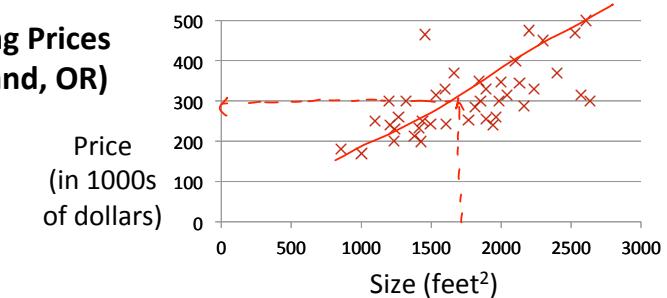


Linear Regression with One Variable

Basic Concepts

Housing Prices (Portland, OR)



Supervised Learning

Given the “right answer” for each example in the data.

Regression Problem

Predict real-valued output

Training set of housing prices (Portland, OR)	Size in feet ² (x)	Price (\$) in 1000's (y)
	$x^{(1)} = 2104$	$y^{(1)} = 460$
	$x^{(2)} = 1416$	$y^{(2)} = 232$
	\vdots	\vdots
	1534	315
	852	178
	$x^{(m)}$	$y^{(m)}$

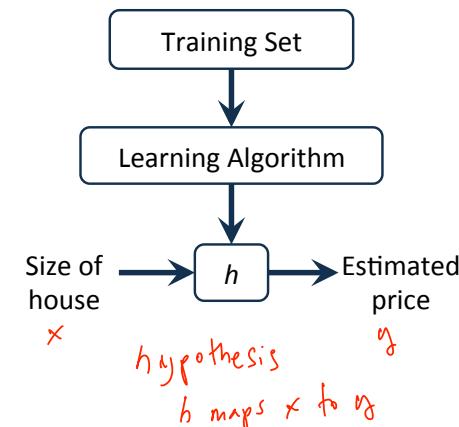
Notation:

m = Number of training examples

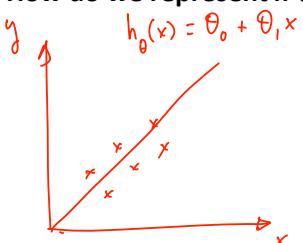
x 's = “input” variable / features

y 's = “output” variable / “target” variable

i th training example $\equiv (x^{(i)}, y^{(i)})$, $i = 1 \dots m$

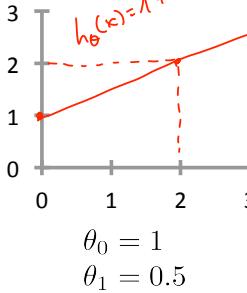
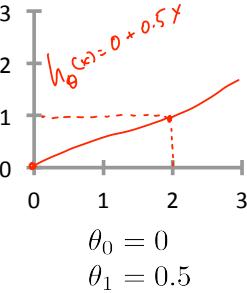
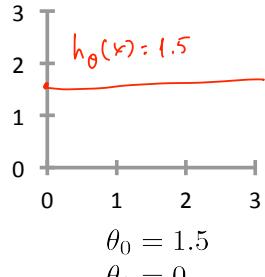


How do we represent h ?



Linear regression with one variable.
Univariate linear regression.

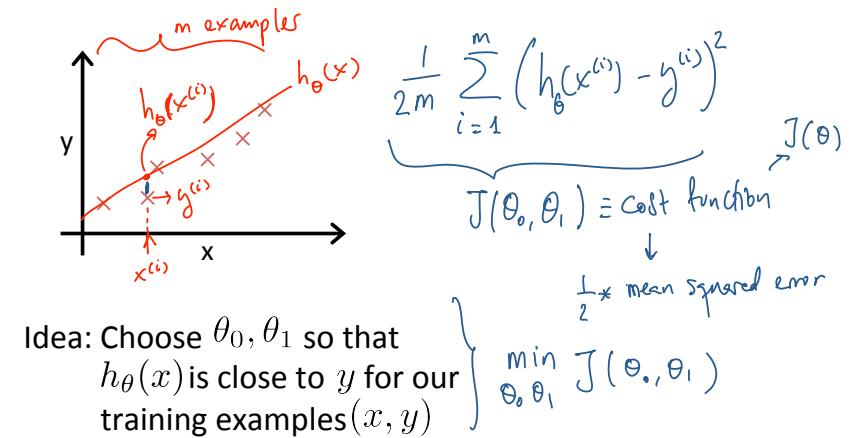
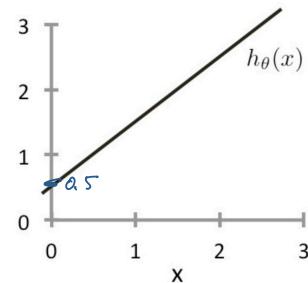
Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$



How to choose θ_0, θ_1 ?

Given the hypothesis $h_\theta(x) = \theta_0 + \theta_1 x$, in the figure, select the right parameters values:

- $\theta_0 = 0, \theta_1 = 1$
- $\theta_0 = 0.5, \theta_1 = 1$
- $\theta_0 = 1, \theta_1 = 0.5$
- $\theta_0 = 1, \theta_1 = 1$



Simplified

Hypothesis:
 $h_\theta(x) = \theta_0 + \theta_1 x$

$$h_\theta(x) = \theta_1 x$$

Parameters:
 θ_0, θ_1

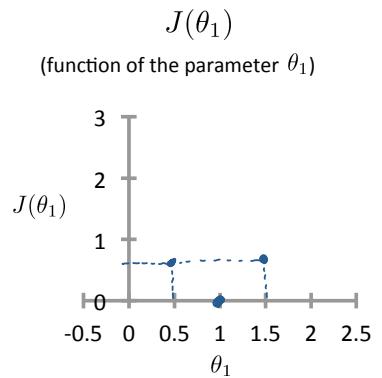
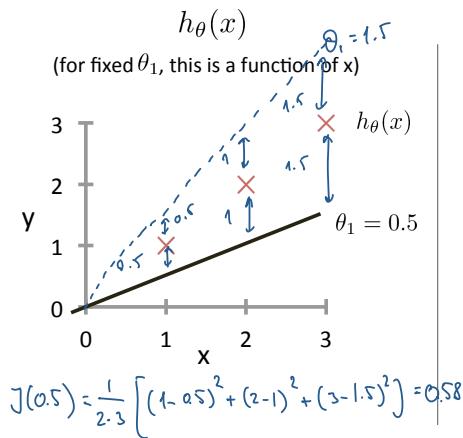
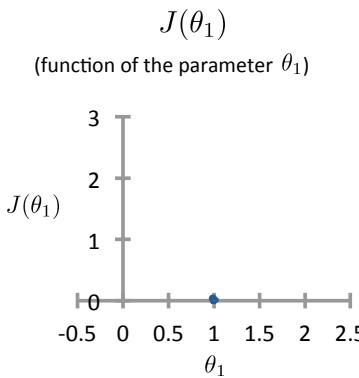
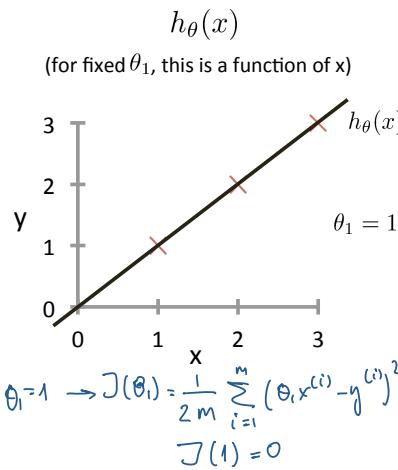
$$\theta_1$$

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

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

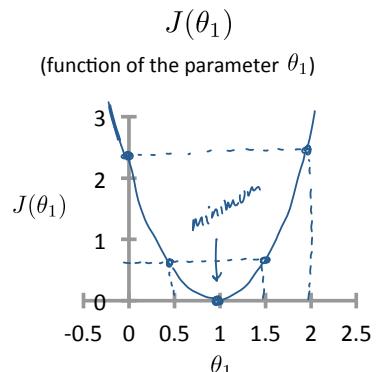
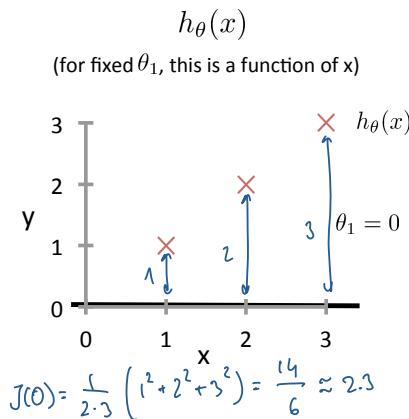
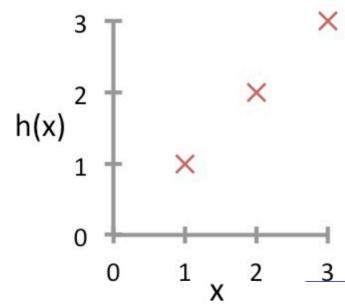
Goal: minimize $J(\theta_0, \theta_1)$

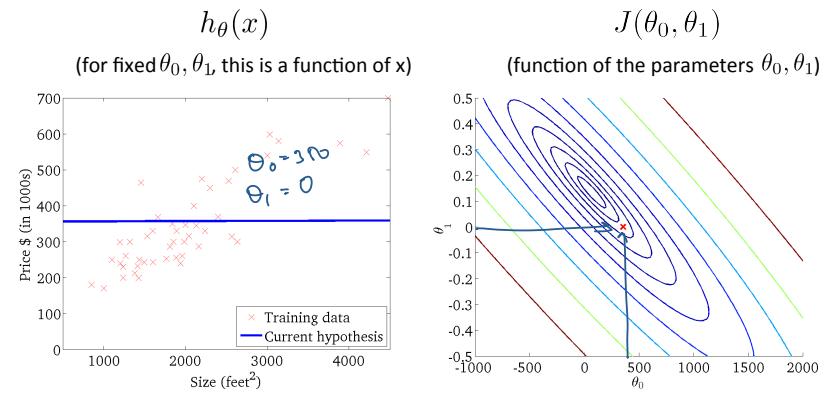
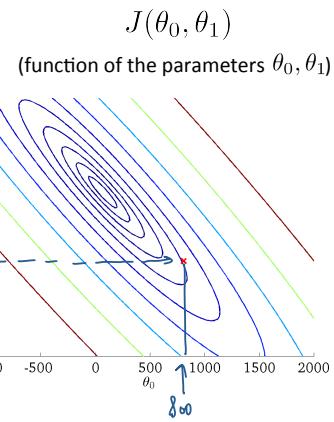
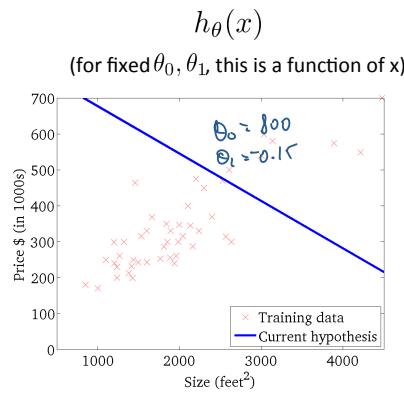
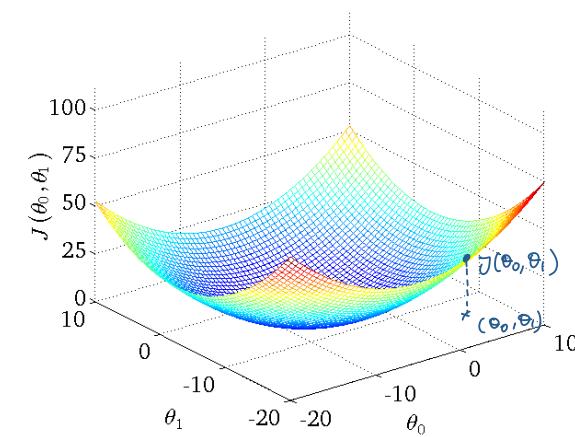
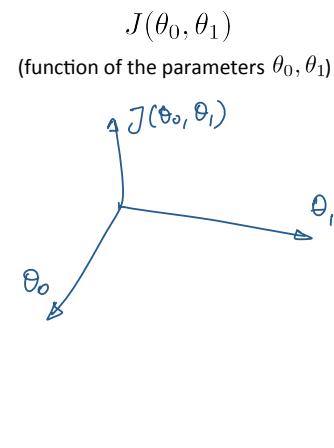
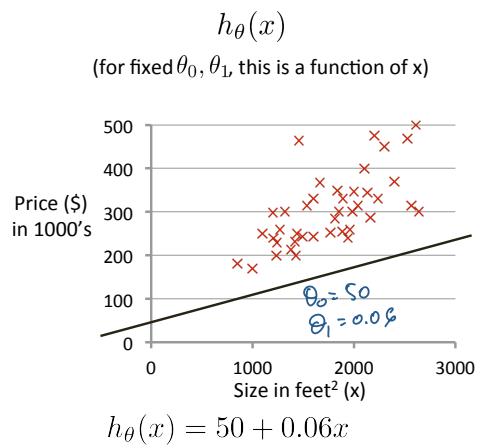
$$\min_{\theta_1} J(\theta_1)$$

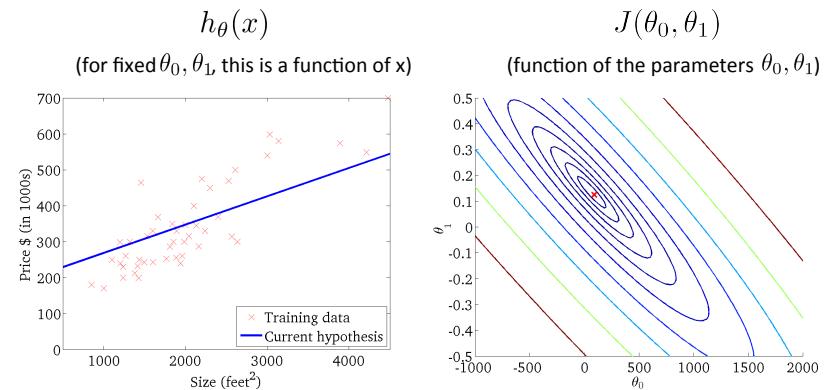
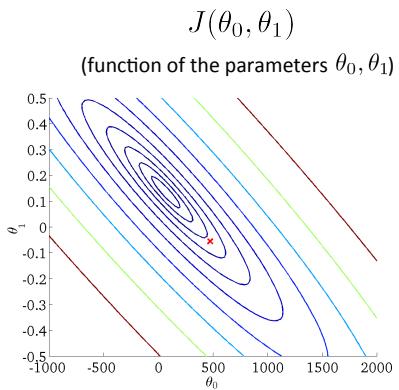
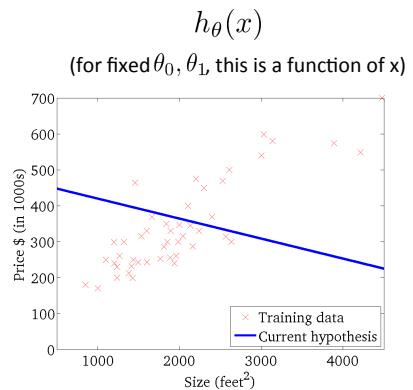


Consider $m=3$ inputs, shown in the figure.
The hypothesis function is $h_\theta(x) = \theta_1 x$
And the cost function is $J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$
Select $J(0)$

- 0
- 1/6
- 1
- 14/6







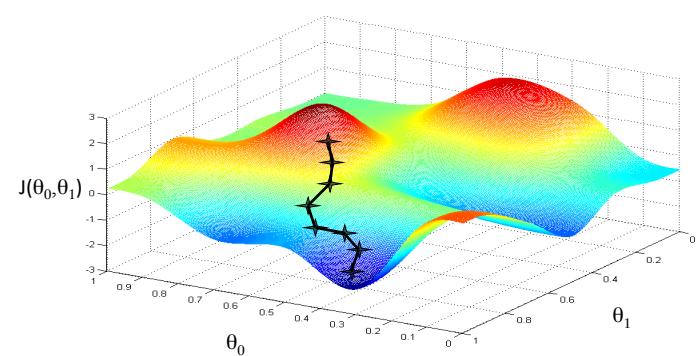
Gradient Descent Algorithm

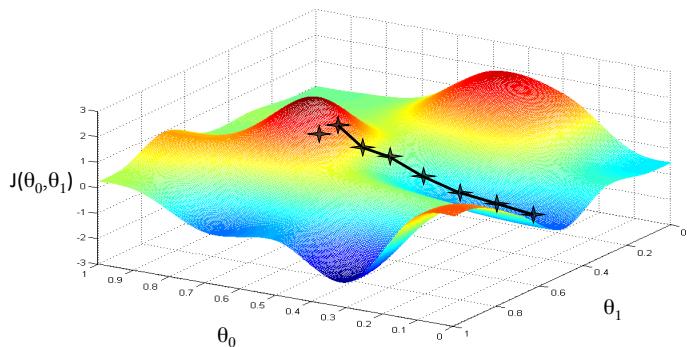
Have some function $J(\theta_0, \theta_1)$

Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Outline:

- Start with some θ_0, θ_1
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum





Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

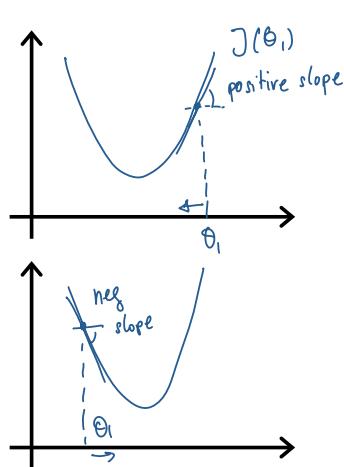
↓
learning rate

Correct: Simultaneous update

$$\begin{aligned} \text{temp0} &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp1} &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 &:= \text{temp0} \\ \theta_1 &:= \text{temp1} \end{aligned}$$

Incorrect:

$$\begin{aligned} \text{temp0} &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_0 &:= \text{temp0} \\ \text{temp1} &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_1 &:= \text{temp1} \end{aligned}$$



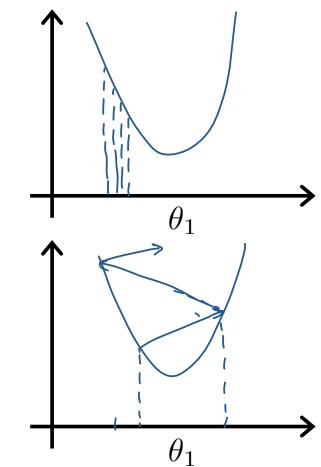
$$\begin{aligned} \theta_1 &:= \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) \\ &\quad \text{---} \\ \theta_1 &:= \theta_1 - \alpha (\text{positive value}) \end{aligned}$$

$$\begin{aligned} \theta_1 &:= \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) \\ &\quad \text{---} \\ \theta_1 &:= \theta_1 - \alpha (\text{negative}) \end{aligned}$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

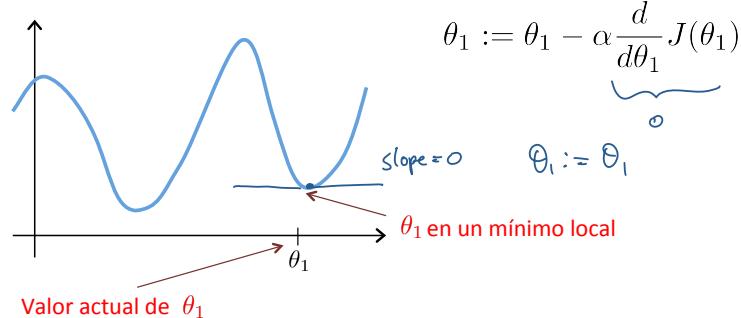
If α is too small, gradient descent can be slow.

If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



Aspectos Prácticos

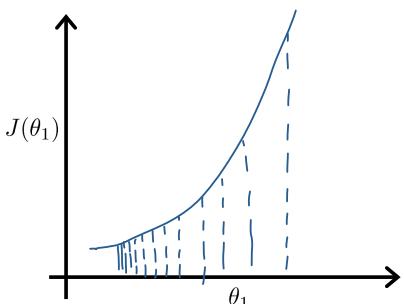
What happens in the gradient descent method if θ_1 is already at a local minimum?



Gradient descent can converge to a local minimum, even with the learning rate α fixed.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.



Gradient descent for linear regression

Gradient descent algorithm

```
repeat until convergence {
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ 
    (for  $j = 1$  and  $j = 0$ )
}
```

Linear Regression Model

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

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

Gradient descent for linear regression

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right) = \\ &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \right) \end{aligned}$$

$$\hat{j} = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

$$\begin{aligned} j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned}$$

Gradient descent for linear regression

Algorithm:

repeat until convergence {

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

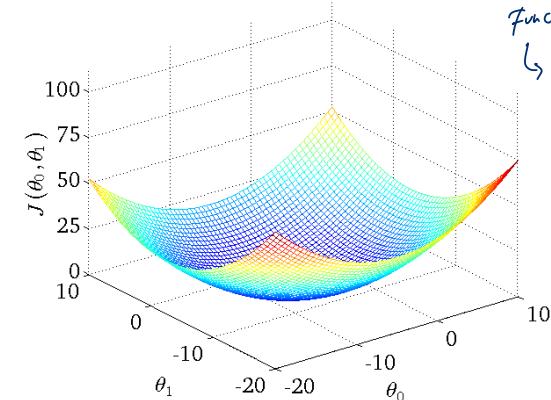
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$
 $\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$
 Update
 θ_0 and θ_1
 simultaneously

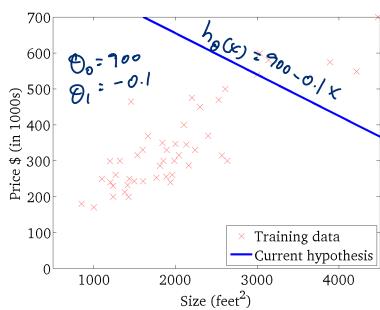
Gradient descent for linear regression

Convex function
 \hookrightarrow 1 global minimum



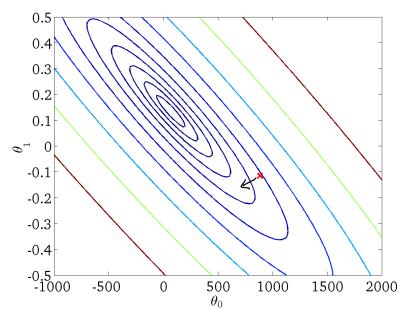
$$h_\theta(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



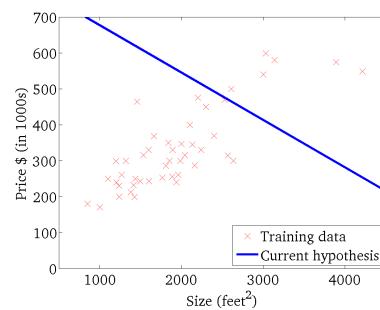
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



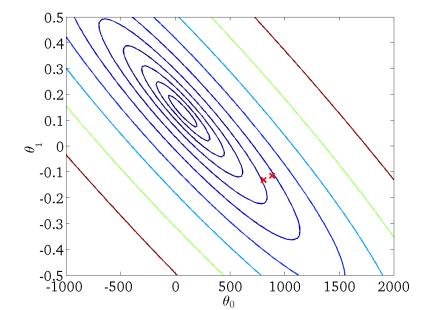
$$h_\theta(x)$$

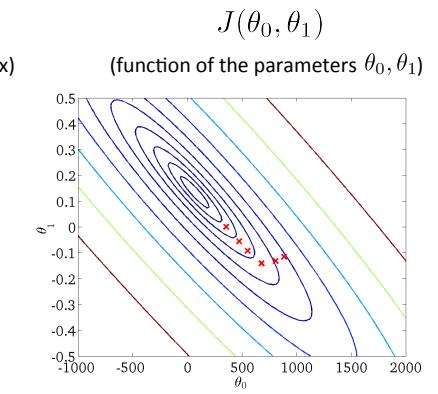
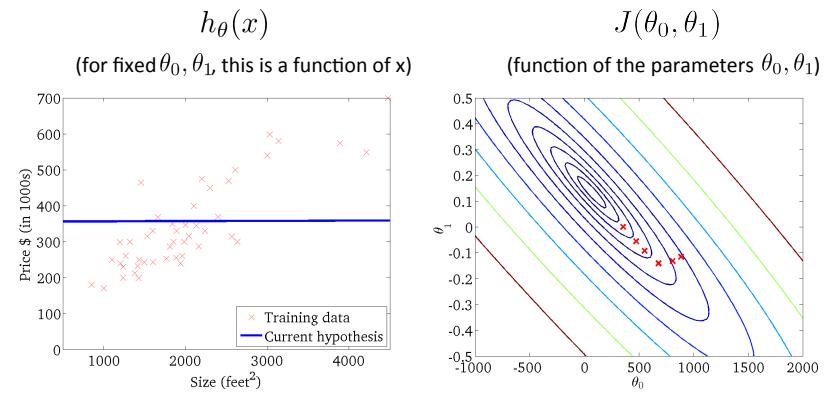
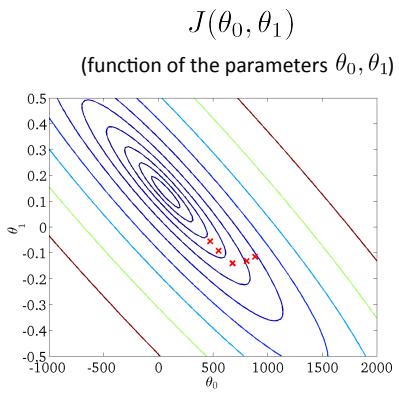
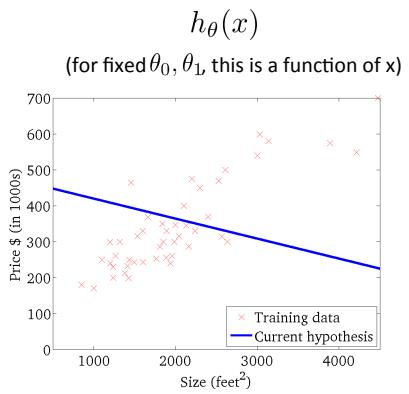
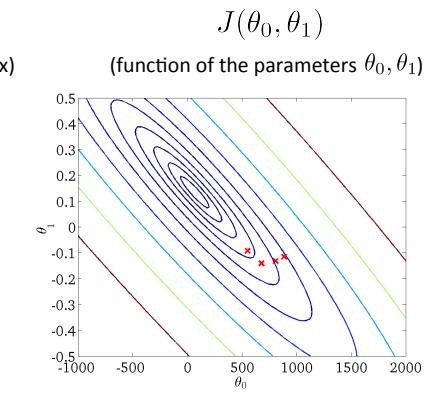
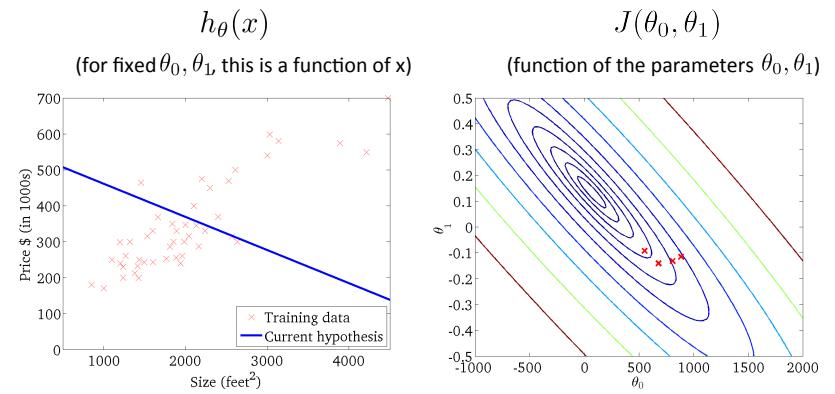
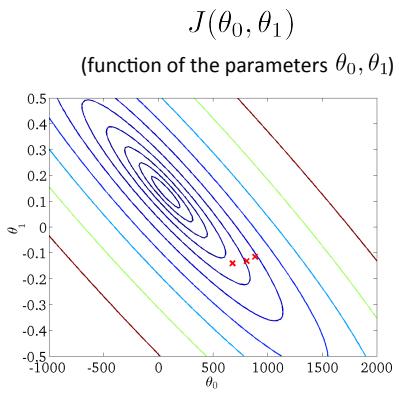
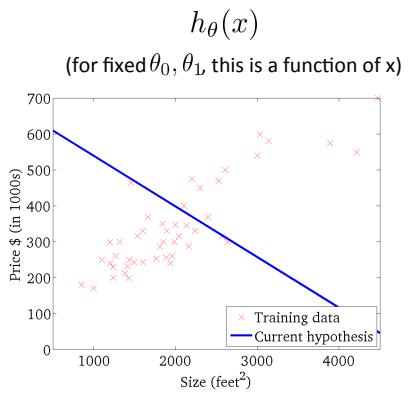
(for fixed θ_0, θ_1 , this is a function of x)

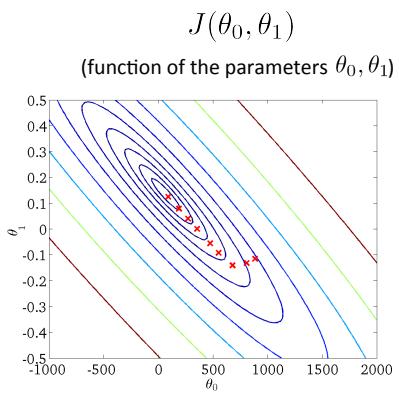
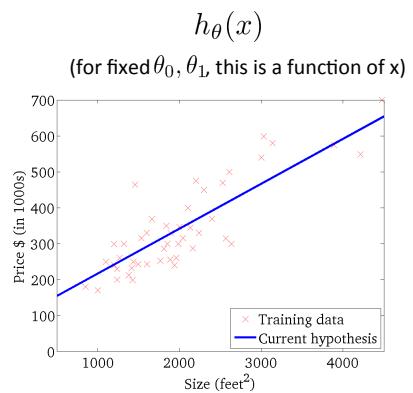
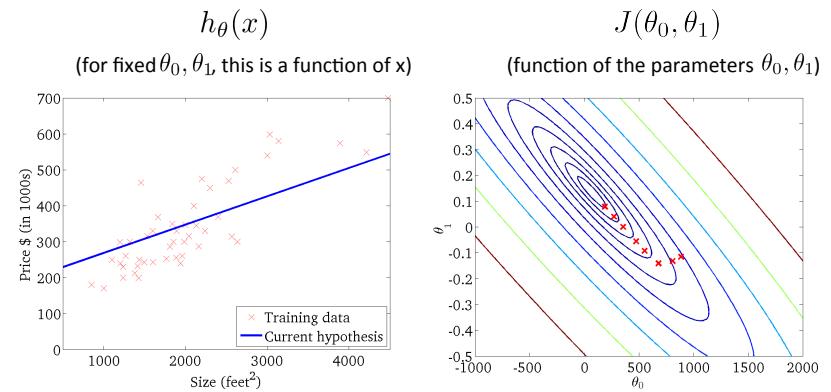
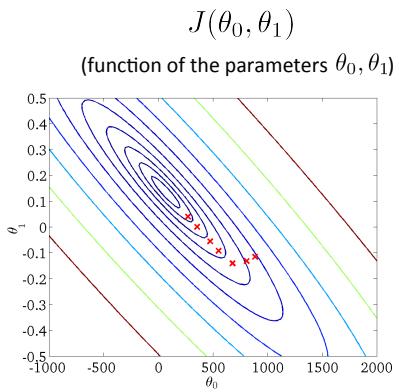
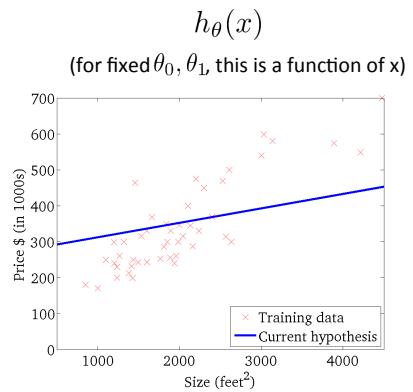


$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)







Which of the following statements are correct?

- In order for the gradient descent to converge, we must gradually decrease over time. *the learning rate α*
- The gradient descent ensures that an overall minimum is reached for any J function $J(\theta_0, \theta_1)$
- The gradient descent can converge even if α it remains constant (although α shouldn't be too large or the algorithm may diverge)
- For function $J(\theta_0, \theta_1)$ associated with linear regression, there are no local optimum, there is only one global optimal.

Multiple features

Size (feet ²) x_1	Number of bedrooms x_2	Number of floors x_3	Age of home (years) x_4	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

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

Multiple features (variables).

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

Notation:
 n = number of features example: $n=4$ $x_1^{(2)} = 1416$
 $x^{(i)}$ = input (features) of i^{th} training example. $x_2^{(2)} = 3$
 $x_j^{(i)}$ = value of feature j in i^{th} training example. $x_3^{(2)} = 2$
 \vdots

Hypothesis:

$$\text{Previously: } h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$\text{Now: } h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

For convenience of notation, define: $x_0 = 1$

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

$$\theta^T \cdot x = [\theta_0 \ \theta_1 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

Multivariate linear regression.

$$\text{Hypothesis: } h_{\theta}(x) = \underline{\theta^T x} = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

Parameters: $\theta_0, \theta_1, \dots, \theta_n \rightarrow \theta$ an $n+1$ dimensional vector

Cost function:

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

Gradient descent:

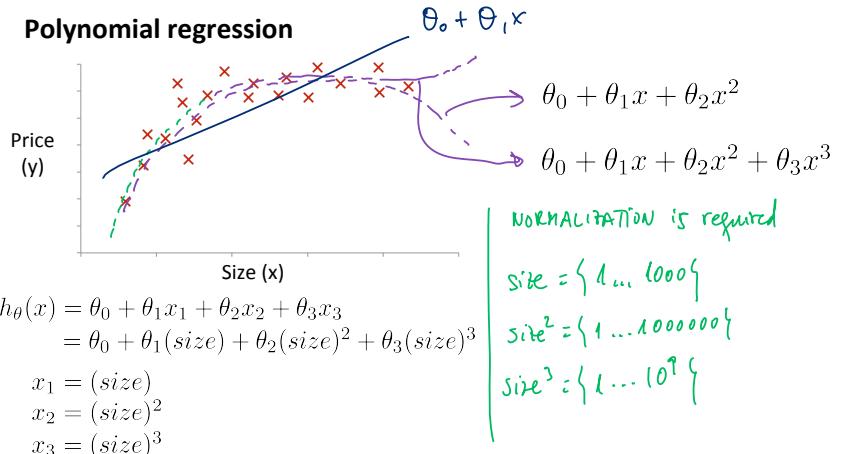
$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \\ \end{array} \right. \quad \text{(simultaneously update for every } j = 0, \dots, n \text{)}$$

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

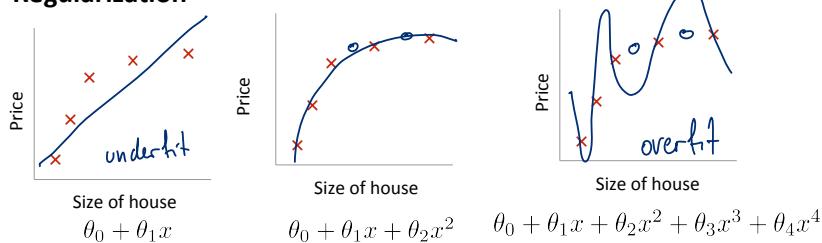
Summary:

- If α is too small: slow convergence.
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge.

To choose α , try $\dots, 0.001, \frac{0.003}{x}, \frac{0.01}{x}, \frac{0.03}{x}, \frac{0.1}{x}, \dots, 1, \dots$



Regularization



Assume we penalize the size of θ_3, θ_4

$$\min_{\theta} \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2 \right]$$

Regularization

Smaller values for $\theta_0, \theta_1, \dots, \theta_n$

- "simpler" hypothesis
- Less prone to overfitting

Example:

- Variables: x_1, x_2, \dots, x_{100}
- Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

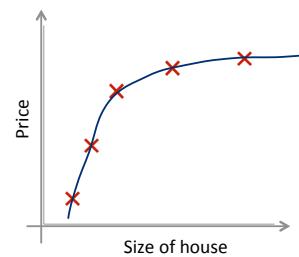
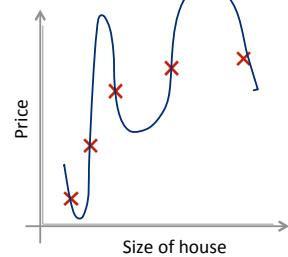
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \cdot \sum_{j=1}^n \theta_j^2 \right]$$

other option + $\lambda \cdot \sum_{j=1}^n |\theta_j|$

Regularization

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

$$\min_{\theta} J(\theta)$$



In regularized linear regression, we choose θ to minimize

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

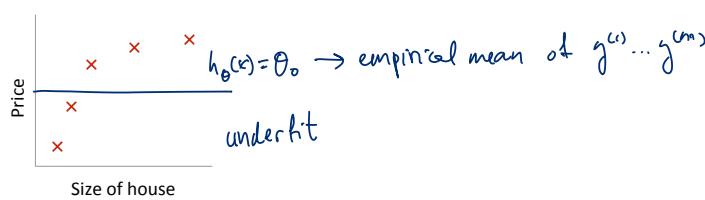
When λ is set to an extremely high value, for example, $\lambda = 10^{10}$ what happens?

- The algorithm works well, increasing λ can't hurt it.
- The algorithm fails to eliminate overfitting.
- The algorithm suffers underfitting. (It can't even adjust the training data.)
- The descent of the gradient doesn't manage to converge.

In regularized linear regression, we choose θ to minimize

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

When λ is set to an extremely high value, for example, $\lambda = 10^{10}$ what happens? $\theta_1 \approx 0, \theta_2 \approx 0, \theta_3 \approx 0, \theta_4 \approx 0$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Regresión lineal regularizada

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

$$\begin{aligned} \min_{\theta} J(\theta) \\ j=0 \quad \frac{\partial J(\theta)}{\partial \theta_0} &= \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \overset{1}{\underset{x_0^{(i)}}{\sim}} \right] \\ j \neq 0 \quad \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \lambda \theta_j \right] \quad j=1 \dots n \end{aligned}$$

↳ this is the only change

Gradient descent

Repeat {

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

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = \cancel{X}, 1, 2, 3, \dots, n)$$

$$\theta_j := \underbrace{\theta_j (1 - \alpha \frac{\lambda}{m})}_{\text{Reduce } \theta_j} - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{update with gradient descent step}}$$

↓ common factor

$$\begin{cases} m' & \left\{ \begin{array}{l} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m')} \\ y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m')} \end{array} \right\} \rightarrow \text{trainning set} \\ t & \left\{ \begin{array}{l} x^{(m'+1)} \\ \vdots \\ x^{(t)} \\ y^{(m'+1)} \\ \vdots \\ y^{(t)} \end{array} \right\} \rightarrow \text{test set} \end{cases} \quad m' \gg t$$

$$\text{trainning set } (x^{(i)}, y^{(i)}) \quad i = 1 \dots m'$$

$$\text{test set } (x^{(i)}, y^{(i)}) \quad i = m'+1 \dots m$$

For a given λ

$$J^{\lambda}_{\text{train}}(\theta) = \frac{1}{2m'} \left(\sum_{i=1}^{m'} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right)$$

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} J^{\lambda}_{\text{train}}(\theta)$$

$$MSE_{\text{train}}(\lambda) = \frac{1}{m'} \left(\sum_{i=1}^{m'} (h_{\theta^*}(x^{(i)}) - y^{(i)})^2 \right)$$

$$MSE_{\text{test}}(\lambda) = \frac{1}{t} \left(\sum_{i=m'+1}^t (h_{\theta^*}(x^{(i)}) - y^{(i)})^2 \right)$$

Regularization factor adjustment

$$\lambda^* = \underset{\lambda}{\operatorname{arg\,min}} MSE_{\text{test}}(\lambda)$$

