

Lecture02

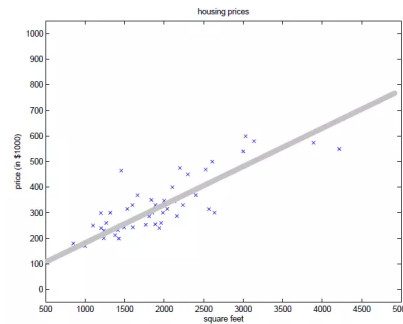
大纲 (outline)

- Linear regression
- Batch/stochastic gradient descent
- Normal equation

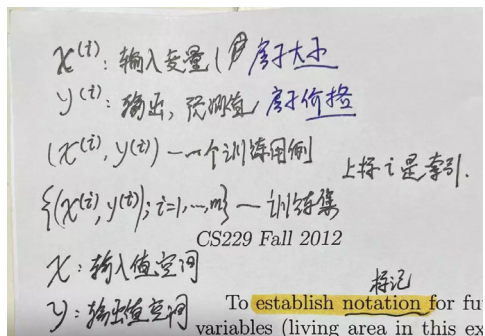
Supervised Learning

房价例子

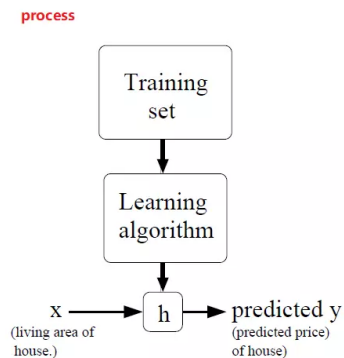
Living area (feet ²)	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
⋮	⋮



Notation



Process



训练集通过 学习算法 可以生成一个 假设H函数, 通过这个函数, 可以对 输入的房屋大小 预测出 房屋价格。

which we also know the number of bedrooms in each house

Living area (feet ²) x_1	#bedrooms x_2	Price (1000\$) y
2104 $x_1^{(1)}$	3	400
1600	3	330
2400 $x_1^{(3)}$	3	369
1416	2	232
3000	4	540
\vdots	\vdots	\vdots

how to represent the Hypothesis?

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

因为假设 h 是由两个参数 θ 和 x 共同决定的所以写成左边那种形式，然后写成 y 近似为 x 的线性函数，之后用 θ 对 x 到 y 线性函数空间进行参数化。

$x_0 = 1$ (this is the **intercept term**), so that
 截距 $x_0 = 1$

$$h(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

$\rightarrow = [\theta_0 \ \theta_1 \ \theta_2] \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$

把 θ 和 x 都看成“向量”

where on the right-hand side above we are viewing θ and x both as vectors, and here n is the number of input variables (not counting x_0).

n 是输入变量的个数

$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$ $\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$

\vec{x} 房子面积 \rightarrow 房子面积
 \vec{x} 房间数目 \rightarrow 房间数目

x : input features,

x_1 : 房屋大小

x_2 : 房间数目

术语:

θ : 参数

m : 训练实例个数，表格中一行代表一个 training example

x : input/features,

y : output/target variable

(x, y) : training example

$(x^{(i)}, y^{(i)})$: 第 i 个训练实例

	Living area (feet ²)	Price (1000\$)
例:	2104 $x_1^{(1)}$	400
	1600 $x_1^{(2)}$	330

n : features, 特征个数，这里的特征个数为2 (area、bedrooms)

how choose parameters?

选择参数，让房子预估值($h_{\theta}(x)$) 离 真实价格(y)很近

定义一个函数 $J(\theta)$ 表示这个接近程度

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

我们要做的事：找到参数 θ 使得 $J(\theta)$ 最小，从而 房子预估值 与 真实值 就越接近

how to implement this algorithm?

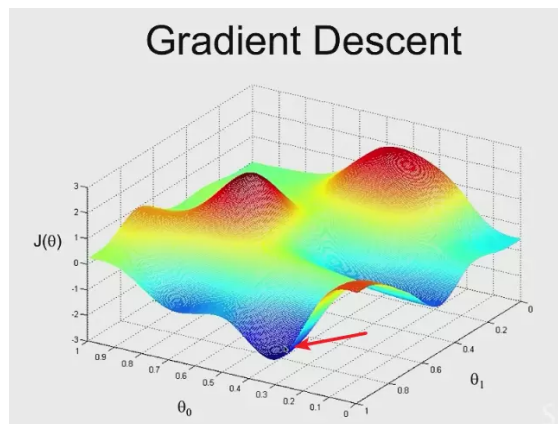
Gradient descent (梯度下降)

start with some initial θ ,

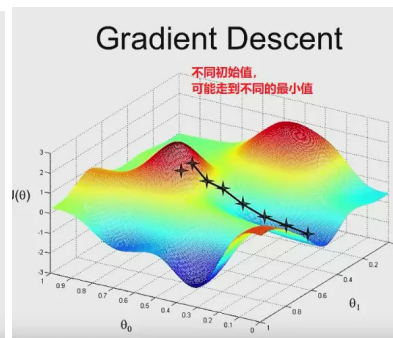
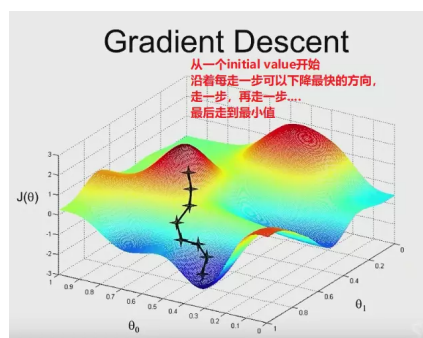
keep changing θ to reduce $J(\theta)$

解释：从一个初始值开始，然后重复修改 θ 值，以使减小 $J(\theta)$ ，最终收敛到最小值

梯度可视化：



在这个例子中， θ 有两个参数，就是找到 适合的 θ_0 和 θ_1 使得 $J(\theta)$ 最小，（这里箭头指向的这个值就是了）



转一圈找到走一步就可以下降最大的方向

formalize the gradient descent algorithm

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

这里的 $:=$ 为赋值运算符(assign), $=$ 是判断运算符

每一步梯度下降，都要选取一个 θ_j ，这里 j 取值为 $0, 1, 2, \dots, n$ ，（ n 为features个数）

α 为 learning rate,

后面是 $J(\theta)$ 对 θ_j 求偏导数

求偏导的过程

the definition of J . We have:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \quad (1) \\ &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \quad (2) \\ &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h_\theta(x) - y) x_j \end{aligned}$$

$\because h_\theta(x) = \sum_{i=0}^n \theta_i x_i$

$\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_j x_j + \dots + \theta_n x_n - y$

只有 $\theta_j x_j$ 这一项可对 θ_j 求偏导，其他均为常数 $\rightarrow 0$

$\frac{\partial}{\partial \theta_j} (\theta_j x_j) = x_j$

综上，表达式如下

1个training example $\theta_j := \theta_j - \alpha (h_\theta(x) - y) x_j$

m个 training examples $\theta_j := \theta_j - \alpha \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$

$x^{(i)}$ 是第 i 个 训练实例的 input features

$y^{(i)}$ 是第 i 个 训练实例的 target label

Batch Gradient Descent

batch的含义：

例如房价中有49个training example，则就可以把这49个例子看做成一个 batch of data

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (j=0,1,2,\dots,n \quad n \text{ 是 features})$$

}

n: features, n=2, (房子大小、房间个数)

m: number of training examples, 训练实例个数

用房价例子举例，n=2 (area、bedrooms)，m=3, 3个训练用例

Living area (feet ²)	#bedrooms	Price (1000\$)
2104 $x_1^{(1)}$	3	400
1600	3	330
2400 $x_1^{(3)}$	3	369

对于参数 θ_0 :

Repeat until convergence{

$$\theta_0 := \theta_0 - \alpha \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$:= \theta_0 - \alpha ((h_\theta(x^{(1)}) - y^{(1)}) x_0^{(1)} + (h_\theta(x^{(2)}) - y^{(2)}) x_0^{(2)} + (h_\theta(x^{(3)}) - y^{(3)}) x_0^{(3)})$$

}

对于参数 θ_1 :

Repeat until convergence{

$$\theta_1 := \theta_1 - \alpha \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

}

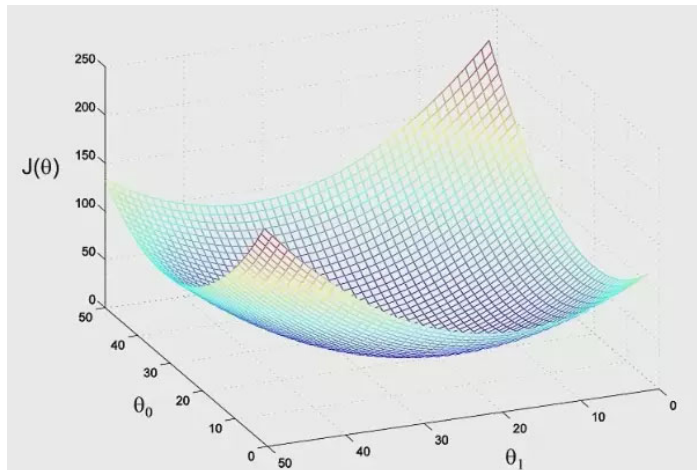
对于参数 θ_2 :

Repeat until convergence{

$$\theta_2 := \theta_2 - \alpha \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

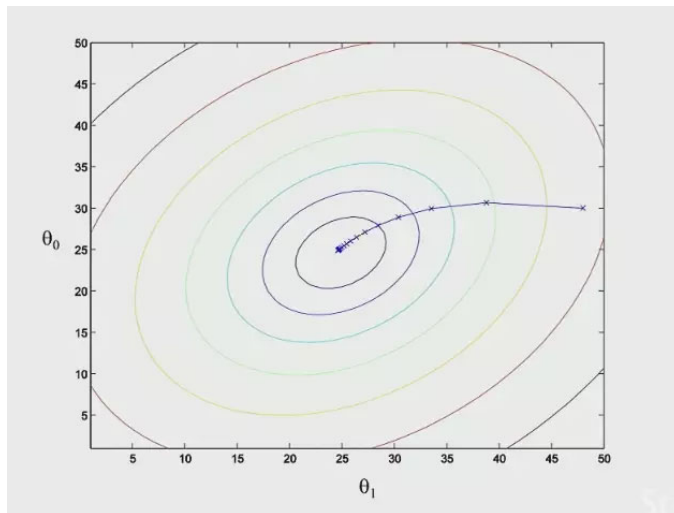
}

如果把 θ_j 定义为 若干平方项求和，那么它就是一个二次函数，那么 θ_j 就是下面这个形式：



θ_j 没有局部最优，局部最优就是全局最优

看这个函数的另一种方法是看 轮廓线，对上面这个大碗进行水平分层，



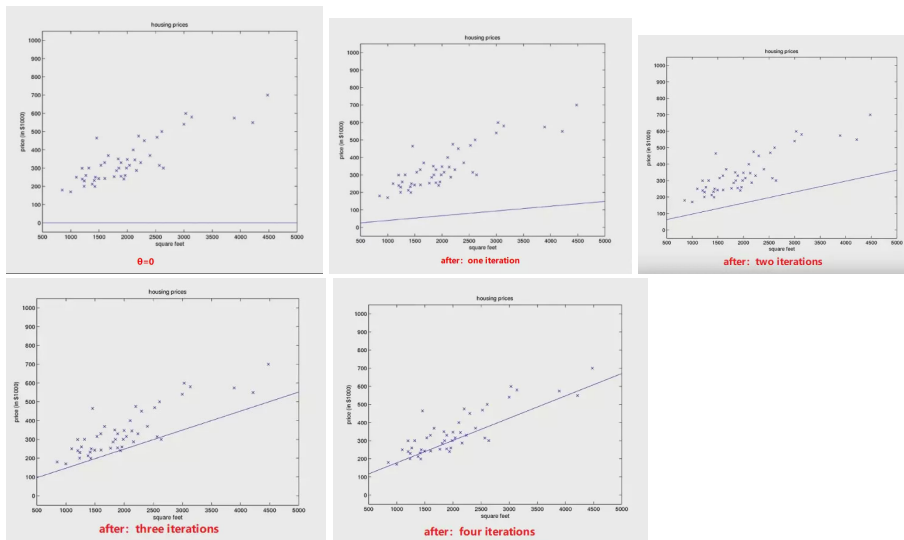
the choice of the learning rate α

设置的很大，跨的步很大，可能会错过最小的地方；

设置的很小，需要迭代很多次，算法会很慢。

设置经验：

随机尝试几个值，看如果哪个值可以最有效的推动 θ_j 的下降，如果你看到 θ_j 在增加而不是减少，就说明你设置的 α 太大了。Ng的做法是在指数尺度上面尝试。



BGD的缺点:

当数据集很大的时候, 更新一次 θ , 就要对整个数据集进行求和一次

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (j=0,1,2,\dots,n \quad n \text{ 是 features})$$

} 当数据集很大, 即m很大的时候, 更新一次 θ_j , 就要求和一次m个数据!

Stochastic Gradient Descent (随机梯度下降)

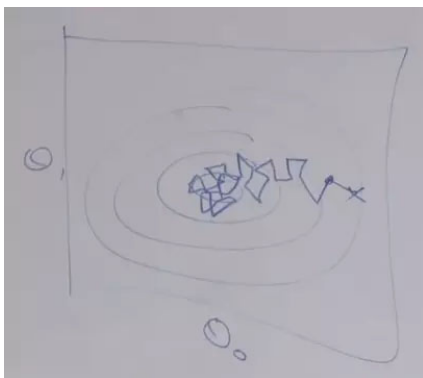
Repeat {

for i=1 to m, {

$$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (\text{for every } j)$$

}

}



在某个地方初始化参数之后, 看第一个房子, 修改参数提高预测当前房子价格的准确性, 只对当前这个房子的数据进行了拟合, 虽然改善了参数, 但是并没有朝最直接的方向走下去。每次改参数都是为了更好的拟合当前那个房子, 若干次之后, 它在最小值附近震荡, 但不会收敛于最小值。

数据集很大时, 通常使用SGD (随机梯度下降算法)

Q: 可以从SGD切换到BGD吗?

A: 可以, 但是当数据集很大时, BGD太慢了, 而且SGD生成的参数性能也不错。

数据集很大时，用SGD；

数据集很小时，用BGD，可以得到全局最优，而不用振荡。

如果使用的算法是线性回归，有一种方法可以直接解出参数theta的最佳值，直接一步跳到全局最优，而不需要使用迭代算法。

2 The Normal equations

Matrix derivatives (矩阵求导)

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

例：

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$
$$f(A) = \frac{3}{2}A_{11} + 5A_{12}^2 + A_{21}A_{22}$$
$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} \end{bmatrix}$$
$$\Rightarrow \begin{bmatrix} \frac{3}{2} & 10A_{12} \\ A_{22} & A_{21} \end{bmatrix}$$
$$\begin{cases} \frac{\partial f}{\partial A_{11}} = \frac{\partial}{\partial A_{11}} \left(\frac{3}{2}A_{11} + 5A_{12}^2 + A_{21}A_{22} \right) = \frac{3}{2} \\ \frac{\partial f}{\partial A_{12}} = \frac{\partial}{\partial A_{12}} \left(\frac{3}{2}A_{11} + 5A_{12}^2 + A_{21}A_{22} \right) = 10A_{12} \\ \frac{\partial f}{\partial A_{21}} = \frac{\partial}{\partial A_{21}} \left(\frac{3}{2}A_{11} + 5A_{12}^2 + A_{21}A_{22} \right) = A_{22} \\ \frac{\partial f}{\partial A_{22}} = \frac{\partial}{\partial A_{22}} \left(\frac{3}{2}A_{11} + 5A_{12}^2 + A_{21}A_{22} \right) = A_{21} \end{cases}$$

trace operator (矩阵的迹)

$$\text{tr} A = \sum_{i=1}^n A_{ii}$$

运算法则

$$\text{tr} A = \text{tr} A^T$$

$$\text{tr} AB = \text{tr} BA$$

$$\text{tr} ABC = \text{tr} CAB = \text{tr} BCA$$

$$\nabla_A \text{tr} AB = B^T$$

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix}$$

$$AB = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} & a_{11}b_{13} + a_{12}b_{23} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} & a_{21}b_{13} + a_{22}b_{23} \\ a_{31}b_{11} + a_{32}b_{21} & a_{31}b_{12} + a_{32}b_{22} & a_{31}b_{13} + a_{32}b_{23} \end{bmatrix}$$

$$\text{tr} AB = a_{11}b_{11} + a_{12}b_{21} + a_{21}b_{12} + a_{22}b_{22} + a_{31}b_{13} + a_{32}b_{23}$$

$$\frac{\partial \text{tr} AB}{\partial a_{11}} = b_{11} \quad \frac{\partial \text{tr} AB}{\partial a_{21}} = b_{12} \quad \frac{\partial \text{tr} AB}{\partial a_{31}} = b_{13}$$

$$\frac{\partial \text{tr} AB}{\partial a_{12}} = b_{21} \quad \frac{\partial \text{tr} AB}{\partial a_{22}} = b_{22} \quad \frac{\partial \text{tr} AB}{\partial a_{32}} = b_{23}$$

$$\Rightarrow \nabla_{\text{tr} AB} = \begin{bmatrix} \frac{\partial \text{tr} AB}{\partial a_{11}} & \frac{\partial \text{tr} AB}{\partial a_{12}} \\ \frac{\partial \text{tr} AB}{\partial a_{21}} & \frac{\partial \text{tr} AB}{\partial a_{22}} \\ \frac{\partial \text{tr} AB}{\partial a_{31}} & \frac{\partial \text{tr} AB}{\partial a_{32}} \end{bmatrix} = \begin{bmatrix} b_{11} & b_{21} \\ b_{12} & b_{22} \\ b_{13} & b_{23} \end{bmatrix} = B^T$$

$$\nabla_{\text{tr} AA^T C} = CA + C^T A$$

矩阵向量的形式重写 $J(\theta)$

m : 训练实例的个数

design X

x: input values

$$X = \begin{bmatrix} - & (x^{(1)})^T & - \\ - & (x^{(2)})^T & - \\ & \vdots & \\ - & (x^{(m)})^T & - \end{bmatrix}$$

y: target values

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

由

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

$$X\theta = \begin{pmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{pmatrix} \cdot \theta = \begin{pmatrix} (x^{(1)})^T \cdot \theta \\ (x^{(2)})^T \cdot \theta \\ \vdots \\ (x^{(m)})^T \cdot \theta \end{pmatrix} = \begin{pmatrix} h_{\theta}(x^{(1)}) \\ h_{\theta}(x^{(2)}) \\ \vdots \\ h_{\theta}(x^{(m)}) \end{pmatrix}$$

$$X\theta - \vec{y} = \begin{pmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{pmatrix}$$

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

$$\text{由 } d^2d = a_1^2 + a_2^2 + \dots + a_n^2 = \sum_{i=1}^n a_i^2$$

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

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

$$= \left(\frac{1}{2} \cdot \nabla_{\theta} (X\theta - \vec{y})^T (X\theta - \vec{y}) \right) \rightarrow \text{拆开}$$

$$= \frac{1}{2} \cdot \nabla_{\theta} [(X\theta)^T \cdot X\theta - (X\theta)^T \cdot \vec{y} - \vec{y}^T \cdot X\theta + \vec{y}^T \vec{y}]$$

$$= \frac{1}{2} \cdot \nabla_{\theta} [\theta^T X^T X\theta - \theta^T X^T \vec{y} - \vec{y}^T \cdot X\theta + \vec{y}^T \vec{y}]$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr}(\theta^T X^T X\theta - \theta^T X^T \vec{y} - \vec{y}^T \cdot X\theta + \vec{y}^T \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} [\text{tr} \theta^T X^T X\theta - \text{tr} \theta^T X^T \vec{y} - \text{tr} \vec{y}^T X\theta + \text{tr} \vec{y}^T \vec{y}]$$

$$= \frac{1}{2} \nabla_{\theta} [\text{tr} \theta^T X^T X\theta - 2 \text{tr} \vec{y}^T X\theta + \vec{y}^T \vec{y}]$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr} \theta^T X^T X\theta - \nabla_{\theta} \text{tr} \vec{y}^T X\theta + 0$$

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

$$= X^T X\theta - X^T \vec{y}$$

为了求最小值，让导数等于0，

即： $X^T X\theta = X^T \vec{y}$

则求得： $\theta = (X^T X)^{-1} X^T \vec{y}$