

Lecture 4

Regularization

Dr. Le Huu Ton

Outline



Overfitting Problem



Regularization



Regularization with Linear Regression



Regularization with Logistic Regression

Outline



Overfitting Problem



Regularization

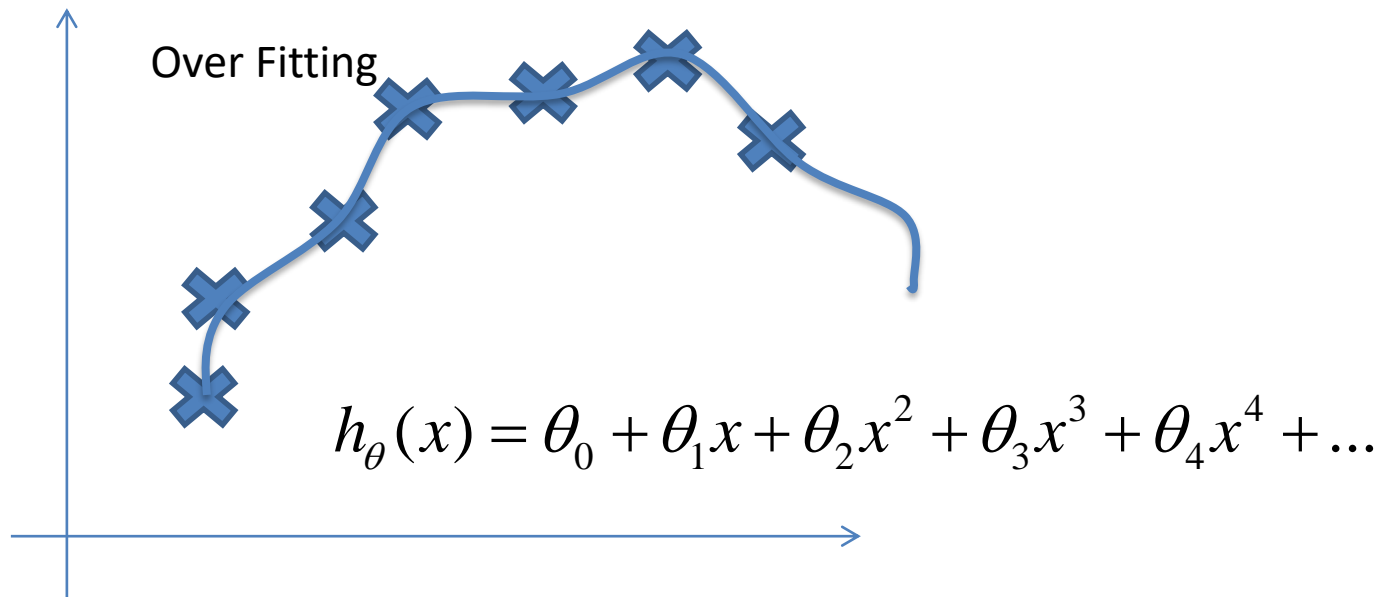
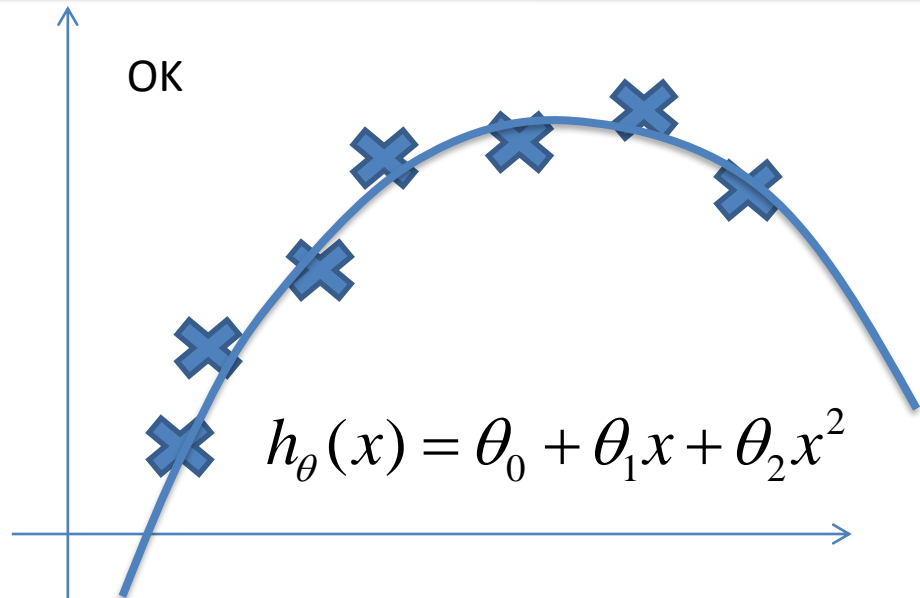
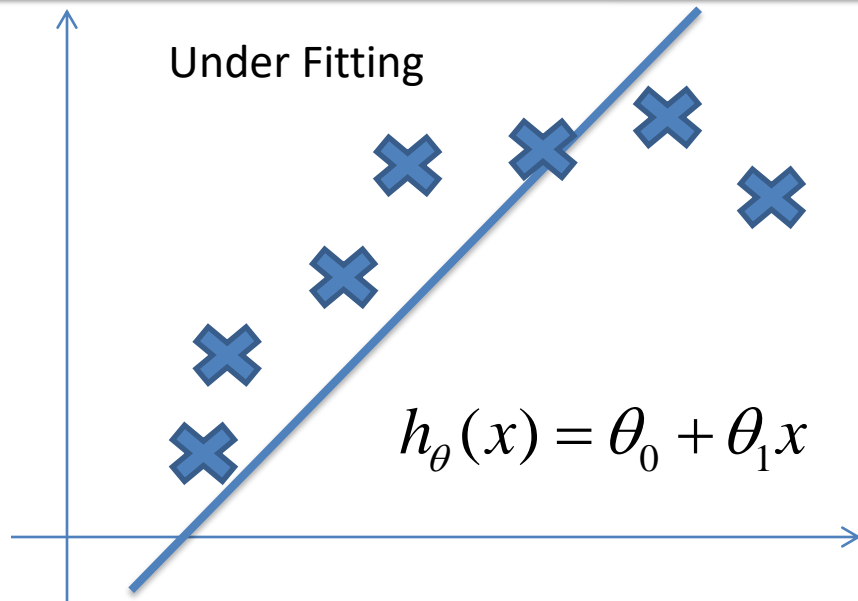


Regularization with Linear Regression

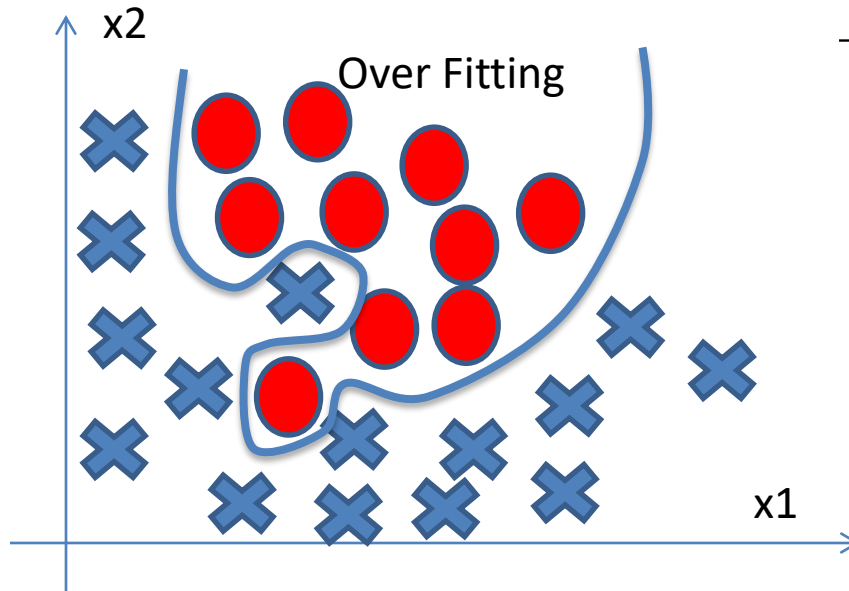
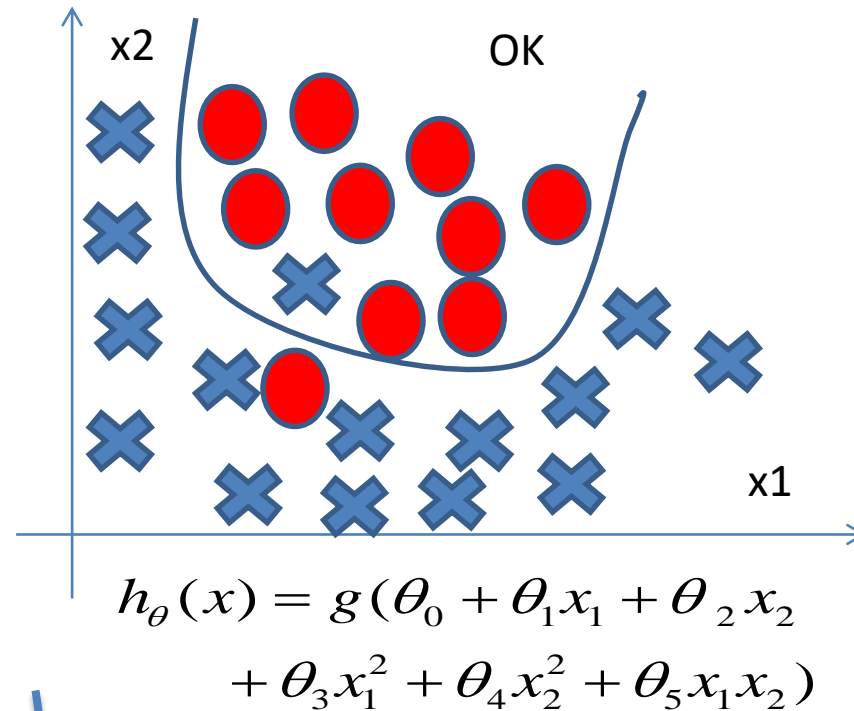
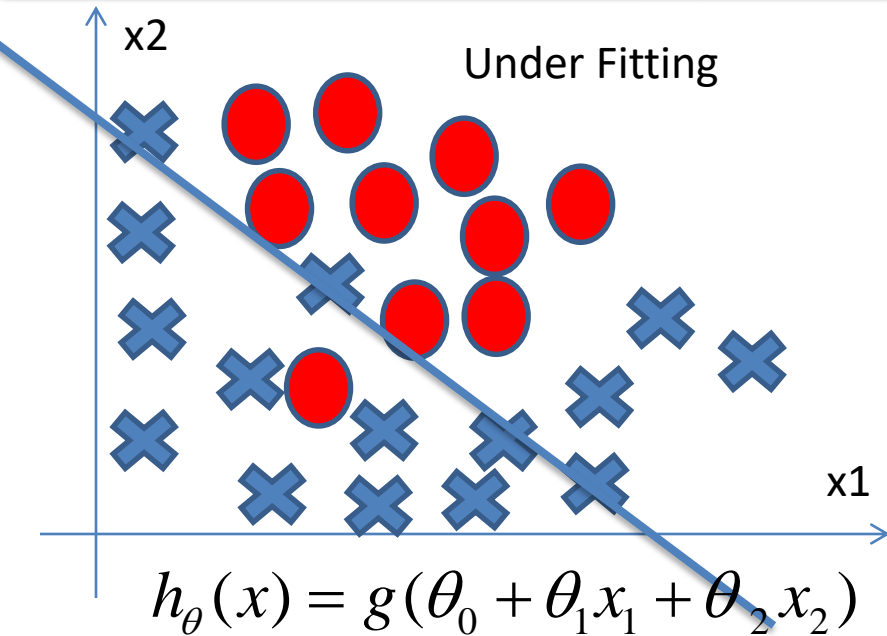


Regularization with Logistic Regression

Overfitting Problem



Overfitting Problem



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1^3 + \theta_6 x_2^3 + \dots)$$

Overfitting Problem

Under fitting:

Under fitting refers to a model that can neither model the training data nor generalize to new data.

An under fit machine learning model is not a suitable model and will be obvious as it will have poor performance on the training data.

Over Fitting :

Overfitting happens when a model learns the detail and noise in the training data to the extent that it negatively impacts the performance on the model on new data.

Outline



Overfitting Problem



Regularization



Regularization with Linear Regression



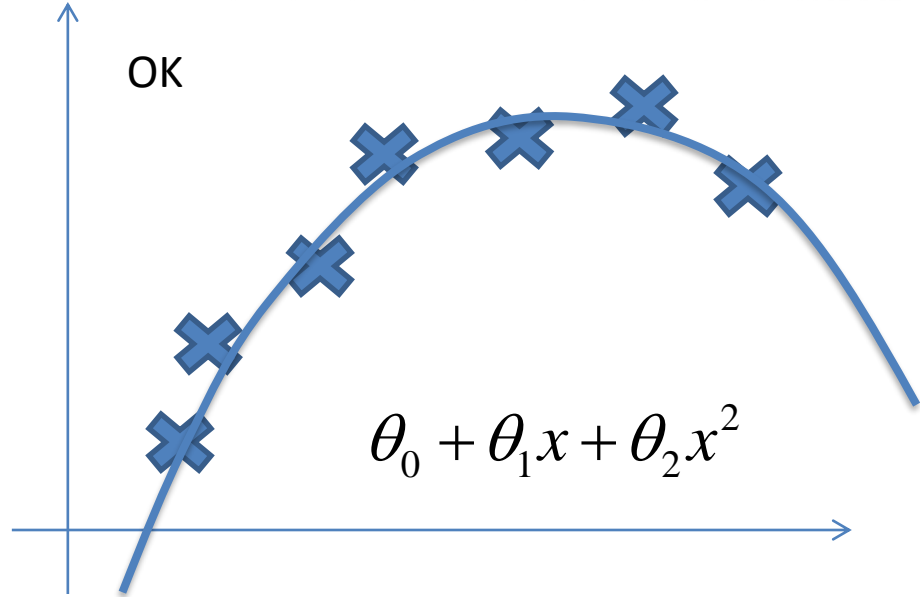
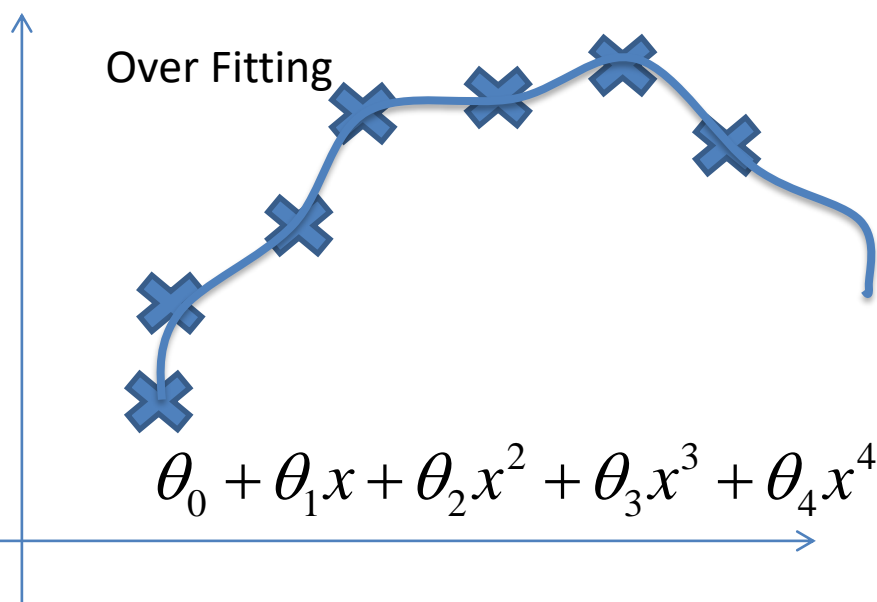
Regularization with Logistic Regression

Regularization

Regularization is a *technique* used in an attempt to solve the **overfitting** problem.

Regularization is done by reduce the magnitude of some coefficient θ_j

Overfitting Problem



Regularization: reduce value of θ_3 and θ_4

Minimize the cost function

$$E(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 9999\theta_3 + 9999\theta_4$$
$$\Rightarrow \theta_3 \approx 0, + \theta_4 \approx 0$$

Regularization

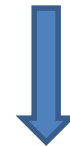
Small values of coefficients $\theta_0, \theta_1, \dots, \theta_n$

⇒ Simpler hypothesis $h(x)$

⇒ Less prone to overfitting

Regularization: Add a regularization component into the cost function

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



Regularization component

Regularization

Question:

What if λ is set by a extremely large number (too large for our problem), which of the following statement is correct:

1. The algorithm works fine
2. Algorithm fail to eliminate overfitting
3. Algorithm results in under fitting
4. Gradient descent will fail to converge

Outline



Overfitting Problem



Regularization



Regularization with Linear Regression



Regularization with Logistic Regression

Regularization with Linear Regression

Regularization:

Minimize the Cost Function

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

Gradient descent:

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

Regularization with Linear Regression

Gradient Descent:

Repeat until converged:

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

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_j^i - \frac{\alpha \lambda}{m} \theta_j \quad \forall j = 1 : n$$

$$\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_j^i \quad \forall j = 1 : n$$

$\}$

Regularization with Linear Regression

Normal Equation without regularization:

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

Normal Equation with regularization

$$\theta = (X^T X + \lambda \begin{vmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 \\ \dots & \dots & 1 & \dots \\ 0 & 0 & 0 & 1 \end{vmatrix})^{-1} X^T Y$$

Outline



Overfitting Problem



Regularization



Regularization with Linear Regression



Regularization with Logistic Regression

Regularization with Logistic Regression

Logistic Regression: Minimize the cost function

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

Gradient descent:

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

Regularization with Logistic Regression

Gradient Descent:

Repeat until converged:

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

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m [(h(x^{(i)}) - y^{(i)}) x_0^i] - \frac{\lambda}{m} \theta_j \quad \forall j = 1 : n$$

$$\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_0^i \quad \forall j = 1 : n$$

}

Regularization with Logistic Regression

Newton's Method with Regularization

$$\theta^{t+1} := \theta^t - H^{-1} \Delta_{\theta} E$$

$$\Delta_{\theta} E = \begin{bmatrix} \frac{\partial}{\partial \theta_0} E(\theta) \\ \dots \\ \frac{\partial}{\partial \theta_n} E(\theta) \end{bmatrix} = \begin{bmatrix} \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) x_0^i \\ \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) x_1^i + \frac{\lambda}{m} \theta_1 \\ \dots \\ \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) x_n^i + \frac{\lambda}{m} \theta_n \end{bmatrix}$$

Regularization with Logistic Regression

Hessian Matrix:

$$H = \frac{1}{m} \sum_{i=1}^m \left[h(x^{(i)})(1-h(x^{(i)}))x^{(i)}(x^{(i)})^T \right] + \lambda \begin{vmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 \\ \dots & \dots & 1 & \dots \\ 0 & 0 & 0 & 1 \end{vmatrix}$$

Regularization with Logistic Regression

When using regularized logistic regression, which of these is the best way to monitor whether gradient descent is working correctly?

- ☐ Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right]$ as a function of the number of iterations, and make sure it's decreasing.
- ☐ Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right] - \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ as a function of the number of iterations, and make sure it's decreasing.
- ☐ Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ as a function of the number of iterations, and make sure it's decreasing.
- ☐ Plot $\sum_{j=1}^n \theta_j^2$ as a function of the number of iterations, and make sure it's decreasing.

References

<http://openclassroom.stanford.edu/MainFolder/CoursePage.php?course=MachineLearning>

Andrew Ng Slides:

https://www.google.com.vn/url?sa=t&rct=j&q=&esrc=s&source=web&cd=2&cad=rja&uact=8&sqi=2&ved=0ahUKEwjNt4fdvMDPAhXIn5QKHZO1BSgQFggfMAE&url=https%3A%2F%2Fdatajobs.com%2Fdata-science-repo%2FGeneralized-Linear-Models-%5BAndrew-Ng%5D.pdf&usg=AFQjCNGq37q2uVFcpGhNqH-5KZSIJ_HSxg&sig2=vnCEvyvKQGCuryttAPcokw&bvm=bv.134495766,d.dGo

At one iteration $\theta_0 = 1$, $\theta_1 = 2$, $\theta_2 = 1$, $\alpha = 6$, regularization term $\lambda = 10$

What is the value of θ_0 , θ_1 , θ_2 after that iteration

Size (m ²)	N ^o of floors	Price (billion VND)
30	3	2
40	4	3
20	2	2
50	4	5
40	3	3
20	1	2

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_j^i \quad \forall j = 1:n$$

Homework

Exercise:

Starting with θ_0 and θ_1 equal to 0. $\alpha = 0.01$, Regularization term $\lambda = 10$. Calculate the value of coefficient after first iteration using gradient.

Calculate the Hessian matrix and Derivative vector if Newton methods is used

Price	Location	Output Value
2.5	<i>Thanh Xuan</i>	0
3.5	<i>Thanh Xuan</i>	0
5.6	<i>Hoan Kiem</i>	1
2.2	<i>Thanh Xuan</i>	0
6.9	<i>Hoan Kiem</i>	1
9.6	<i>Hoan Kiem</i>	1

$$\Delta_{\theta} E = \begin{bmatrix} \frac{\partial}{\partial \theta_0} E(\theta) \\ \dots \\ \frac{\partial}{\partial \theta_n} E(\theta) \end{bmatrix} = \begin{bmatrix} \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) x_0^i \\ \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) x_1^i + \frac{\lambda}{m} \theta_1 \\ \dots \\ \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) x_n^i + \frac{\lambda}{m} \theta_n \end{bmatrix}$$

$$H = \frac{1}{m} \sum_{i=1}^m \left[h(x^{(i)}) (1 - h(x^{(i)})) x^{(i)} (x^{(i)})^T \right] + \lambda \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 \\ \dots & \dots & 1 & \dots \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

After one iteration $\theta_0 = 1$,
 $\theta_1 = 2$, $\theta_2 = 1$, $\alpha = 6$,
 regularization term $\lambda = 10$

What is the value of θ_0 ,
 θ_1 , θ_2 before that
 iteration

Size (m ²)	N ^o of floors	Price (billion VND)
30	3	2
40	4	3
20	2	2
50	4	5
40	3	3
20	1	2

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_j^i \quad \forall j = 1:n$$