

Machine Learning - Week 2

赵燕

目录

1	Multivariate Linear Regression	2
1.1	Multiple Features	2
1.2	Gradient Descent for Multiple Variables	3
1.3	Gradient Descent in Practice I-Feature Scaling	4
1.4	Gradient Descent in Practice II-Learning Rate	5
1.5	Features and Polynomial Regression	7
2	Computer Parameters Analytically	10
2.1	Normal Equation	10
2.2	Normal Equation Noninvertibility	12

1 Multivariate Linear Regression

1.1 Multiple Features

根据之前的例子，我们对房价模型增加更多的特征，例如房间数楼层等，构成了一个含有多个变量的模型，模型中的特征值为 (x_1, x_2, \dots, x_n) ：

Multiple features (variables).

Size (feet ²) 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
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

Notation:

- $\rightarrow n$ = number of features $n = 4$
- $\rightarrow x^{(i)}$ = input (features) of i^{th} training example.
- $\rightarrow x_j^{(i)}$ = value of feature j in i^{th} training example.

Handwritten notes:

- $m = 47$ (number of training examples)
- $x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$ (feature vector for the 2nd example)
- $x_3^{(2)} = 2$ (value of the 3rd feature for the 2nd example)

图 1: 多变量房价模型

注释:

(1) n : 特征的数量

(2) $x^{(i)}$: 第 i 个训练样本，是特征矩阵中的第 i 行，是一个向量 (vector)；例如：上图中的

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \quad (1)$$

(3) $x_j^{(i)}$: 特征矩阵中第 i 行的第 j 个特征值，也就是第 i 个训练样本实例的第 j 个特征，如上图的 $x_3^{(2)} = 2$ 。

多变量的假设函数:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad (2)$$

这个公式中有你 $n+1$ 个参数和 n 个变量，为了能使公式简化，引入 $x_0=1$ ，则公式转化为:

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad (3)$$

此时模型中的参数是一个n+1维的变量，任何一个训练实例都是n+1维的向量，特征矩阵X的维度是m*(n+1)。公式可以简化为：

$$h_{\theta}(x) = [\theta_0 \quad \theta_1 \quad \cdot \quad \cdot \quad \cdot \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{bmatrix} = \theta^T X \quad (4)$$

$$h_{\theta}(x) = \theta^T X \quad (5)$$

→ $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$
 For convenience of notation, define $x_0 = 1$. ($x_0^{(i)} = 1$)

$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$ $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$

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

Diagram illustrating the matrix multiplication: θ^T (1×(n+1) matrix) multiplied by x ((n+1)×1 vector) results in $\theta^T x$.

图 2: 多变量假设函数

1.2 Gradient Descent for Multiple Variables

与单变量线性回归类似，在多变量线性回归中也构建一个代价函数，则这个代价函数是所有建模误差的平方和，即：

$$J(\theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (6)$$

其中：

$$h_{\theta}(x) = \theta^T X = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad (7)$$

多变量梯度算法定义：n=1和n_i=1的情况其实是一样的

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad (8)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \quad (9)$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \quad (10)$$

.....

}

In others words:

repeat until convergence: {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (11)$$

$$\text{for } j := 0, 1, 2, \dots, n \quad (12)$$

}

单变量与多变量的比较: The following image compares gradient descent with one variable to gradient descent with multiple variables:

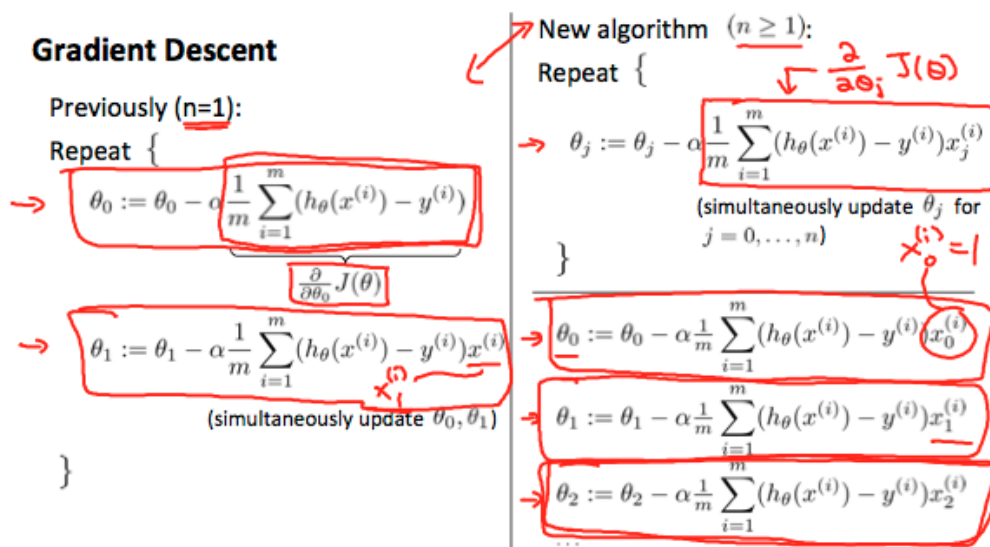


图 3: 单变量与多变量的梯度下降算法对比

1.3 Gradient Descent in Practice I-Feature Scaling

特征缩放可以将梯度下降的速度加快，将梯度下降收敛所需的循环次数减少。

以房价问题为例，假设我们使用两个特征,房屋的尺寸和房间的数量,尺寸的值 为 0-2000 平方英尺,而房间数量的值则是 0-5,以两个参数分别为横纵坐标,绘制代价函数的等高线图，如图所示，代价函数的等高线图很扁，又细又长，梯度下降算法需要非常多次的迭代才能收敛。

解决的额方法是尝试将所有特征的尺度都尽量缩放到-1到1之间。

Feature Scaling

Idea: Make sure features are on a similar scale.

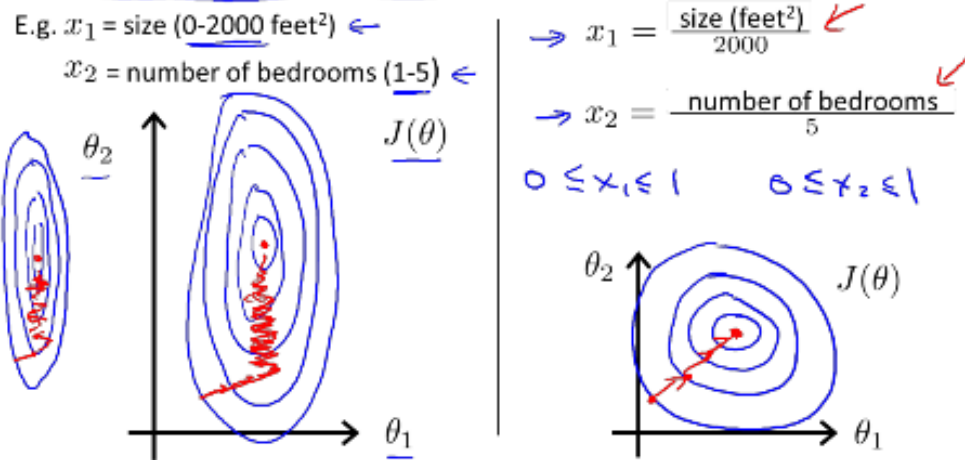


图 4: 多变量假设函数

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:

$$x_i = \frac{x_i - \mu_i}{s_i} \quad (13)$$

μ_i :is the average of all the values for feature (i);

s_i :is the range of values (max - min), or si is the standard deviation(标准差)

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.

For example: if x_i represents housing prices with a range of 100 to 2000 and a mean value of 1000, then:

$$x_i = \frac{\text{price} - 1000}{900} \quad (14)$$

1.4 Gradient Descent in Practice II-Learning Rate

梯度下降算法所需的迭代次数根据模型的不同而不同，我们不能提前预知，我们可以绘制迭代次数个代价函数的图表来观测算法在何时趋于收敛。

也有一些自动测试是否收敛的方法，例如将代价函数的变化值与某个阈值（例如0.001）进行比较，但是通常看如下的图表更好。

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.

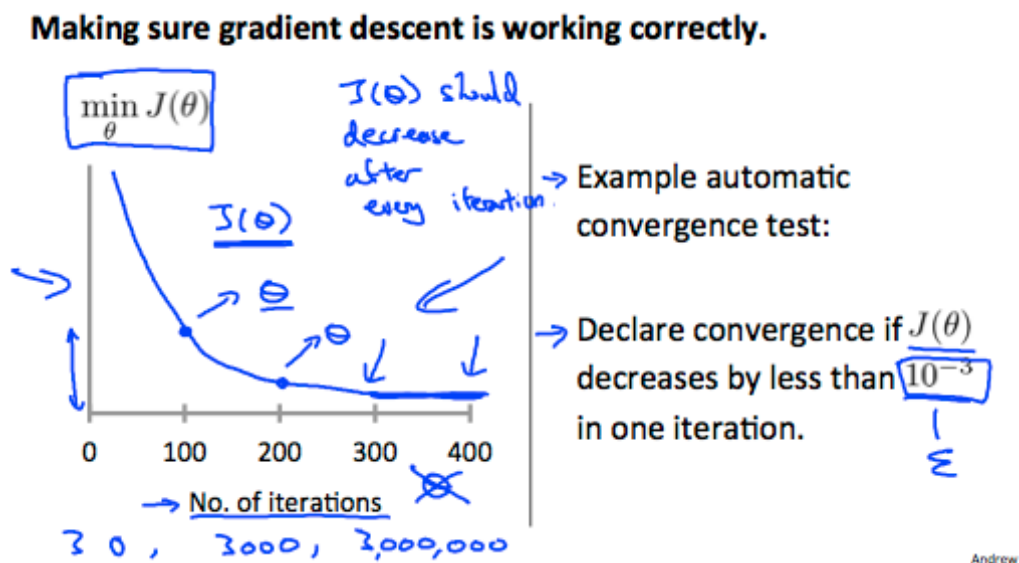


图 5: 迭代次数与代价函数

梯度下降算法的每次迭代受到学习率 α 的影响，如果学习率 α 过小，则达到收敛所需的迭代次数会非常高；如果学习率 α 过大，每次迭代可能不会减小代价函数，可能会越过局部最小值导致无法收敛。

It has been proven that if learning rate α is sufficiently small, then $J(\theta)$ will decrease on every iteration.

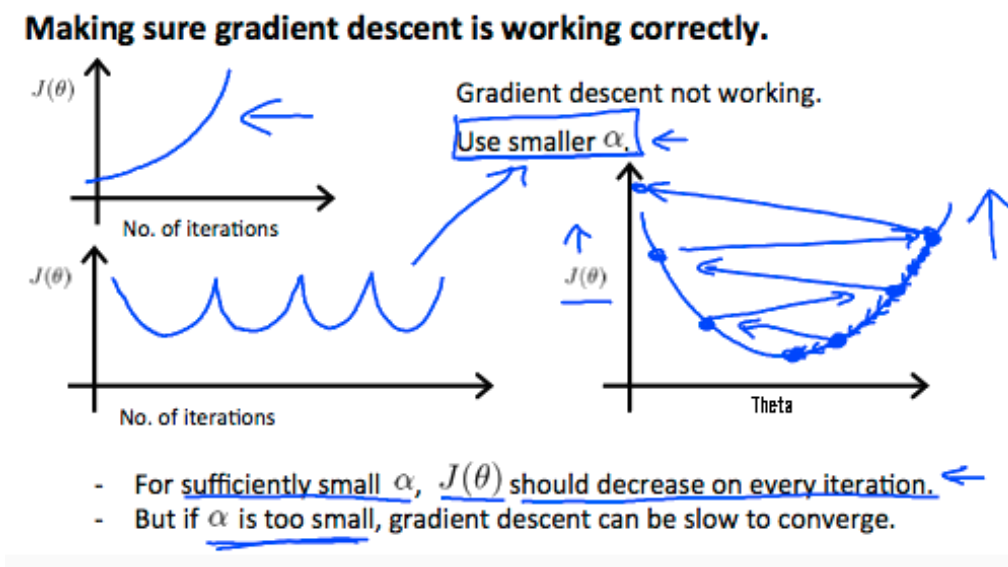


图 6: α 的影响

通常可以考虑尝试这些学习率 α :

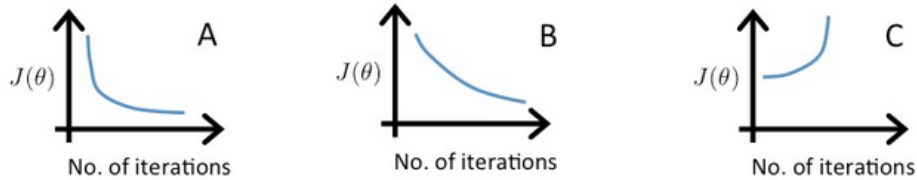
$\alpha = 0.01, 0.03, 0.1, 0.3, 1, 3, 10$

To summarize:

If α is too small: slow convergence.

If α is too large: may not decrease on every iteration and thus may not converge.

Suppose a friend ran gradient descent three times, with $\alpha = 0.01$, $\alpha = 0.1$, and $\alpha = 1$, and got the following three plots (labeled A, B, and C):



Which plots corresponds to which values of α ?

☐ A is $\alpha = 0.01$, B is $\alpha = 0.1$, C is $\alpha = 1$.

☒ A is $\alpha = 0.1$, B is $\alpha = 0.01$, C is $\alpha = 1$.

正确

In graph C, the cost function is increasing, so the learning rate is set too high. Both graphs A and B converge to an optimum of the cost function, but graph B does so very slowly, so its learning rate is set too low. Graph A lies between the two.

图 7: Learning Rate

1.5 Features and Polynomial Regression

举例：房价预测问题

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth} \quad (15)$$

注释：

x_1 :frontage(临街宽度)

x_2 :depth(纵向深度)

x :area=frontage*depth(面积) 则：

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad (16)$$

线性回归并不适用于所有数据，有时我们需要曲线来适应我们的数据，比如一个二次方模型：

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 \quad (17)$$

或者三次方模型：

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 \quad (18)$$

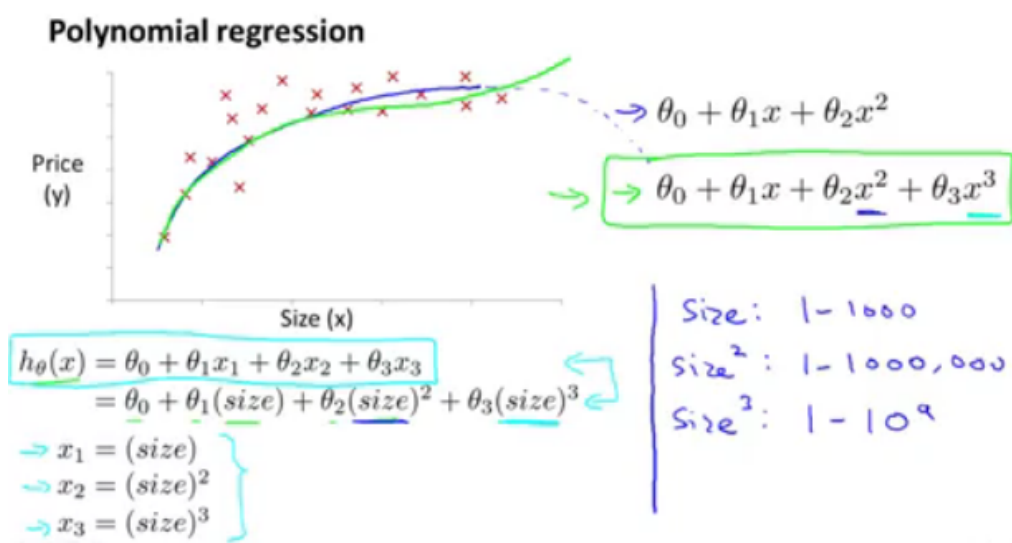


图 8: 房价预测: 多项式回归

通常我们需要先观察数据然后再决定准备尝试怎样的模型。另外我们可以另:

$$x_2 = x_2^2 \quad (19)$$

$$x_3 = x_3^3 \quad (20)$$

从而将模型转化为线性回归模型。

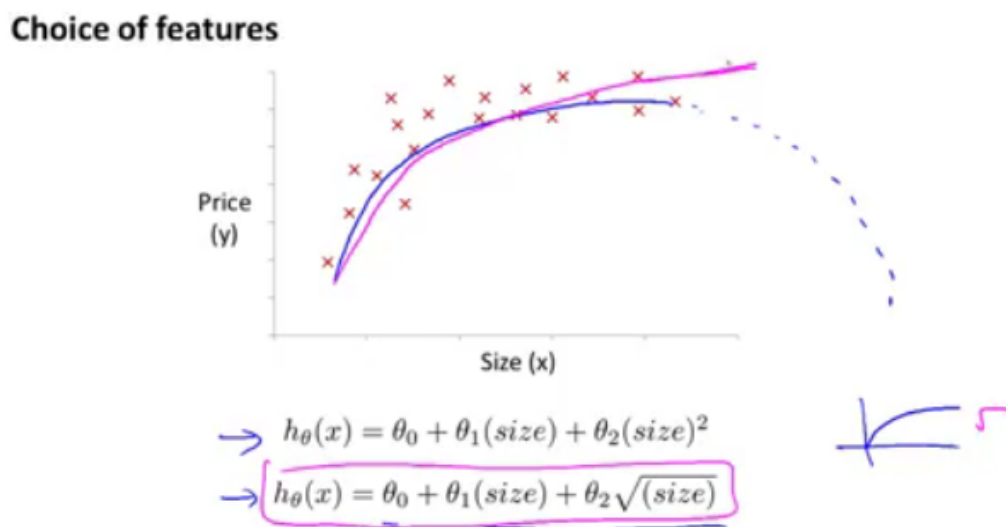


图 9: Choice of features

根据函数的图形特征, 我们可以使:

$$h_{\theta}(x) = \theta_0 + \theta_1 \text{size} + \theta_2 (\text{size})^2 \quad (21)$$

或者:

$$h_{\theta}(x) = \theta_0 + \theta_1 \text{size} + \theta_2 \sqrt{(\text{size})} \quad (22)$$

Polynomial Regression(多项式回归):

Our hypothesis function need 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

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 \quad (23)$$

then we can create additional features based on x_1 , to get the quadratic function

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 \quad (24)$$

or the cubic function

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 \quad (25)$$

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: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

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

注意: 如果我们采用多项式回归, 在运行梯度下降算法前, 特征缩放非常有必要。

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

多项式回归特征缩放练习题:

Suppose you want to predict a house's price as a function of its size. Your model is

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2 \sqrt{(\text{size})}.$$

Suppose size ranges from 1 to 1000 (feet²). You will implement this by fitting a model

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

Finally, suppose you want to use feature scaling (without mean normalization).

Which of the following choices for x_1 and x_2 should you use? (Note: $\sqrt{1000} \approx 32$.)

- ☐ $x_1 = \text{size}, x_2 = 32\sqrt{(\text{size})}$
- ☐ $x_1 = 32(\text{size}), x_2 = \sqrt{(\text{size})}$
- ☒ $x_1 = \frac{\text{size}}{1000}, x_2 = \frac{\sqrt{(\text{size})}}{32}$

正确

图 10: 特征缩放例题

2 Computer Parameters Analytically

2.1 Normal Equation

正规方程：针对某些问题，正规方程有更好的解决方案。

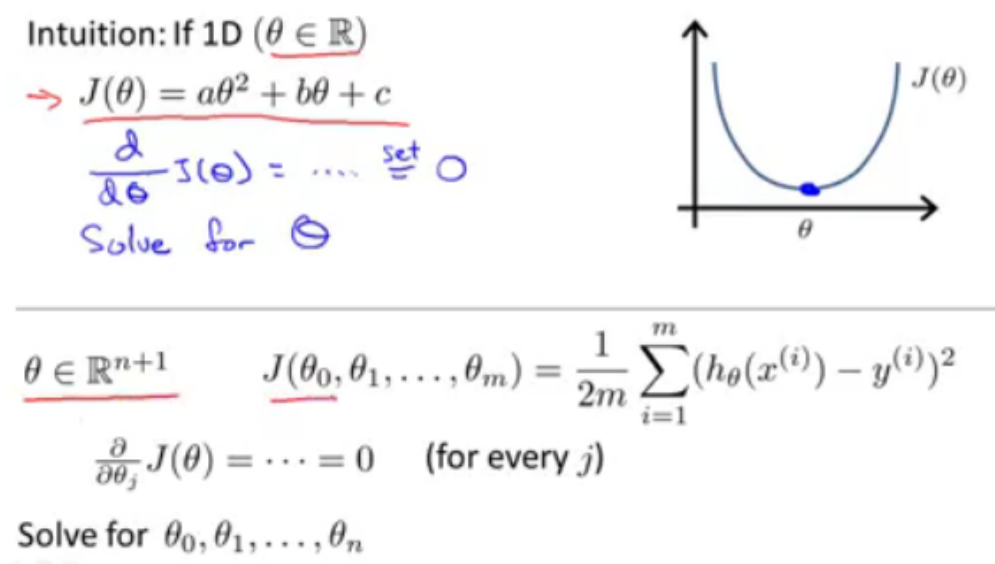


图 11: Intuition

正规方程通过求解下面的方程来找出使得代价函数最小的参数的 θ :

$$\frac{\partial}{\partial \theta_j} J(\theta_j) = 0 \quad (26)$$

假使我们的训练集特征为 X （包含了 $x_0=1$ ），并且我们的训练集结果为向量 y ，则利用正规方程解出向量

$$\theta = (X^T X)^{-1} X^T y \quad (27)$$

注释：上标 T 代表矩阵转置，上标 -1 代表矩阵的逆。设矩阵 $A = X^T X$ ，则： $(X^T X)^{-1} = A^{-1}$

以下数据为例：

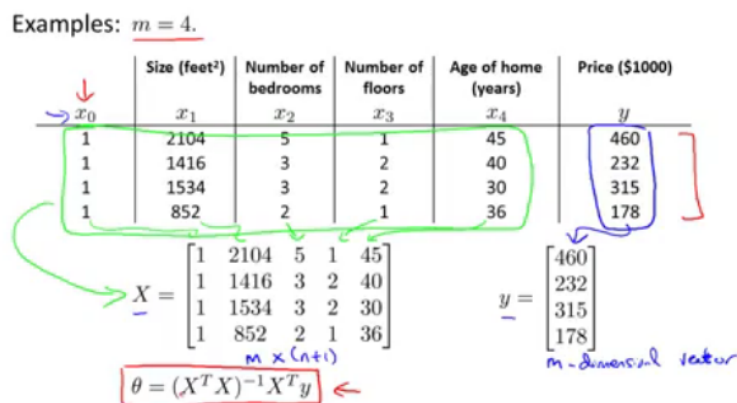


图 12: 房价问题：正规方程

X(0)	X(1)	X(2)	X(3)	X(4)	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

图 13: 数据表格

运用正规方程求解参数:

$$\left(\begin{bmatrix} 1 & 1 & 1 & 1 \\ 2104 & 1416 & 1534 & 852 \\ 5 & 3 & 3 & 2 \\ 1 & 2 & 2 & 1 \\ 45 & 40 & 30 & 36 \end{bmatrix} \times \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \right)^{-1} \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2104 & 1416 & 1534 & 852 \\ 5 & 3 & 3 & 2 \\ 1 & 2 & 2 & 1 \\ 45 & 40 & 30 & 36 \end{bmatrix} \times \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

图 14: 运用正规方程求解参数

在Octave中, 正规方程

$$\theta = (X^T X)^{-1} X^T y \quad (28)$$

写作:

$$\text{pinv}(X' * X) * X' * y \quad (29)$$

注意: 对于那些不可逆的矩阵(通常是因为特征之间不独立,如同时包含英尺为单位的尺寸和米为单位的尺寸两个特征,也有可能是特征数量大于训练集的数量),正规方程方法是不能用的。

梯度下降与正规方程的比较:

梯度下降	正规方程
需要选择学习率 α	不需要
需要多次迭代	一次运算得出
当特征数量 n 大时也能较好适用	需要计算 $(X^T X)^{-1}$ 如果特征数量 n 较大则运算代价大, 因为矩阵逆的计算时间复杂度为 $O(n^3)$, 通常来说当 n 小于 10000 时还是可以接受的
适用于各种类型的模型	只适用于线性模型, 不适合逻辑回归模型等其他模型

图 15: 梯度下降与正规方程的比较

总结: 只要特征变量的数目并不大, 标准方程是一个很好的计算参数 θ 的替代方法。具体地说, 只要特征变量数量小于一万, 我通常使用标准方程法, 而不使用梯度下降法。

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

图 16: Gradient Descent or Normal Equation

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.

2.2 Normal Equation Noninvertibility

When implementing the normal equation in octave we want to use the 'pinv' function rather than 'inv.' The 'pinv' function will give you a value of θ even if $X^T X$ is not invertible.

If $X^T X$ is noninvertible, the common causes might be having :

(1) Redundant features, where two features are very closely related (i.e. they are linearly dependent)

(2) Too many features (e.g. $m \leq n$). In this case, delete some features or use "regularization" (to be explained in a later lesson).

What if $X^T X$ is non-invertible?

- Redundant features (linearly dependent).

E.g. $x_1 = \text{size in feet}^2$

~~$x_2 = \text{size in m}^2$~~

$$x_1 = (3.28)^2 x_2$$

$$1 \text{ m} = 3.28 \text{ feet}$$

$$\rightarrow m = 10 \leftarrow$$

$$\rightarrow n = 100 \leftarrow$$

$$\theta \in \mathbb{R}^{101}$$

- Too many features (e.g. $m \leq n$).

- Delete some features, or use regularization.

↓ later

图 17: Non-invertible

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.