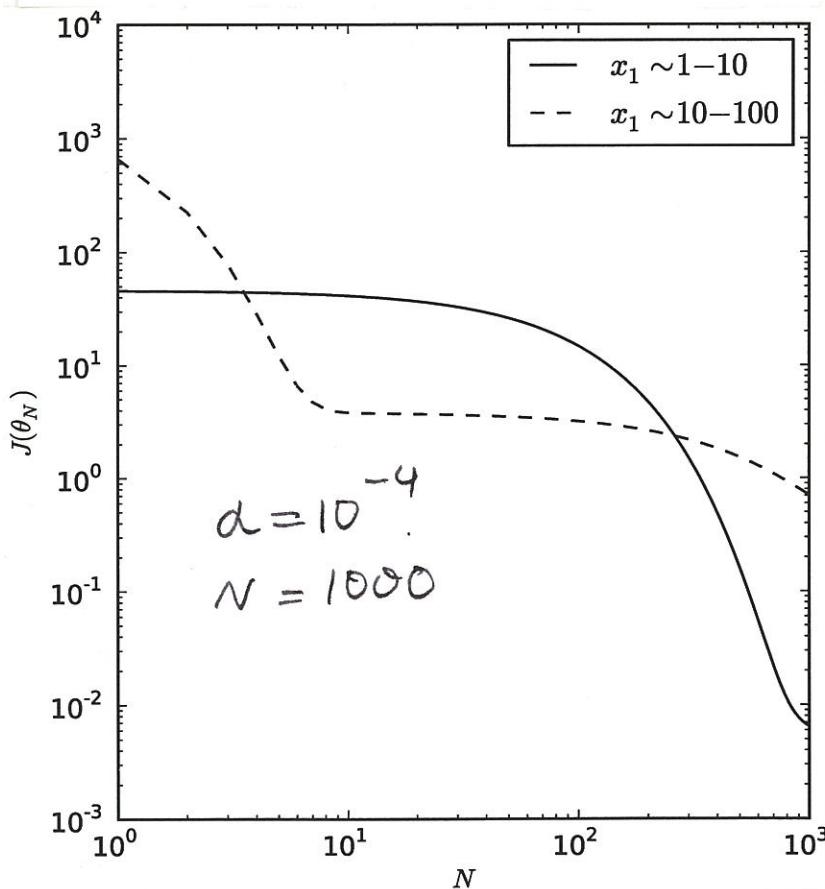
$$x_1 \sim 1-10$$

$$x_2 \sim 1-10$$

③ Construct $y^{(i)} = \theta_1^* x_1^{(i)} + \theta_2^* x_2^{(i)}$ & add a bit of noise to $y^{(i)}$

④ Run Gradient descent on the training set $(x_1, x_2; y)$ using the model $h_\theta(x) = \theta_1 x_1 + \theta_2 x_2$

⑤ Do exactly as above except pick x_1 from a range $x_1 \sim 10-100$.



Conclusion: Gradient Descent takes much longer to converge when the features are not of the same order!

To alleviate this problem, we do what is known as feature scaling:

→ Scale all features so that they are in the similar range (-1, 1).

→ Some ways of doing this:

$$x_j^{(i)} \rightarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad \text{where } \mu_j \text{ & } \sigma_j \text{ are the mean & standard deviation of the feature } j.$$

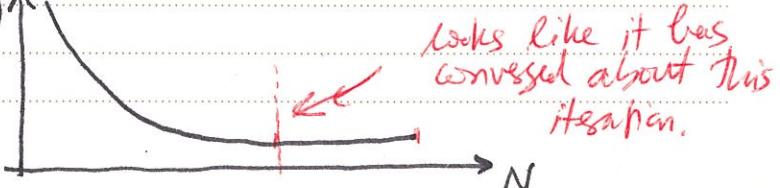
$$x_j^{(i)} \rightarrow \frac{x_j^{(i)} - \mu_j}{S_j} \quad \text{where } S_j \text{ is the range of values of feature } j, \text{ i.e. } \max\{x_j^{(i)}\} - \min\{x_j^{(i)}\}.$$

Lecture 22: Gradient Descent in Practice II - Learning rate

→ A good check of seeing if gradient descent is working correctly is to plot the value of the cost function as a function of the iterations in gradient descent. The bottom figure on the previous page is a plot of this sort. The important thing to note that J should decrease with every iteration. If this is not the case, something is wrong. (of course, provided the learning rate is small enough)

→ If J is increasing with the number of iterations, it usually means that the learning rate α is too large.

→ This plot also helps us to see if gradient descent has converged or not. $J(N) \uparrow$

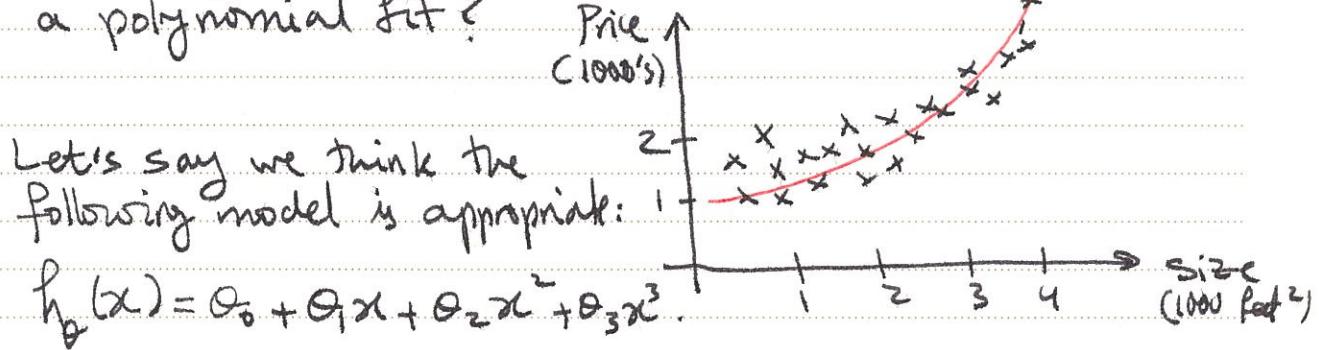


Choosing α → try running gradient descent with different values of α & plot $J(\theta_N)$ - N.

(Example) $\alpha = 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \dots$

Lecture 23: Features and Polynomial Regression

What if we have data which is better modeled as a polynomial fit?



Well, we can then easily use our linear regression techniques by simply letting:

$$\begin{cases} x_1 = x \\ x_2 = x^2 \\ x_3 = x^3 \end{cases}$$

As long as $h_{\theta}(x)$ is linear in θ , we can use linear regression! This is quite powerful.

⇒ Note that feature scaling becomes crucial in such cases. If we take x to be the size of the house in square-feet, then:

$$x_1 \sim 1 - 1000$$

$$x_2 \sim 1 - 10^6$$

$$x_3 \sim 1 - 10^9$$



How can we tell what the best fit to the data should be?
 Well, apparently there are algorithms which automatically find the best fit (e.g. if a 2nd degree polynomial is appropriate, a 3rd degree, a mixing of features, etc.)

↳ My guess is that such an algorithm would try to find a fit which has as low a $J(\theta_{fit})$ as possible, while minimizing the number of features introduced.

↳ Note that we can get pretty fancy: say we have a training set with input variables x_1, x_2 & output variable y . We can apply linear regression to:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots$$

$$\text{where } x_3 = \sqrt{x_1 x_2} ; x_4 = e^{x_1 - x_2} ; \dots$$

Lecture 24: Normal Equation

We can, of course, just minimize the cost function J analytically.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[\sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right]^2$$

$$\frac{\partial J}{\partial \theta_k} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right) x_k^{(i)} = 0$$

$$\sum_{j=0}^n \theta_j \left(\sum_{i=1}^m x_j^{(i)} x_k^{(i)} \right) = \sum_{i=1}^m y^{(i)} x_k^{(i)} \quad \#$$

→ called the
Design Matrix

$$\text{Let } X_{ij} = x_j^{(i)} \text{ & } Y_i = y^{(i)}.$$

$$\text{Then } \sum_{i=1}^m x_j^{(i)} x_K^{(i)} = \sum_{i=1}^m X_{ij} X_{ik} = [X^T X]_{kj}$$

$$\sum_{i=1}^m y^{(i)} x_K^{(i)} = \sum_{i=1}^m X_{ik} Y_i = [X^T Y]_k$$

Then # on previous page reads:

$$\sum_{j=0}^n [X^T X]_{kj} \theta_j = [X^T Y]_k$$

$$\rightarrow [X^T X] \Theta = X^T Y$$

$$\& \quad \boxed{\Theta = [X^T X]^{-1} X^T Y} \quad | \text{ This exact solution is known as the normal equation.}$$

| Example) | x_0 | size (feet ²) | # of bedrooms | Price (\$1000) |
|----------|-------|---------------------------|---------------|----------------|
| | x_1 | x_2 | x_3 | y |
| | 1 | 2104 | 5 | 460 |
| | 1 | 1416 | 3 | 232 |
| | 1 | 1534 | 3 | 315 |
| | 1 | 852 | 2 | 178 |

$$X = \begin{bmatrix} 1 & 2104 & 5 \\ 1 & 1416 & 3 \\ 1 & 1534 & 3 \\ 1 & 852 & 2 \end{bmatrix} \quad \& \quad Y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

Note that X is $m \times (n+1)$ dimensional

In the language we established on pg 17:

$$\mathbb{X} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(m)})^T \end{bmatrix} \quad \text{where } \mathbf{x}^{(i)} = \begin{bmatrix} 1 & \mathbf{x}_1^{(i)} \\ & \vdots \\ & \mathbf{x}_m^{(i)} \end{bmatrix}$$

2) Andrew claims that feature scaling is not necessary here. I can see why, but it would be nice if the elements of \mathbb{X} are of the same order of magnitude...

2) Gradient Descent vs. Normal Equation.

* Obviously, whenever there is an analytic solution, life becomes much easier. The only computational challenge in using the normal equation is inverting $(\mathbb{X}^T \mathbb{X})_{mn}$.

* Inversion is an order n^3 operation & for large n , gradient descent works much better.

* Rough rule: $\begin{cases} n \leq 10000 \rightarrow \text{Normal Equation} \\ n > 10000 \rightarrow \text{Gradient Descent.} \end{cases}$

Remarks: $n = \# \text{ of features!}$

* Indeed, gradient descent is used to solve a large system of linear equations of the form $A\mathbf{x} = b$, when inverting A directly becomes very expensive.

The optimal learning rate for linear regression

→ This is not discussed in the course, but when we're using a model linear in the fitting parameters θ , there's an optimal value of α .

→ Say we start at $\hat{\theta} = \begin{bmatrix} \hat{\theta}_0 \\ \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_n \end{bmatrix}$. The gradient of $J(\theta)$ at $\hat{\theta}$

is $\hat{dJ} = \begin{bmatrix} \frac{\partial J(\hat{\theta})}{\partial \theta_0} \\ \frac{\partial J(\hat{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\hat{\theta})}{\partial \theta_n} \end{bmatrix}$. Our next move would be to $\hat{\theta} - \alpha \hat{dJ}$. What is the optimal α ? Well, we can pick the α which would lead to the lowest value of J ! In other words, consider

$J(\hat{\theta} - \alpha \hat{dJ})$ & find the value of α which minimizes it: $\frac{d}{d\alpha} J(\hat{\theta} - \alpha \hat{dJ}) = 0$

$$J(\hat{\theta} - \alpha \hat{dJ}) = \frac{1}{2m} \sum_{i=1}^m [(\hat{\theta} - \alpha \hat{dJ})^T \mathbf{x}^{(i)} - y^{(i)}]^2$$

$$\frac{d}{d\alpha} J(\hat{\theta} - \alpha \hat{dJ}) = \frac{1}{m} \sum_{i=1}^m [(\hat{\theta} - \alpha \hat{dJ})^T \mathbf{x}^{(i)} - y^{(i)}] \times -\hat{dJ}^T \mathbf{x}^{(i)}$$

$$= 0$$

$$\sum_{i=1}^m (\hat{\theta}^T \mathbf{x}^{(i)}) (\hat{dJ}^T \mathbf{x}^{(i)}) = \sum_{j,k=1}^n \hat{\theta}_j \hat{dJ}_k \sum_{i=1}^m x_j^{(i)} x_k^{(i)}$$

$$= \hat{\theta}^T \mathbf{X}^T \mathbf{X} \hat{dJ} = (\mathbf{X} \hat{\theta})^T \mathbf{X} \hat{dJ}$$

$$\rightarrow (\mathbf{X} \hat{\theta})^T \mathbf{X} \hat{dJ} - \alpha (\mathbf{X} \hat{dJ})^T \mathbf{X} \hat{dJ} - \hat{dJ}^T \mathbf{X}^T \mathbf{y} = 0$$

$$\rightarrow \alpha = \frac{(\mathbf{X} \hat{dJ})^T [\mathbf{X} \hat{\theta} - \mathbf{y}]}{(\mathbf{X} \hat{dJ})^T (\mathbf{X} \hat{dJ})}$$

It just remains to find $\hat{\nabla} J$:

$$\begin{aligned}\frac{\partial J}{\partial \theta_k} &= \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_k^{(i)} \\ &= \frac{1}{m} [(x^T x) \theta_0 - x^T y]_k\end{aligned}$$

$$\hat{\nabla} J = \frac{1}{m} [x^T x \theta - x^T y]$$

$$\begin{aligned}\rightarrow \alpha &= \frac{1}{m} [x x^T x \hat{\theta} - x x^T y]^T [x \hat{\theta} - y] \\ &= \frac{1}{m^2} [x x^T x \hat{\theta} - x x^T y]^T [x x^T x \hat{\theta} - x x^T y] \\ &= \frac{m [x \hat{\theta} - y]^T x x^T [x \hat{\theta} - y]}{[x \hat{\theta} - y]^T x x^T x x^T [x \hat{\theta} - y]}\end{aligned}$$

$$\text{Let } M \equiv x x^T \quad \& \quad G(\theta) \equiv x \theta - y$$

$$\rightarrow \alpha = \frac{G(\theta)^T M G(\theta)}{G(\theta)^T M^2 G(\theta)}$$

Gradient Descent: repeat until converge {

$$\theta := \theta - \frac{G(\theta)^T M G(\theta)}{G(\theta)^T M^2 G(\theta)} x^T G(\theta)$$

}

Gradient Descent for Linear Regression.

Note that M only has to be computed once!

Lecture 25: Normal Equation Noninvertibility

Could it be that $A = X^T X$ is non-invertible?

Yes. This could have two reasons:

1. Redundant features

Say we have 2 feature x_1 & x_2 such that

$$x_2^{(i)} = \alpha x_1^{(i)}. \text{ Then: } A_{11} = \sum_{i=1}^m (x_1^{(i)})^2$$

$$A_{12} = \alpha A_{11} \text{ & } A_{22} = \alpha^2 A_{11}$$

Also assume we're not fitting for θ_0 . Then

$$A = \sum_{i=1}^m (x_i^{(i)})^2 \begin{bmatrix} 1 & \alpha \\ \alpha & \alpha^2 \end{bmatrix} \Rightarrow \det A = 0.$$

2. Too many features (i.e. $m \leq n$)

This means that we're trying to fit for too many parameters given a small training set. Of course, this would cause degeneracies which translate into saying that A is not invertible.

* Lectures 26 - 32 are about using Octave.