

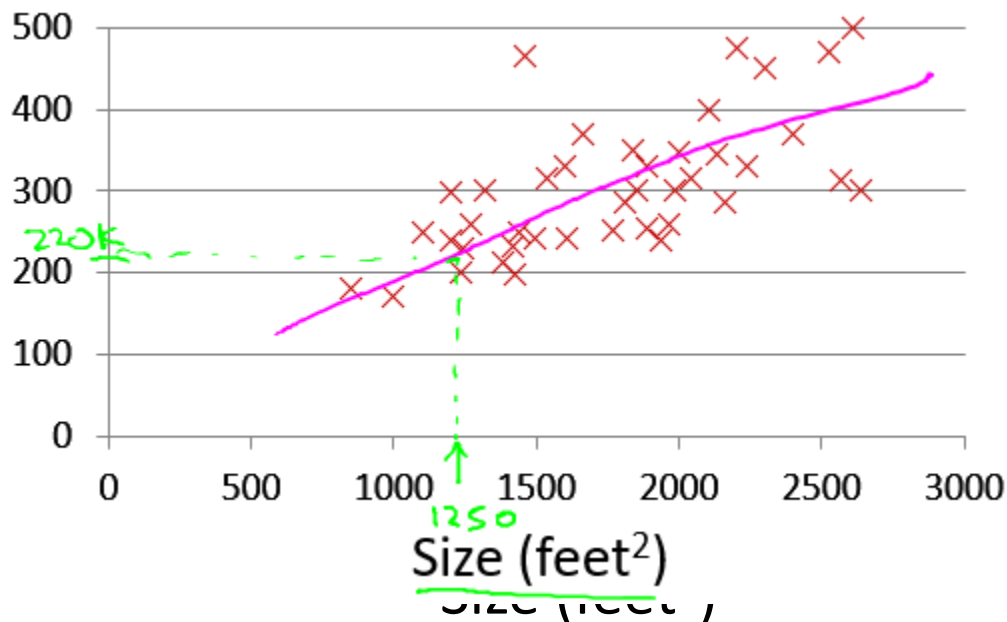
Machine Learning

Linear regression
with one variable

Model
representation

Housing Prices (Portland, OR)

Price
(in 1000s
of dollars)



Supervised Learning

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

Regression Problem

Predict real-valued output

Classification: Discrete-valued output

Training set of housing prices (Portland, OR)

Size in feet ² (x)	Price (\$) in 1000's (y)
→ 2104	460
1416	232
→ 1534	315
852	178
...	...

$m = 42$

Notation:

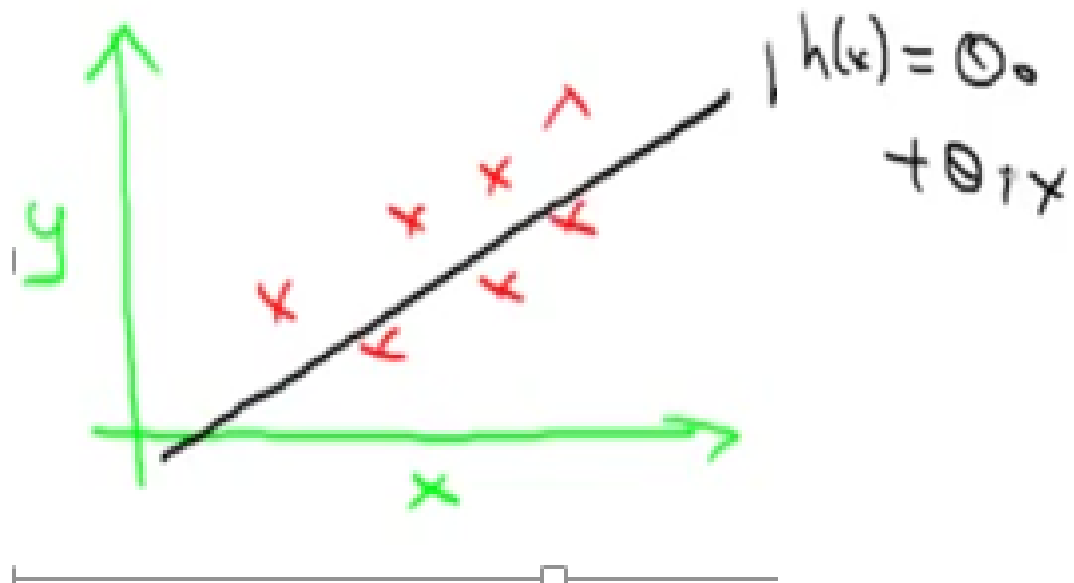
m = Number of training examples

x 's = "input" variable / features

y 's = "output" variable / "target" variable

(x, y) - one training example
 $(x^{(i)}, y^{(i)})$ - i^{th} training example

$$\begin{cases} x^{(1)} = 2104 \\ x^{(2)} = 1416 \\ y^{(1)} = 460 \end{cases}$$



Training Set

Learning Algorithm

Size of house

h

Estimated price

x

hypothesis

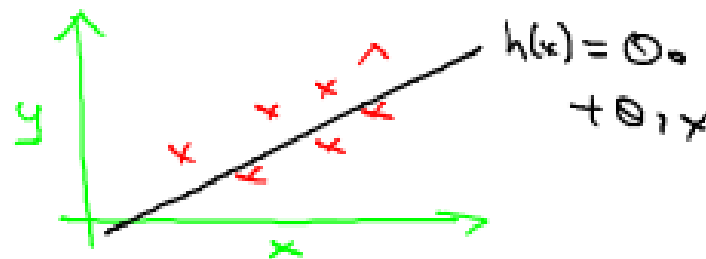
price (estimated value of y)

h maps from x 's to y 's

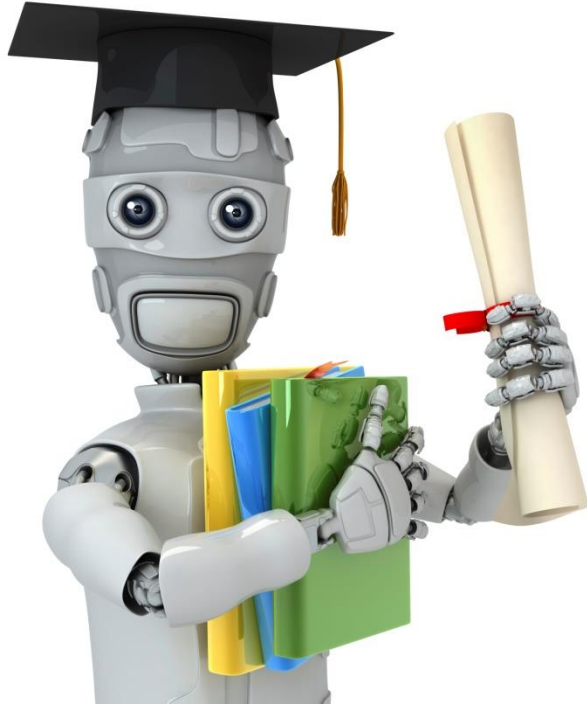
How do we represent h ?

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

Shorthand: $h(x)$



Linear regression with one variable.
Univariate linear regression.



Machine Learning

Linear regression
with one variable

Cost function

Training Set

Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

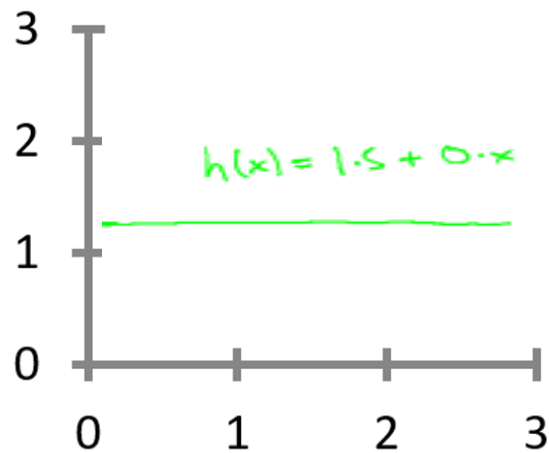
m = 47

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

θ_i 's: Parameters

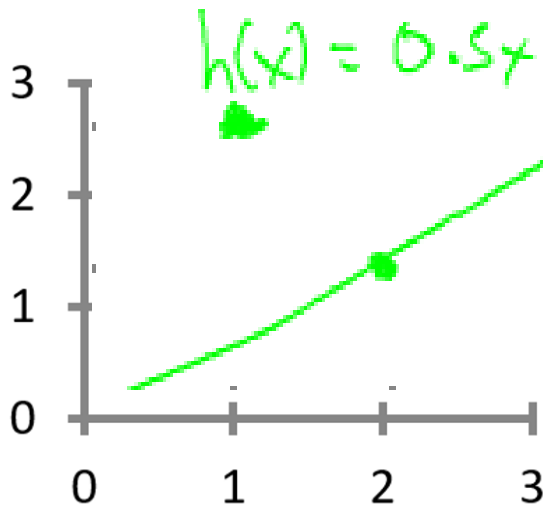
How to choose θ_i 's ?

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



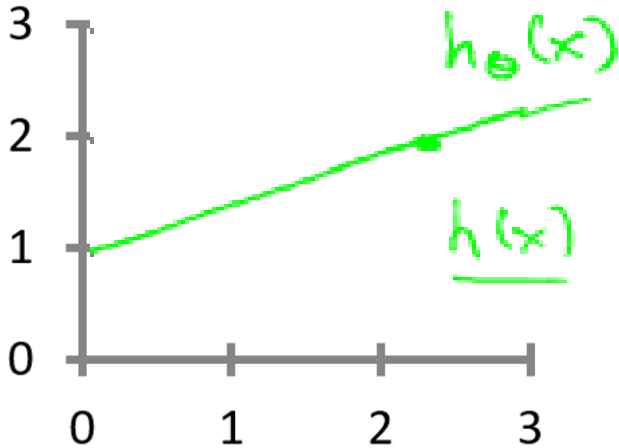
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$



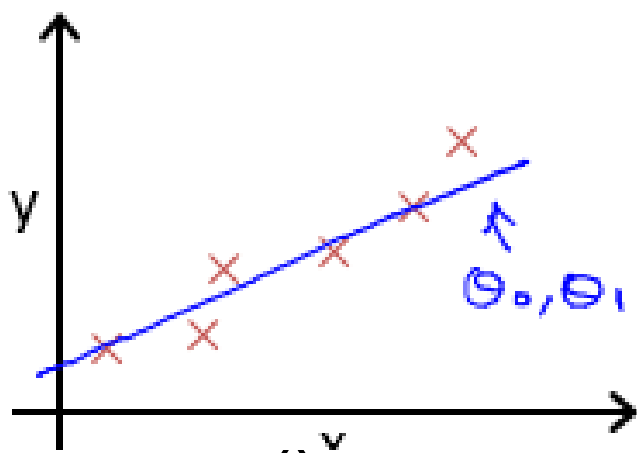
$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

Idea: Choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples (x, y)

Regression is minimization problem

$$\text{minimize}_{\theta_0, \theta_1} (h_{\theta}(x) - y)^2$$



Idea: Choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples (x, y)

$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m \left(\underbrace{h_{\theta}(x^{(i)})}_{\substack{\text{# training examples} \\ h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}}} - y^{(i)} \right)^2$$

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

minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1 Cost function
 Squared error function

Hypothesis:

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

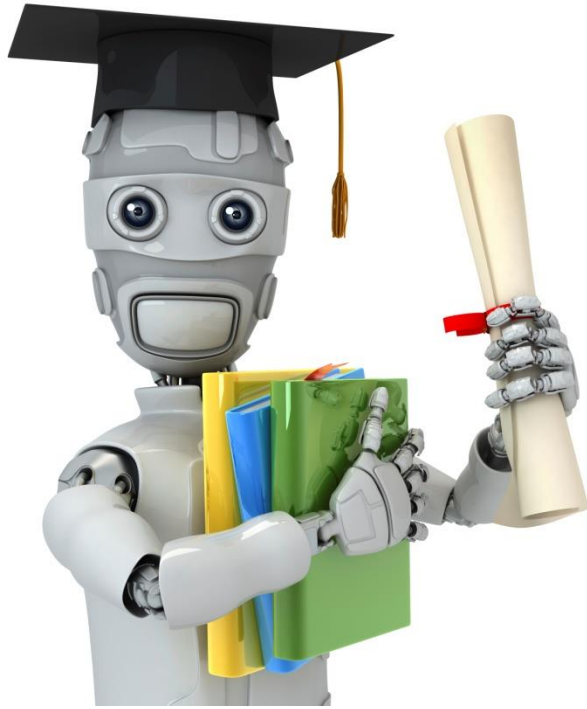
Parameters:

$$\theta_0, \theta_1$$

Cost Function:

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

Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1



Machine Learning

Linear regression
with one variable

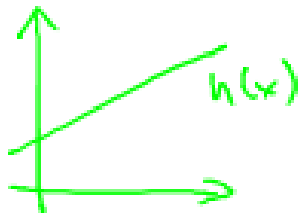
Cost function
intuition I

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$



Cost Function:

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

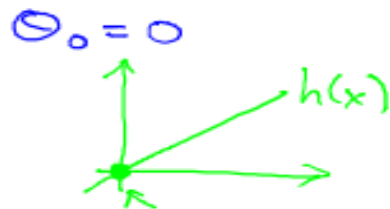
Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

Simplified

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

$$\theta_1$$

$$\theta_1$$

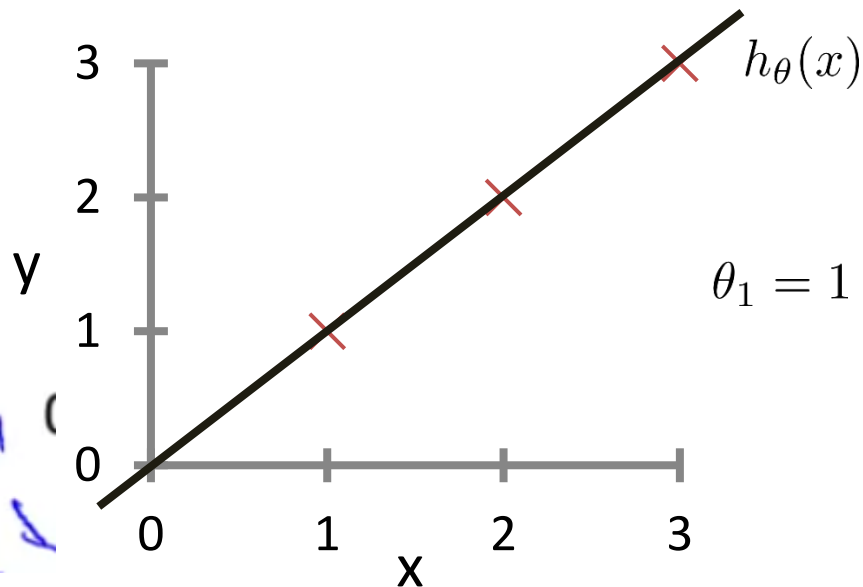


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

minimize $J(\theta_1)$
 θ_1

$$h_{\theta}(x)$$

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

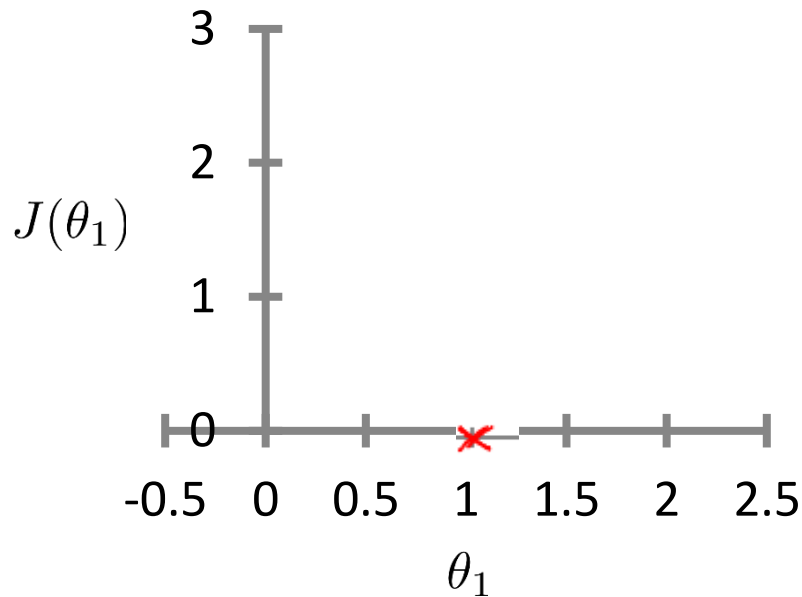


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^m (1 \cdot x^{(i)} - y^{(i)})^2 = \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0^2$$

$$J(\theta_1)$$

(function of the parameter θ_1)

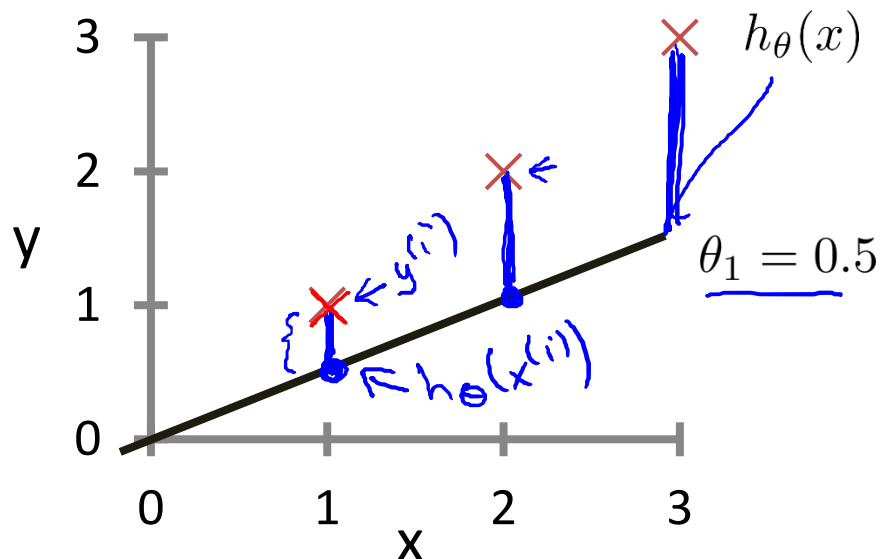


$\theta_1 = 0.5?$

$J(1) = 0$

$$h_{\theta}(x)$$

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

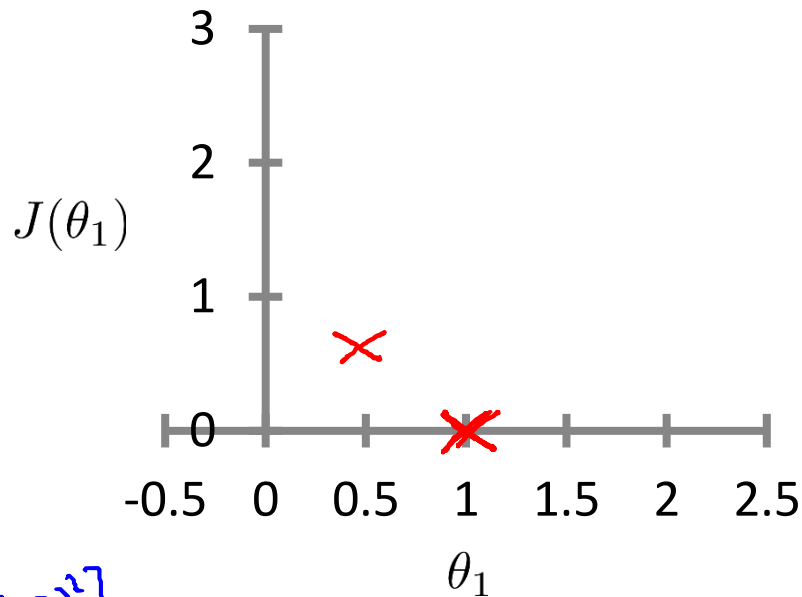


$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2 \times 3} (3.5) = \frac{3.5}{6} \approx \underline{0.58}$$

$$J(\theta_1)$$

(function of the parameter θ_1)

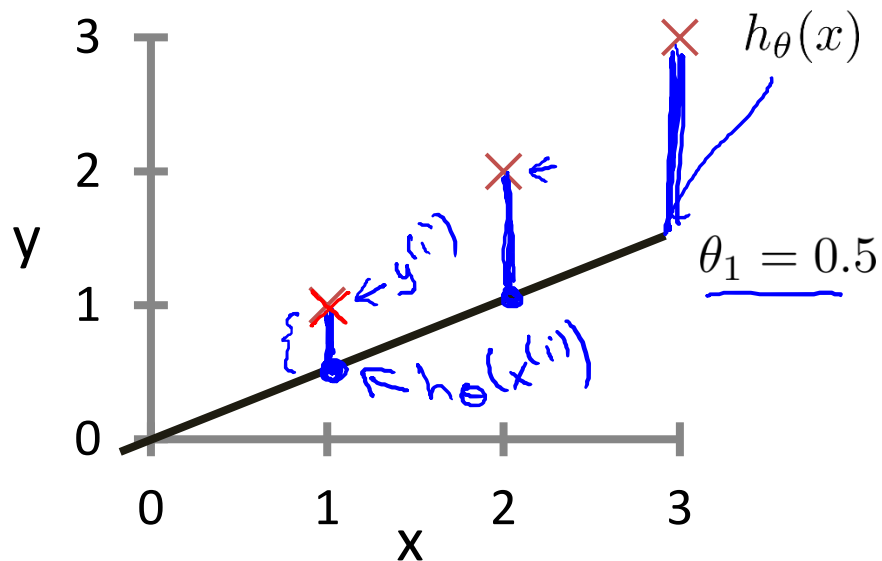


$$\theta_1 = 0?$$

$$J(0) = ?$$

$$h_{\theta}(x)$$

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

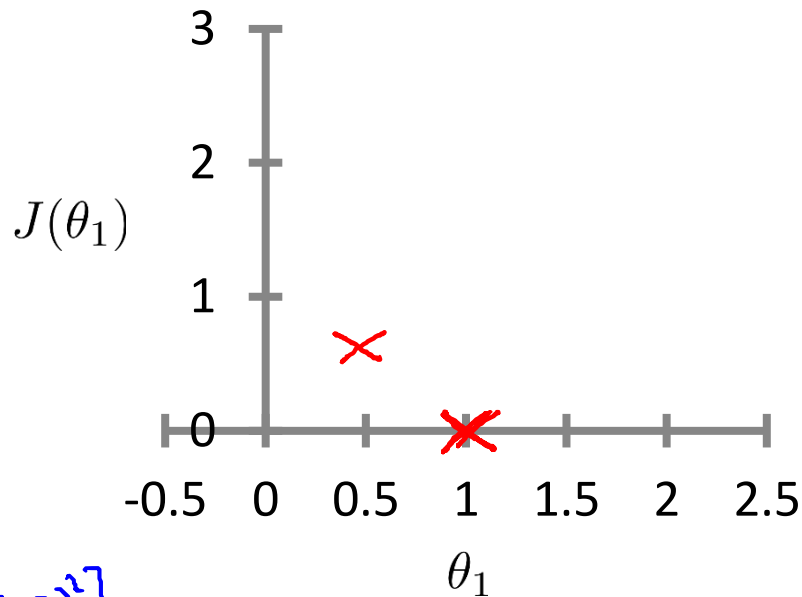


$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2 \times 3} (3.5) = \frac{3.5}{6} \approx \underline{0.58}$$

$$J(\theta_1)$$

(function of the parameter θ_1)

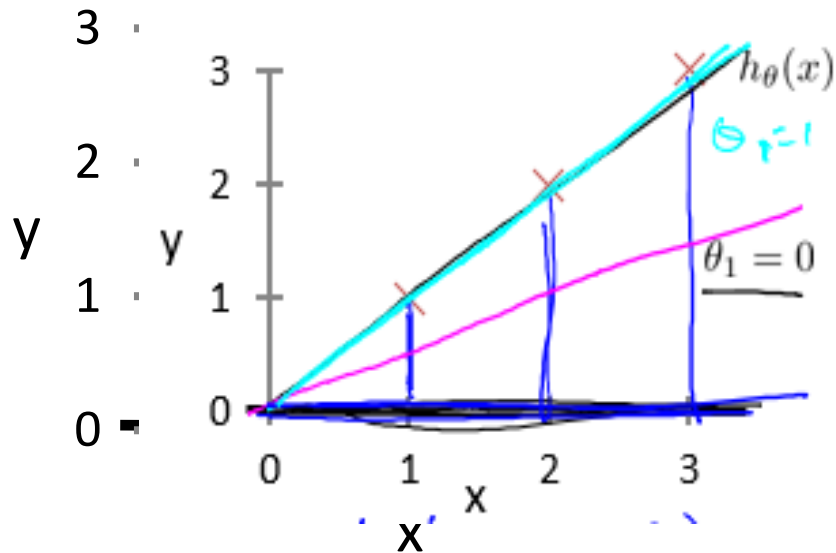


$$\theta_1 = 0?$$

$$J(0) = ?$$

$$h_{\theta}(x)$$

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

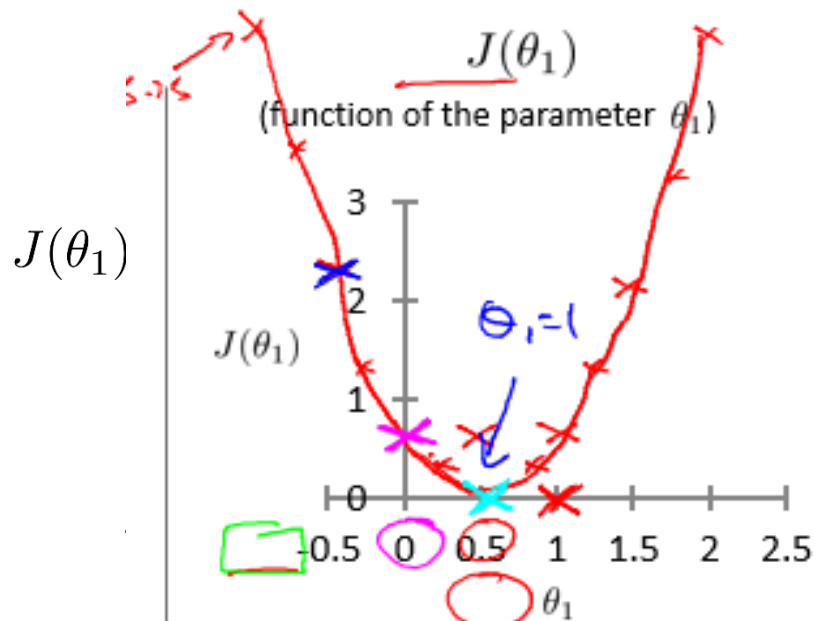


$$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2) = \frac{1}{6} \cdot 14 \approx 2.3$$

$$h(x) = -0.5x$$

$$J(\theta_1)$$

(function of the parameter θ_1)



$$\text{Minimize } J(\theta_1)$$

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

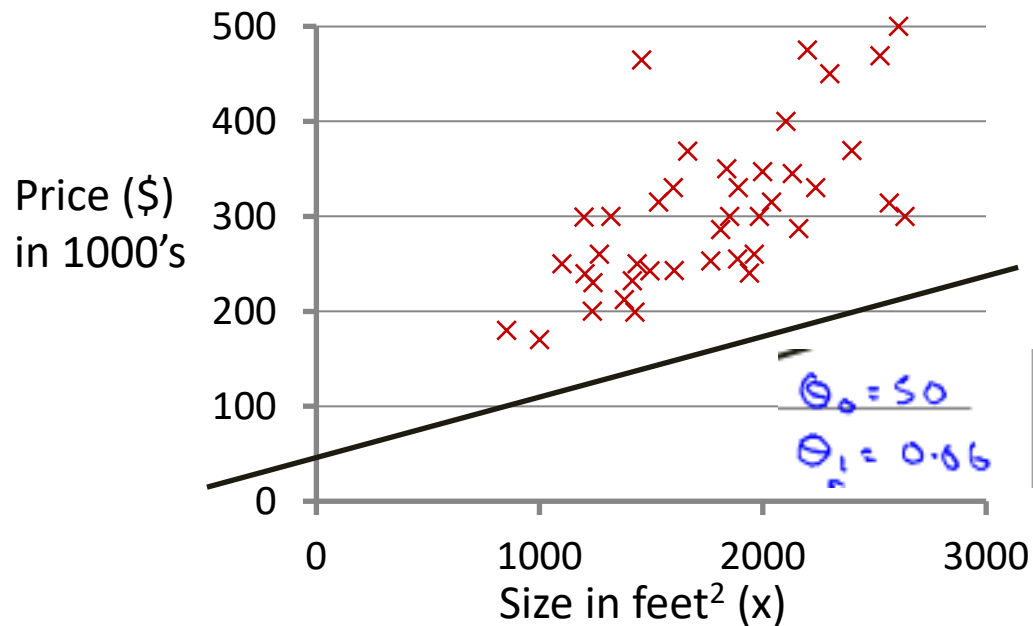
Parameters: θ_0, θ_1

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

Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

$$h_{\theta}(x)$$

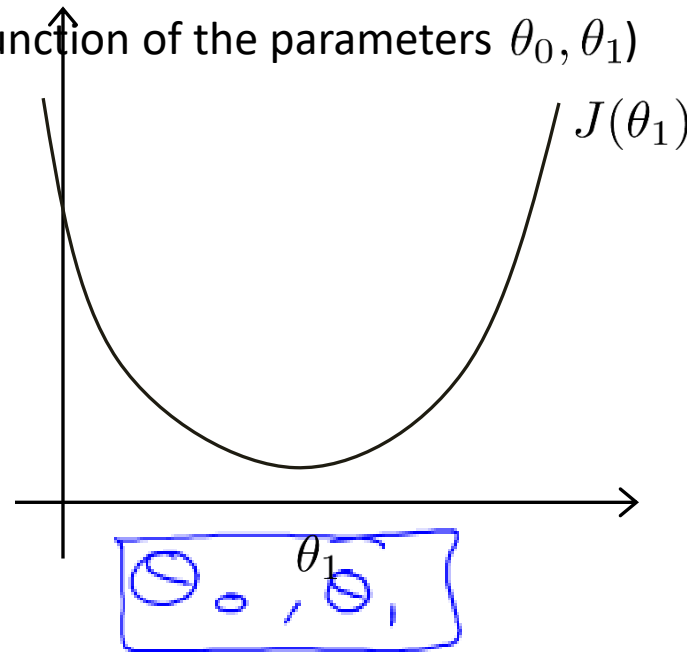
(for fixed θ_0, θ_1 , this is a function of x)



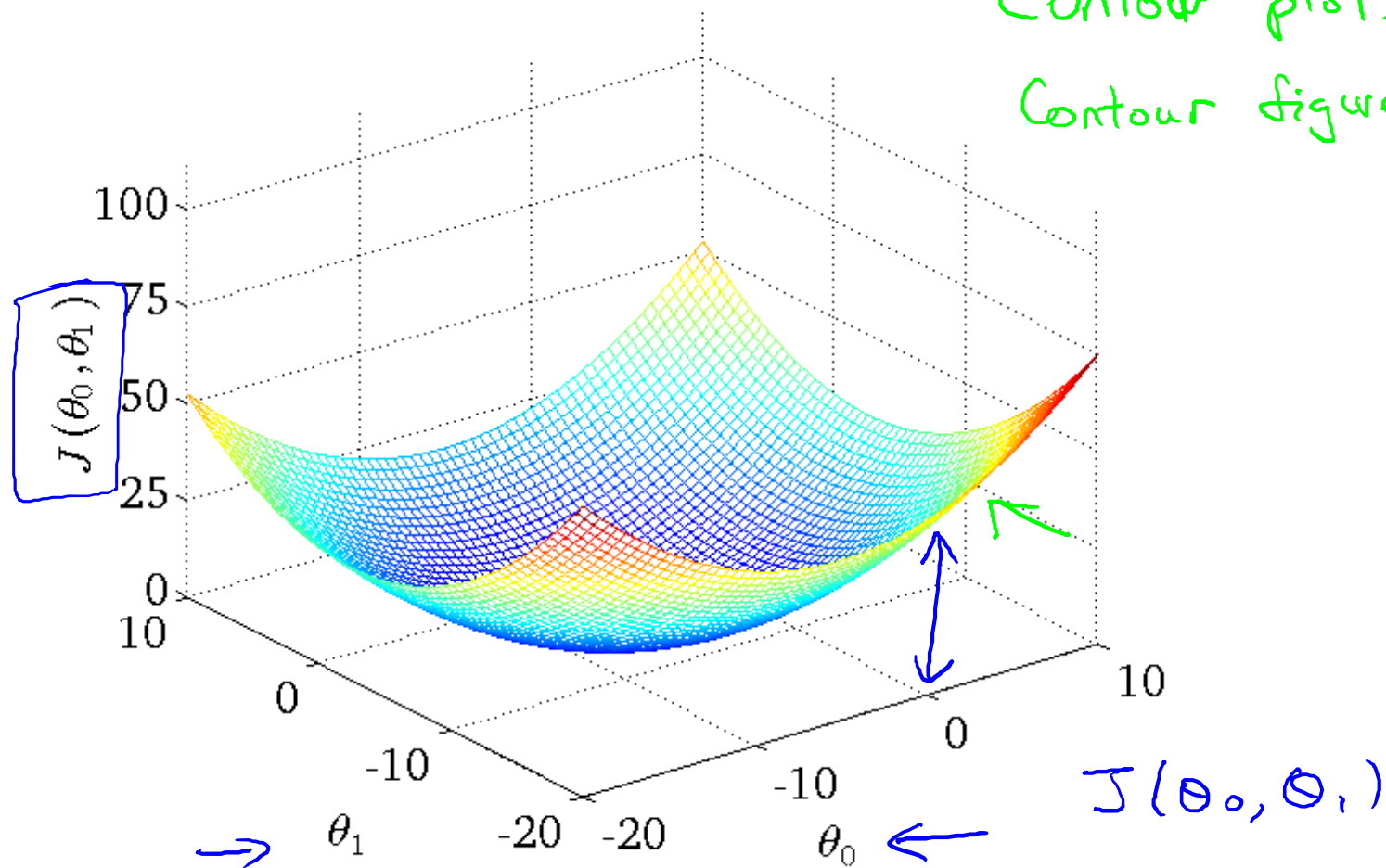
$$h_{\theta}(x) = 50 + 0.06x$$

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

(function of the parameters θ_0, θ_1)

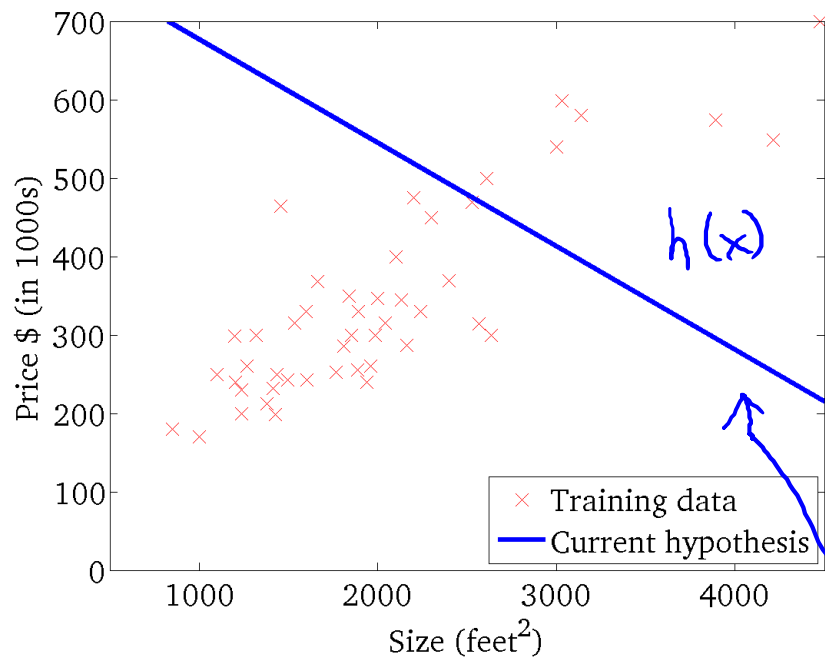


Contour plots
Contour figures -



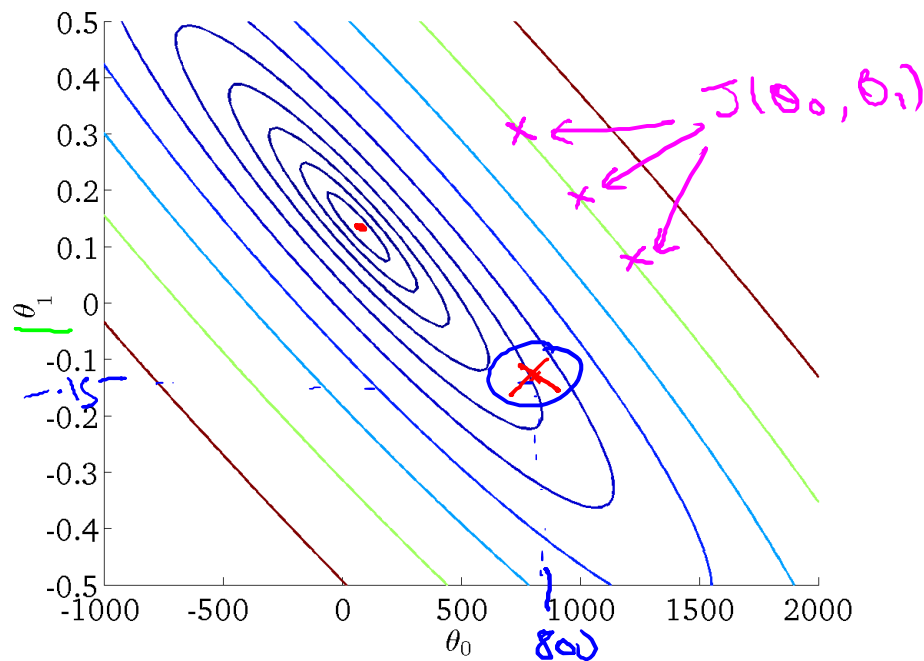
$$h_{\theta}(x)$$

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



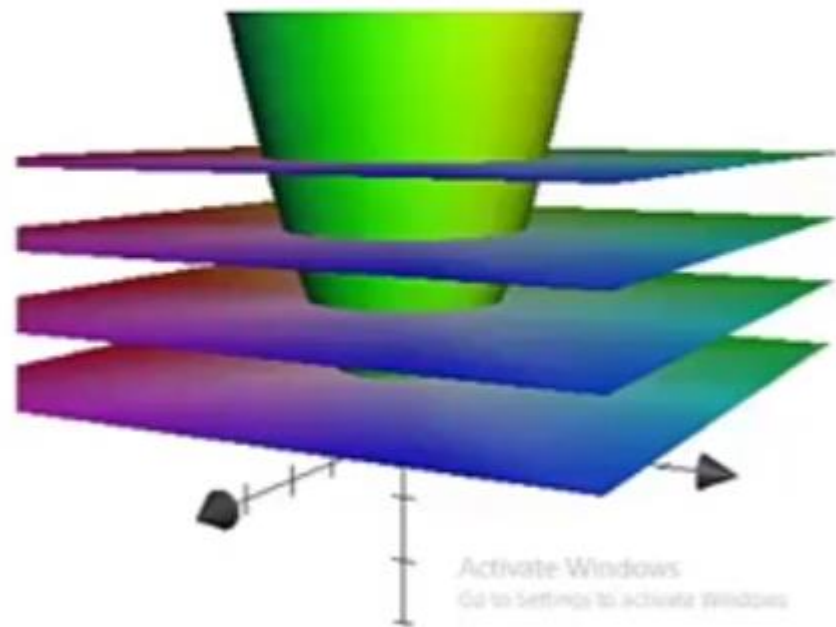
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)

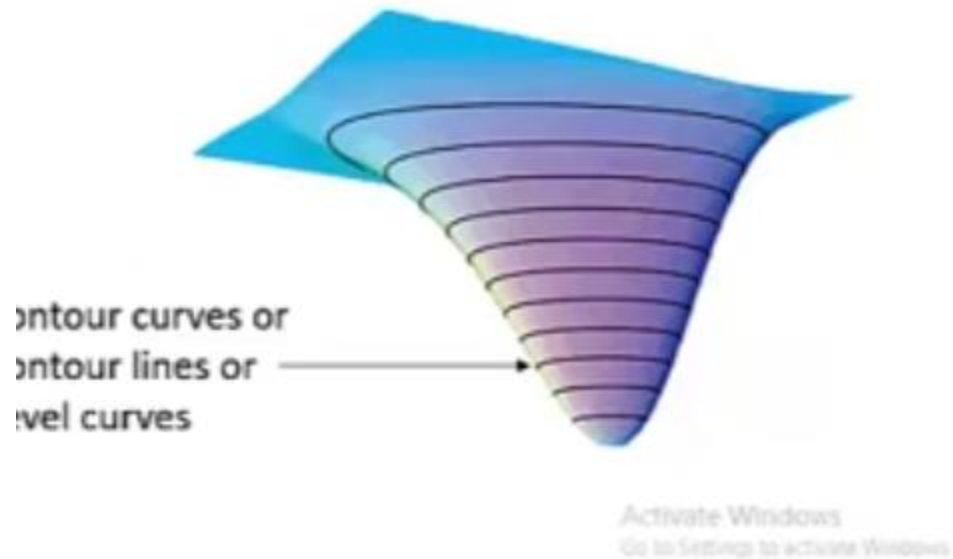


θ_0, θ_1

Cost Function



Cost Function

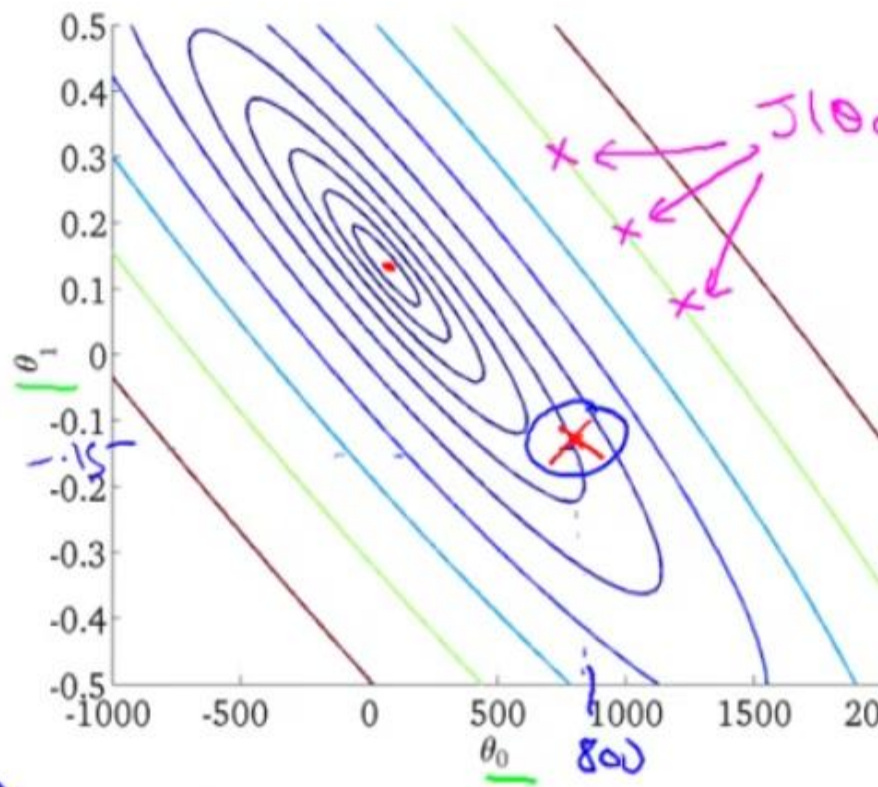
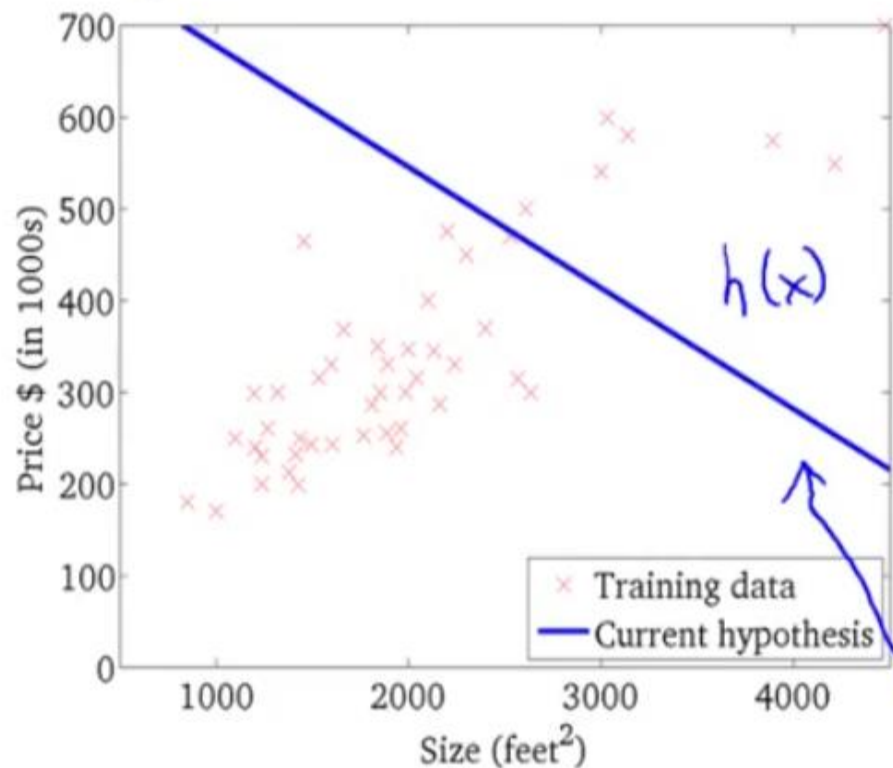


$$h_{\theta}(x)$$

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

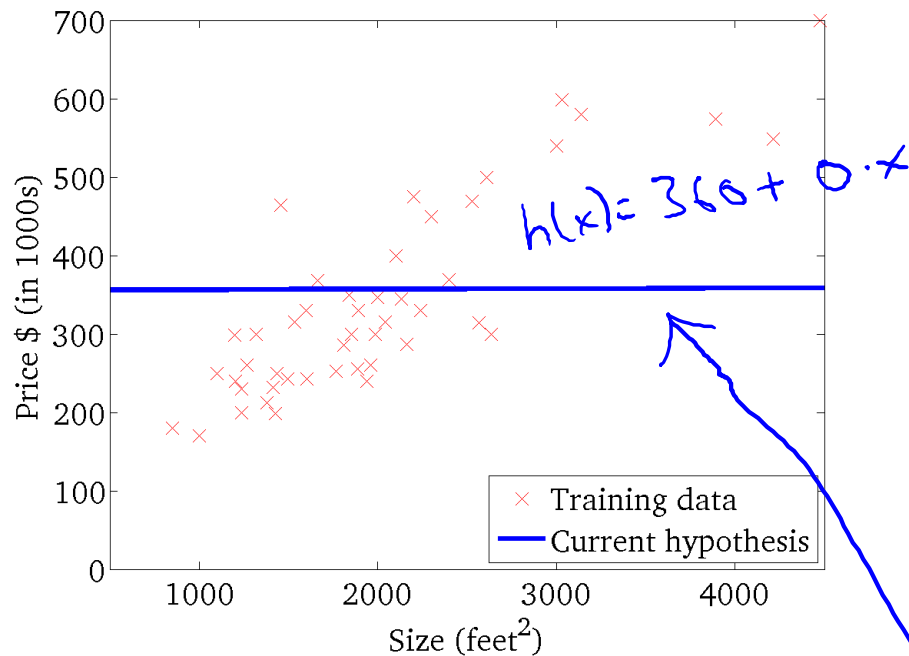
(for fixed θ_0, θ_1 , this is a function of x)

(function of the parameters θ_0, θ_1)



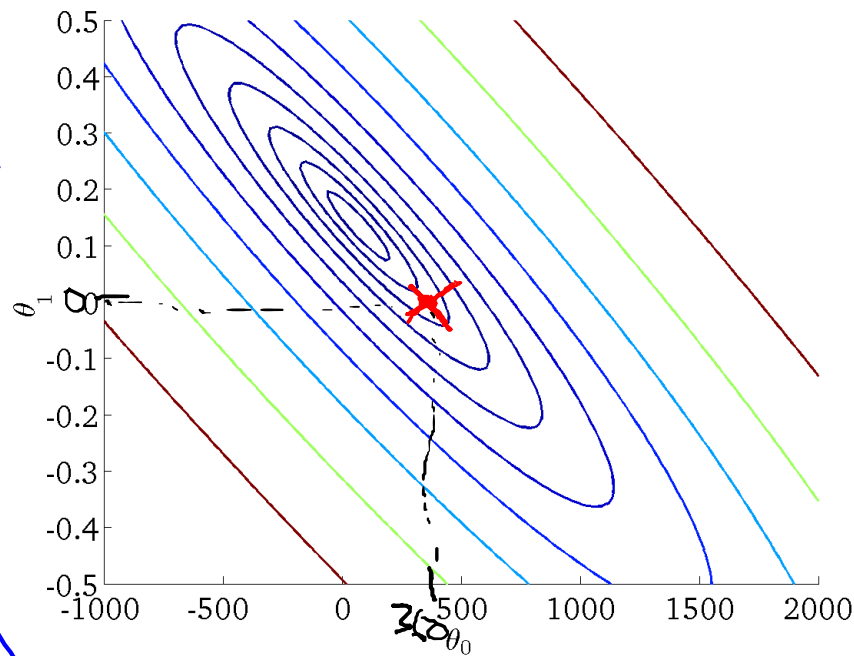
$$h_{\theta}(x)$$

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



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

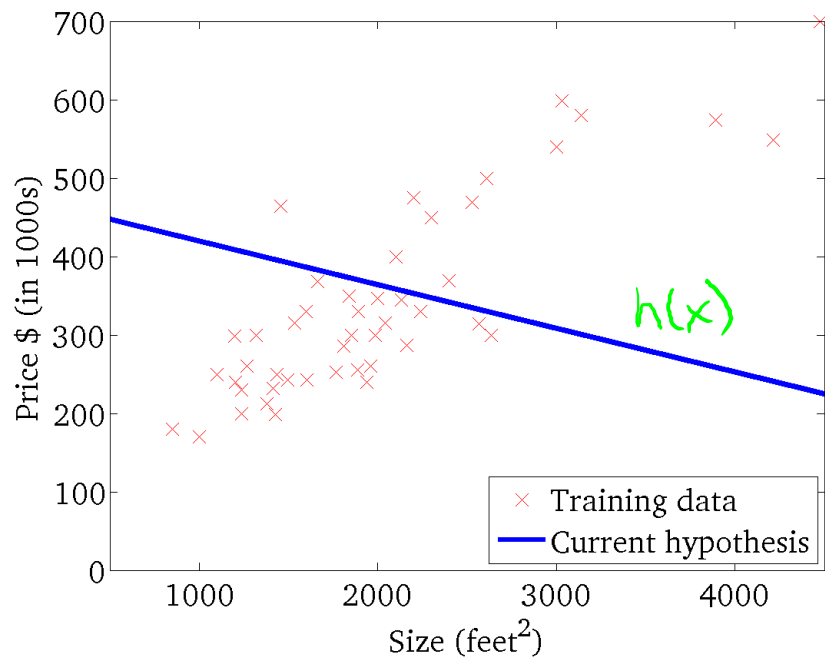
(function of the parameters θ_0, θ_1)



$$\begin{cases} \theta_0 = 360 \\ \theta_1 = 0 \end{cases}$$

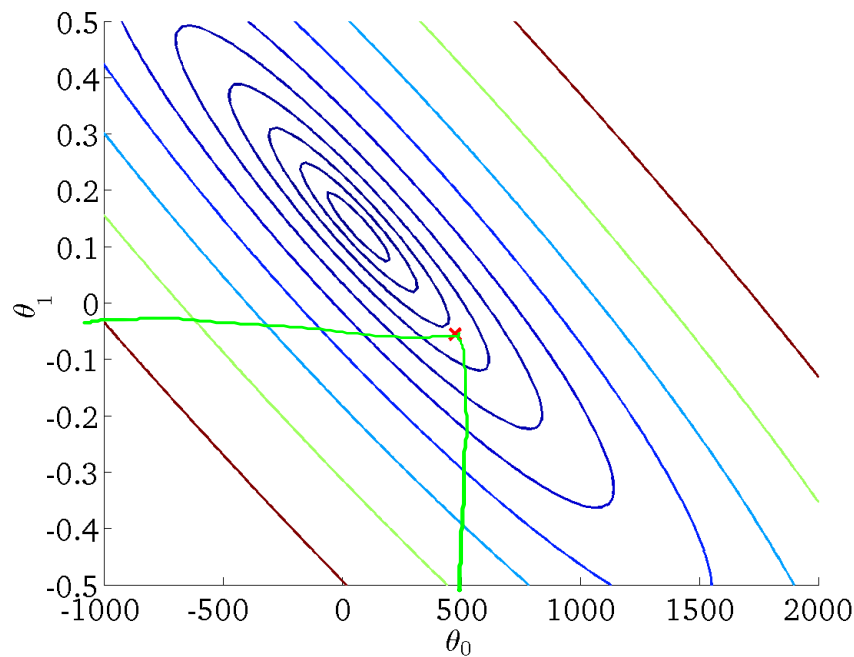
$$h_{\theta}(x)$$

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



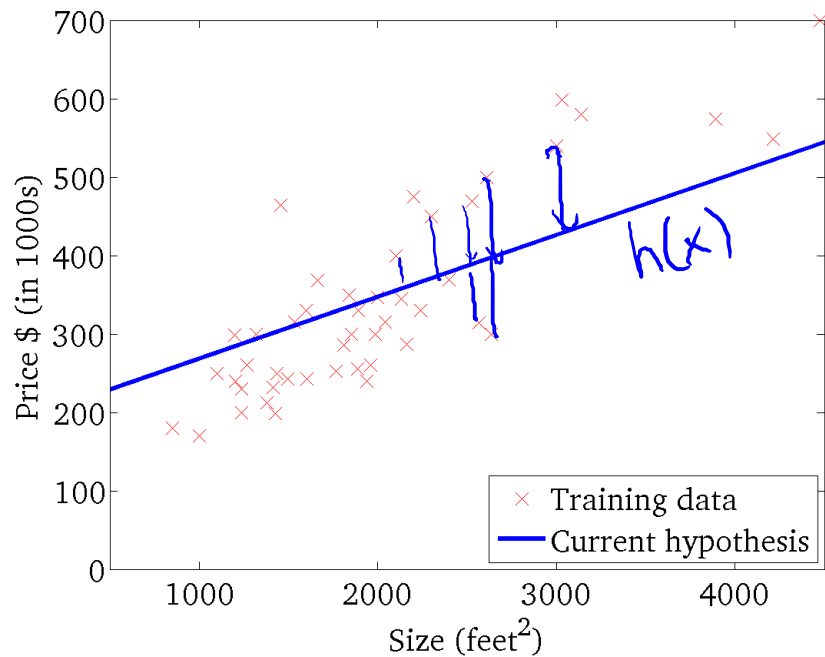
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



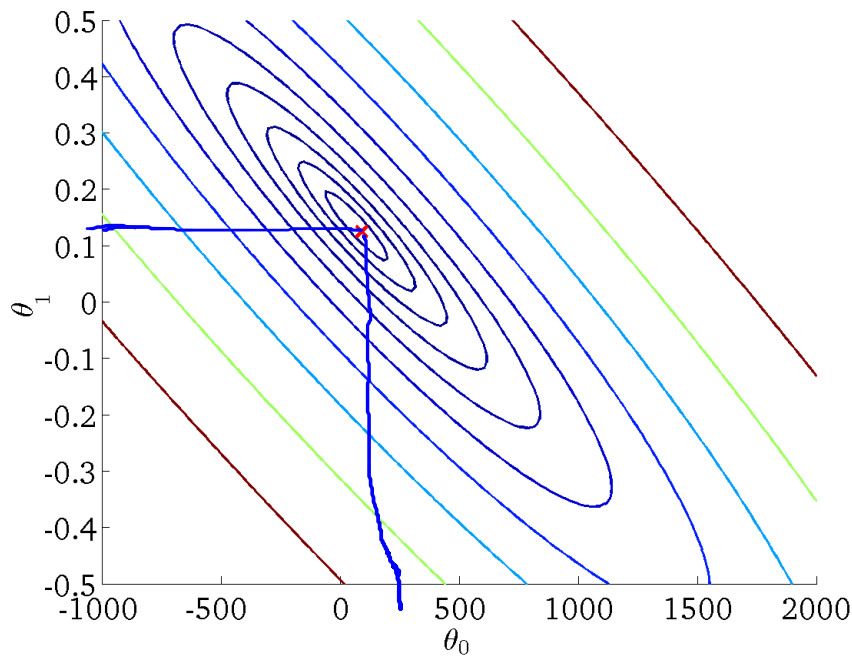
$$h_{\theta}(x)$$

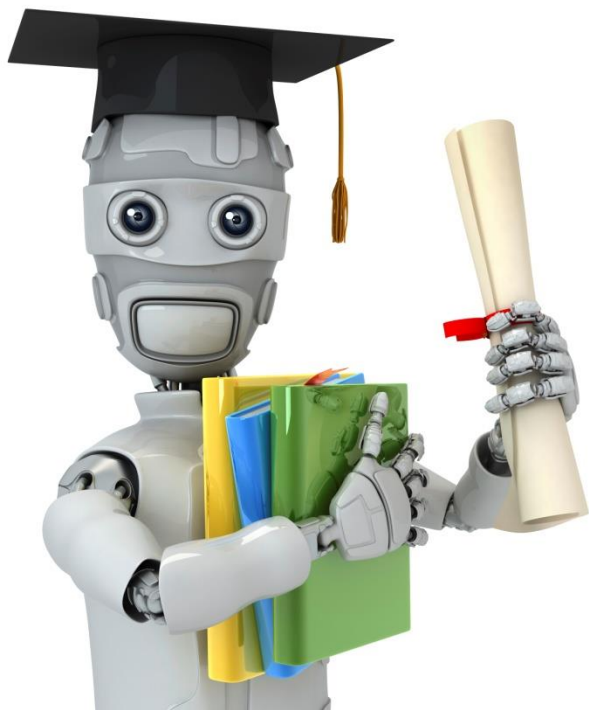
(for fixed θ_0, θ_1 , this is a function of x)



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

(function of the parameters θ_0, θ_1)





Machine Learning

Linear regression
with one variable

Gradient
descent

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

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

Outline:

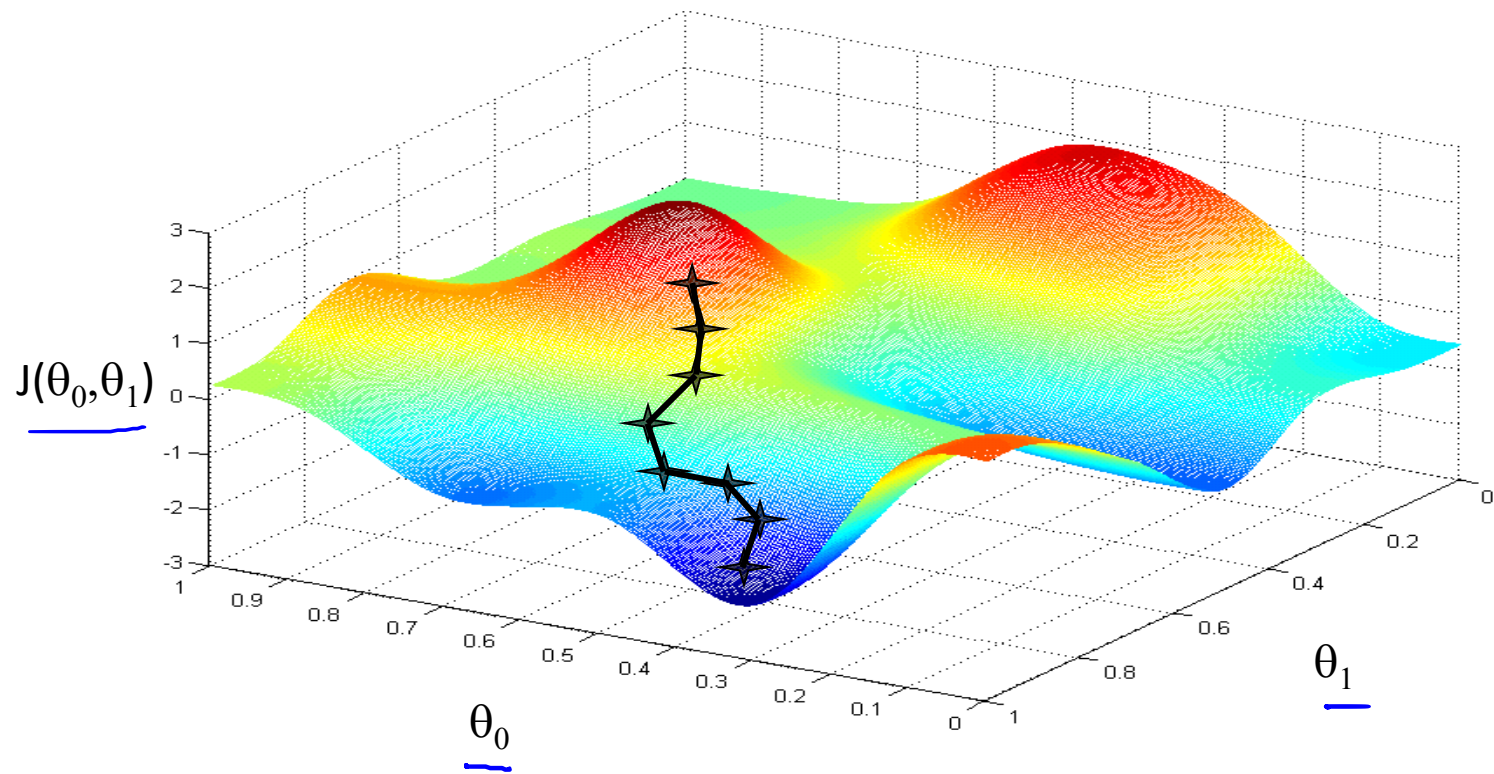
- Start with some θ_0, θ_1
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum

Have some function $J(\theta_0, \theta_1)$ $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

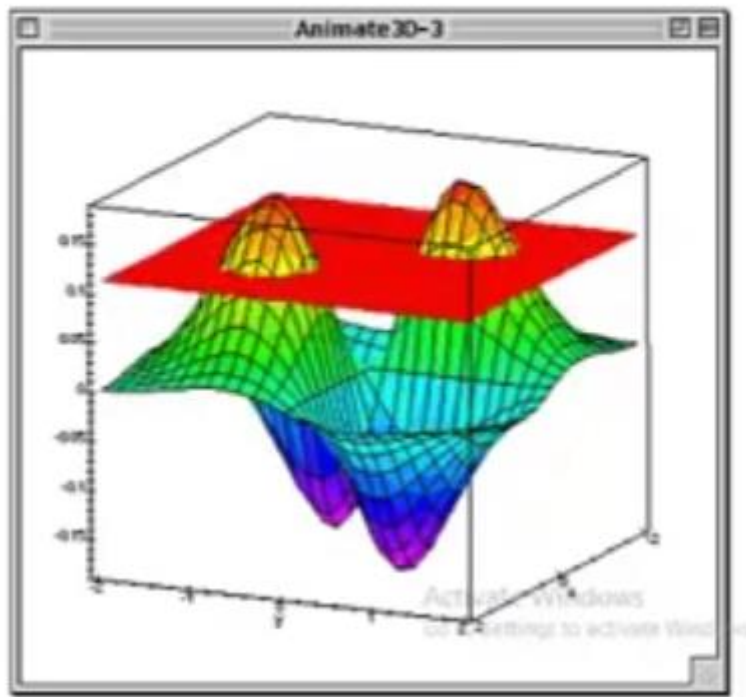
Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ $\min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$

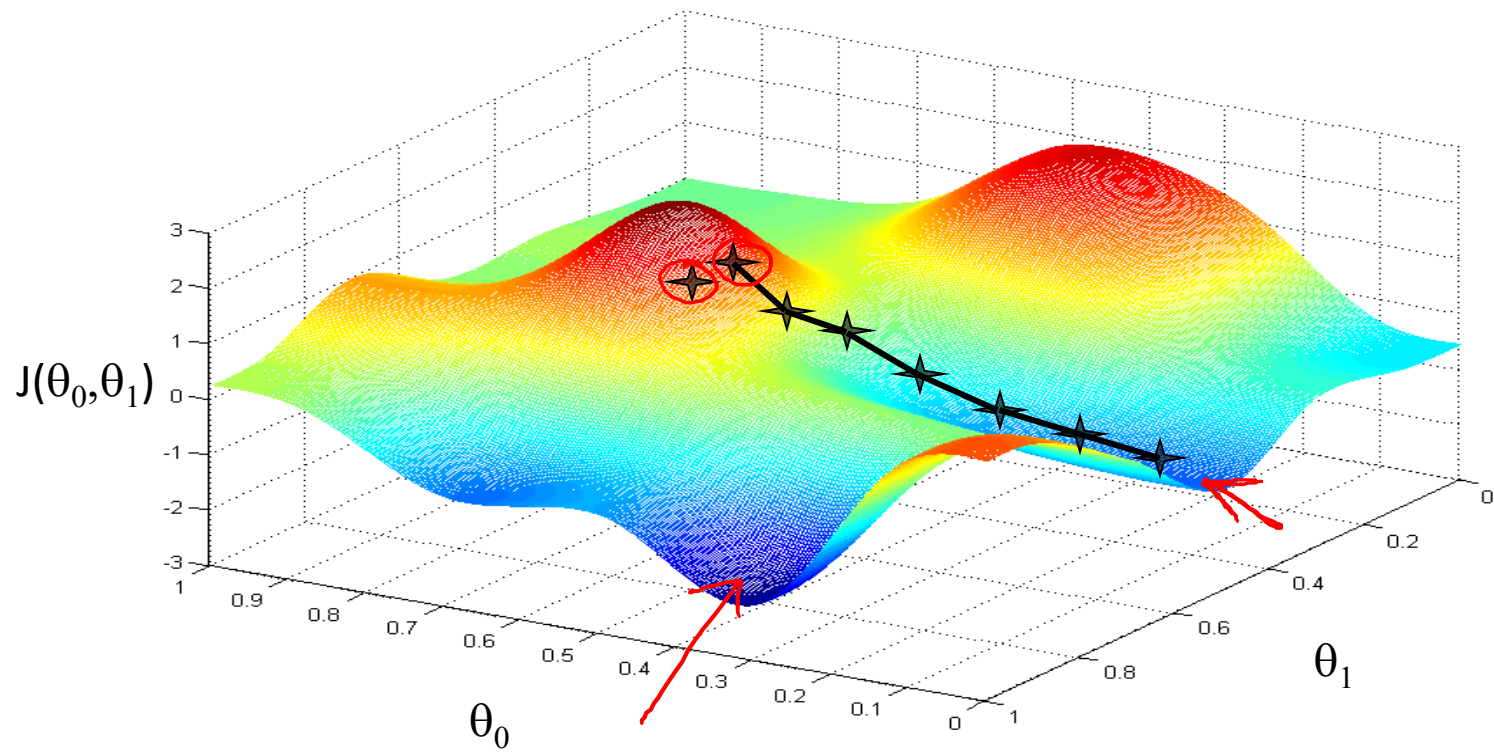
Outline:

- Start with some θ_0, θ_1 (say $\theta_0 = 0, \theta_1 = 0$)
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum



Cost Function





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

Assignment
 $\Rightarrow a := b$
 $\swarrow \quad \searrow$
 $a := a + 1$

Truth assertion
 $a = b \leftarrow$
 $a = a + 1 \times$

Correct: Simultaneous update

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrect:

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp1}$$

Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

learning rate

θ_0, θ_1

(for $j = 0$ and $j = 1$)

Simultaneously update θ_0 and θ_1

Assignment

$$a = b$$

$$a = a + 1$$

Truth assertion

$$a = b$$

$$a = a + 1$$

Correct: Simultaneous update

$$\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\rightarrow \theta_0 := \text{temp0}$$

$$\rightarrow \theta_1 := \text{temp1}$$

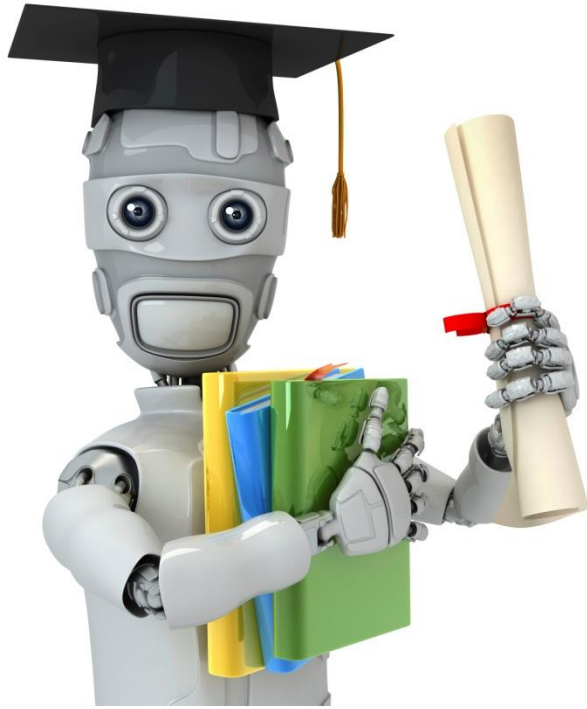
Incorrect:

$$\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\rightarrow \theta_0 := \text{temp0}$$

$$\rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\rightarrow \theta_1 := \text{temp1}$$



Machine Learning

Linear regression
with one variable

Gradient descent
intuition

Gradient descent algorithm

repeat until convergence {

$$\rightarrow \underline{\theta_j} := \underline{\theta_j} - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

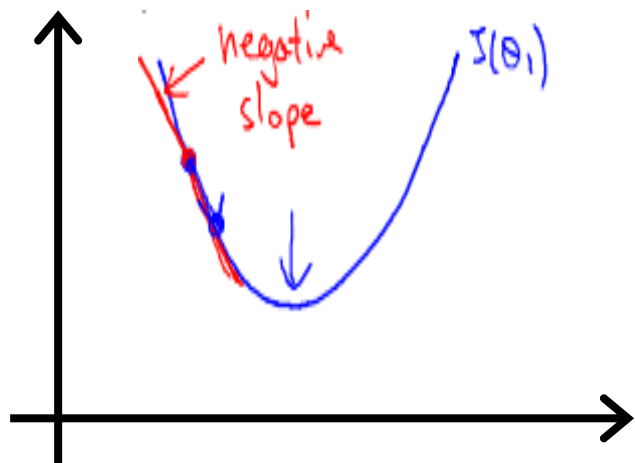
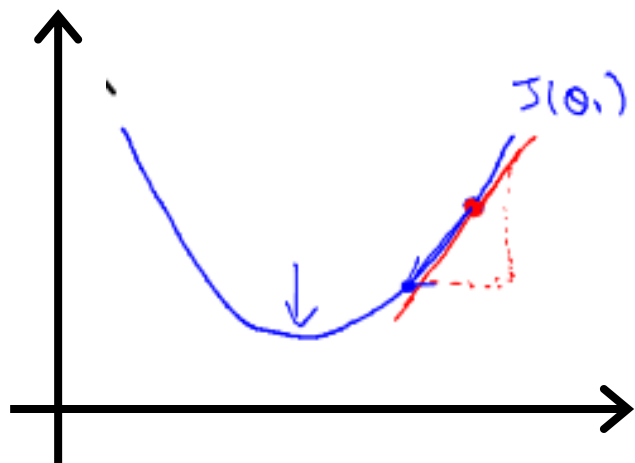
learning
rate

derivative

(simultaneously update
 $j = 0$ and $j = 1$)

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

$$\theta_1 \in \mathbb{R}.$$



$$(\theta_1 \in \mathbb{R})$$

$$\theta_1 := \theta_1 - \underbrace{\alpha}_{\geq 0} \cdot \underbrace{\frac{\partial}{\partial \theta_1} J(\theta_1)}_{\geq 0}$$

$$\theta_1 := \theta_1 - \underline{\alpha} \cdot (\text{positive number})$$

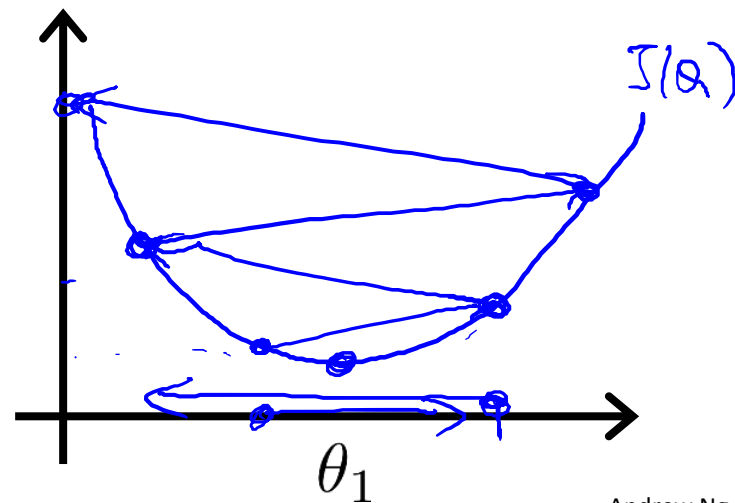
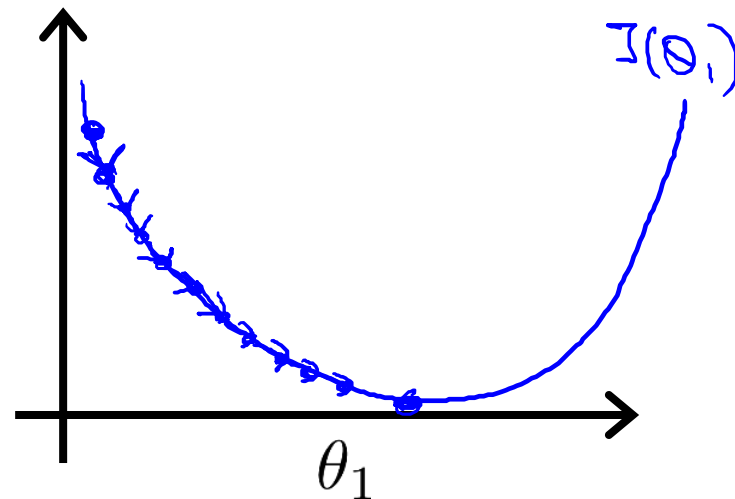
$$\frac{\partial}{\partial \theta_1} J(\theta_1)$$

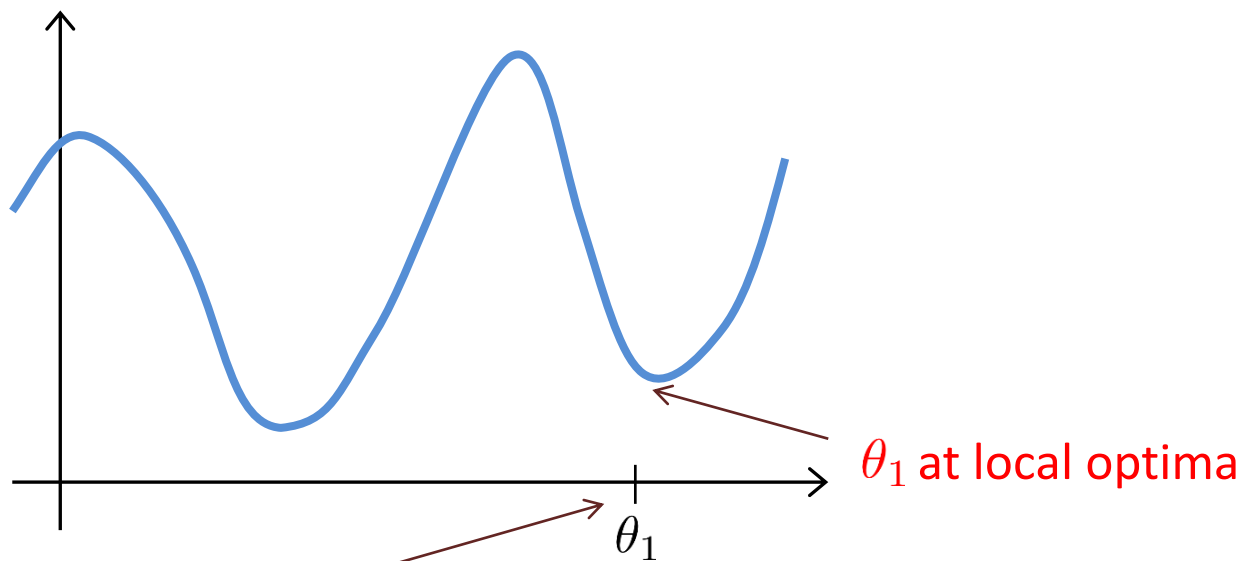
$$\theta_1 := \theta_1 - \underbrace{\alpha}_{\uparrow} (\underbrace{\text{negative number}}_{\uparrow})$$

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

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.





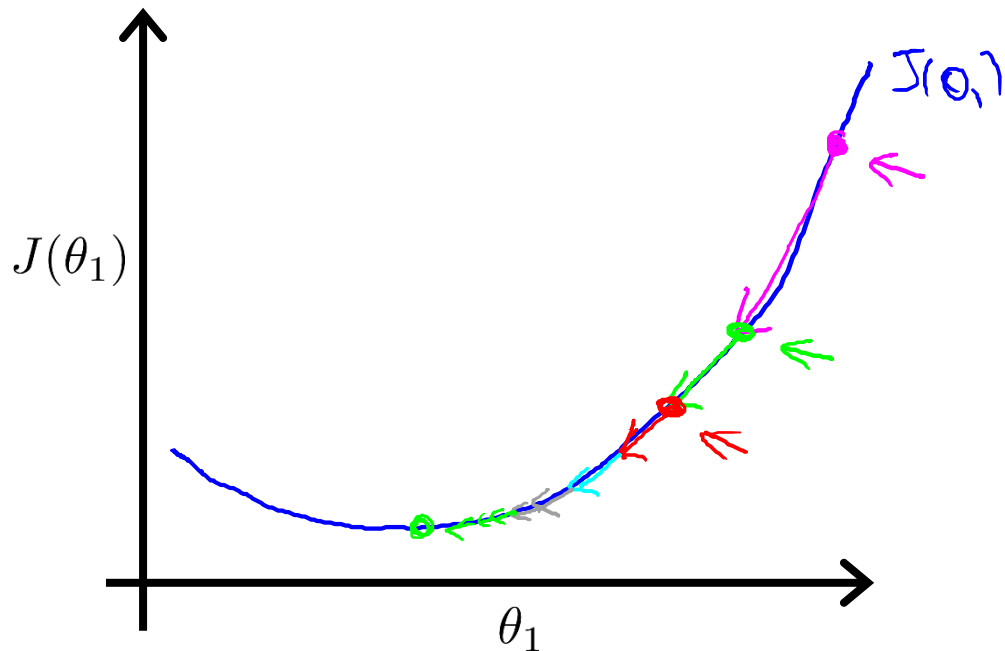
Current value of θ_1

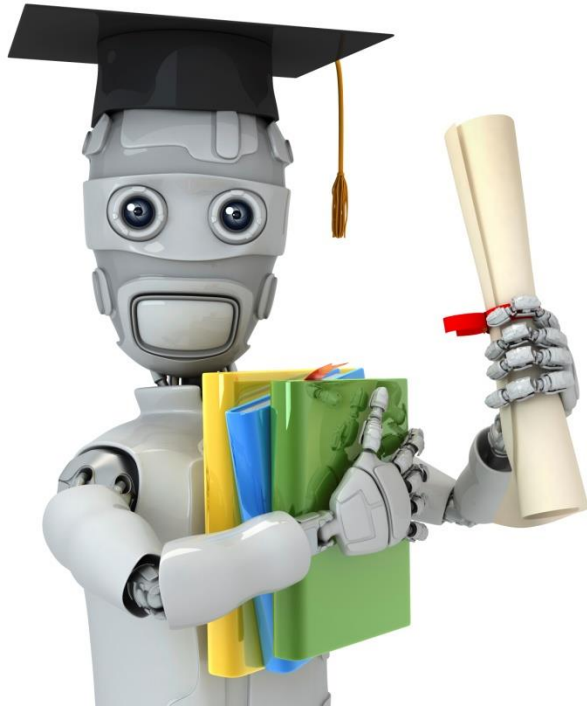
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

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.





Machine Learning

Linear regression with one variable

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

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

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

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

$$= \frac{2}{2\theta_j} \frac{1}{2m} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - y^{(i)})^2$$

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

$$j = 1 : \underline{\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

Gradient descent algorithm

repeat until convergence {

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

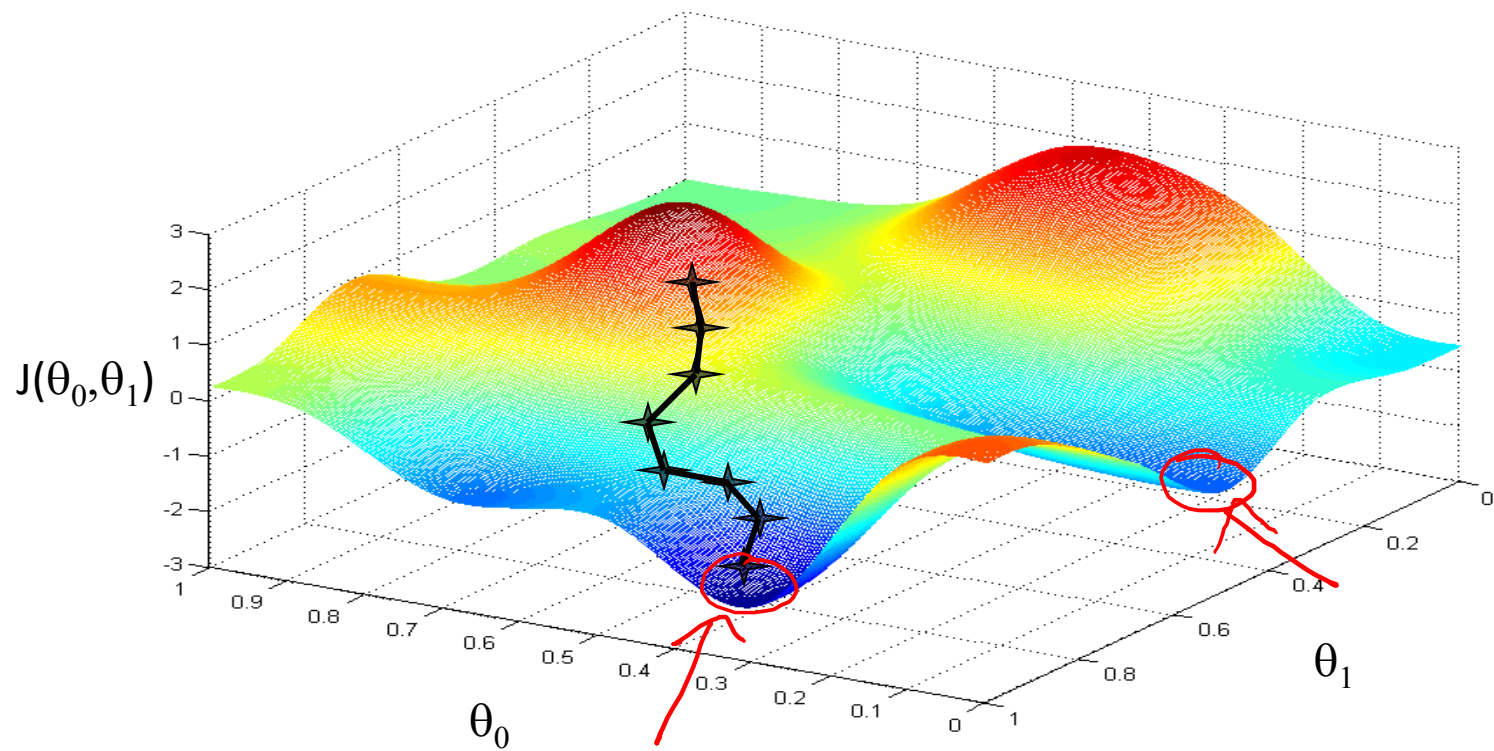
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

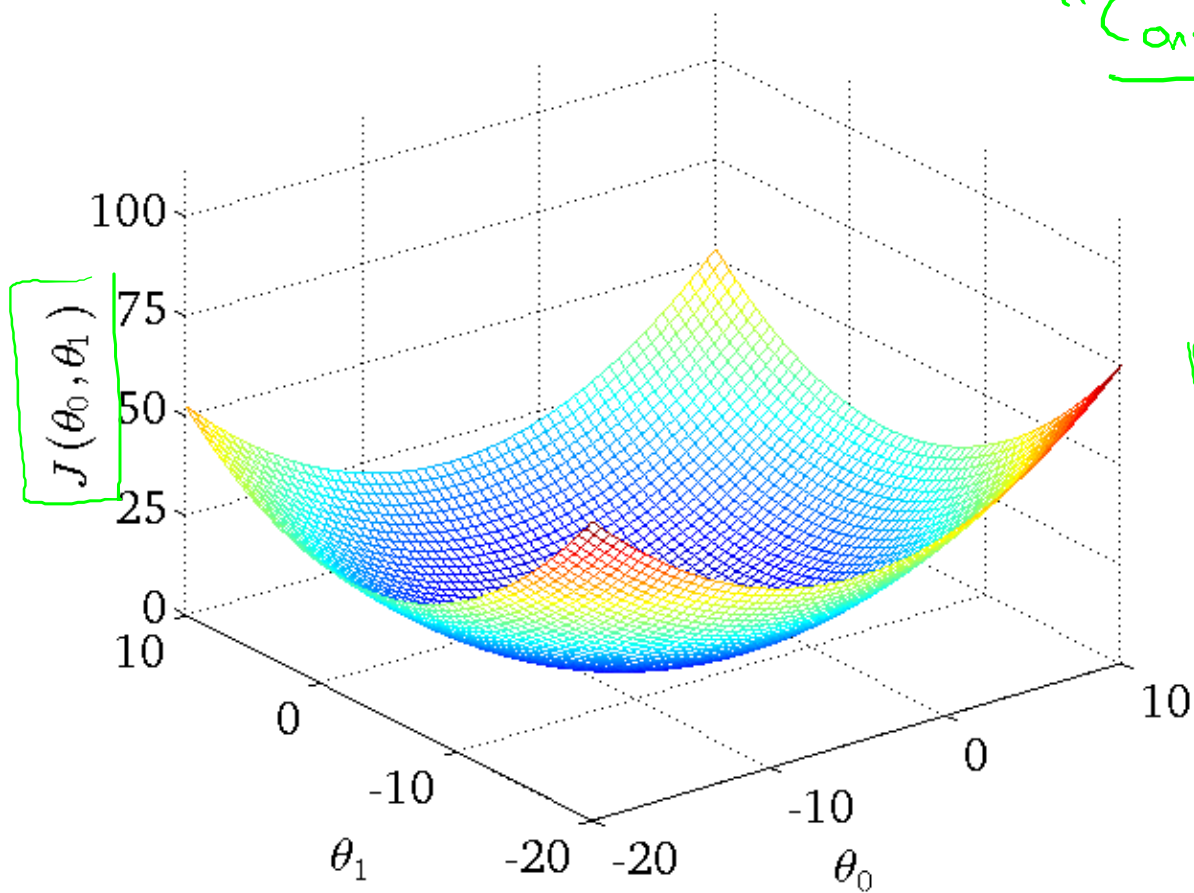
}

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

update
 θ_0 and θ_1
simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$



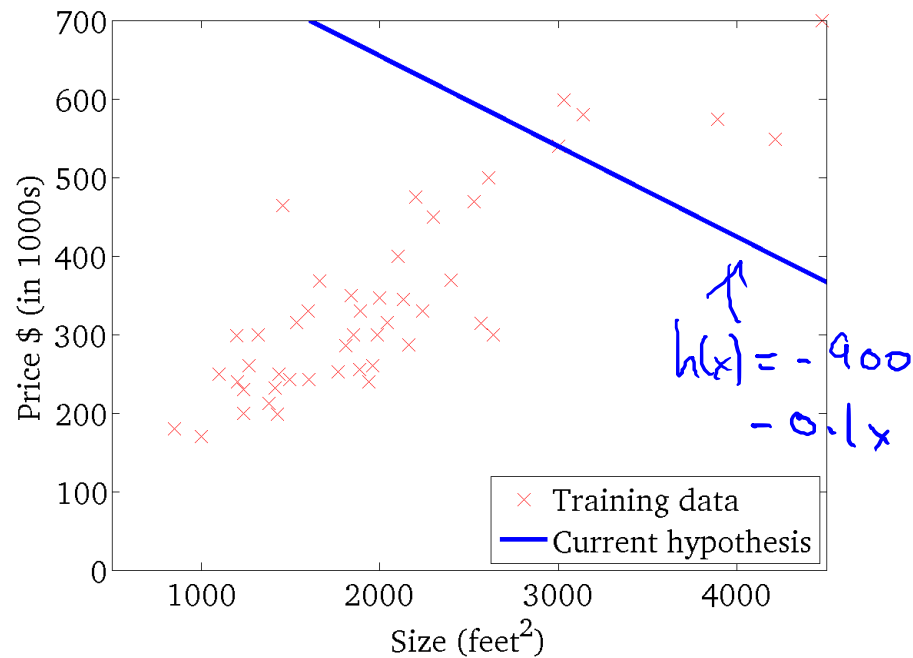


"Convex function"

Bowl-shaped

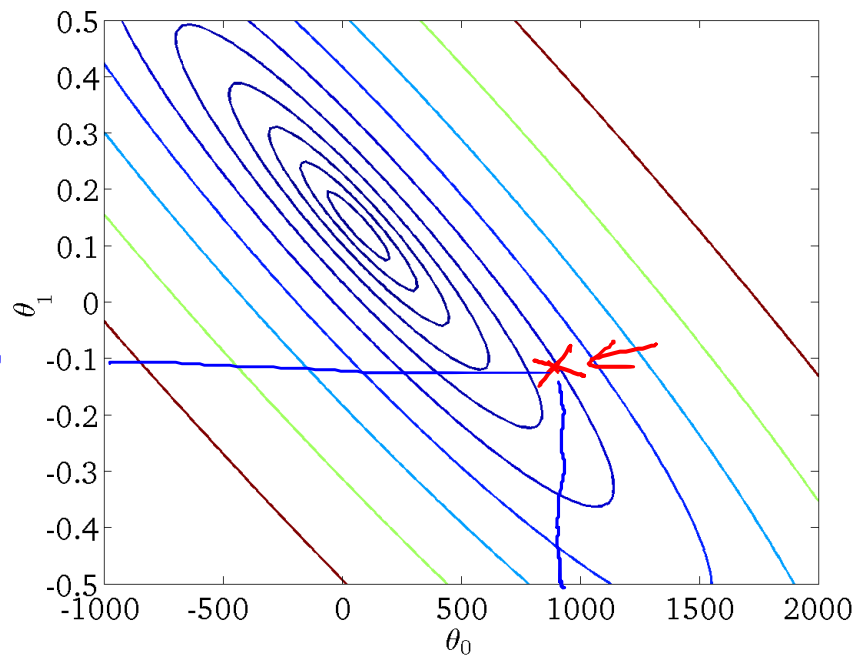
$$\underline{h_{\theta}(x)}$$

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



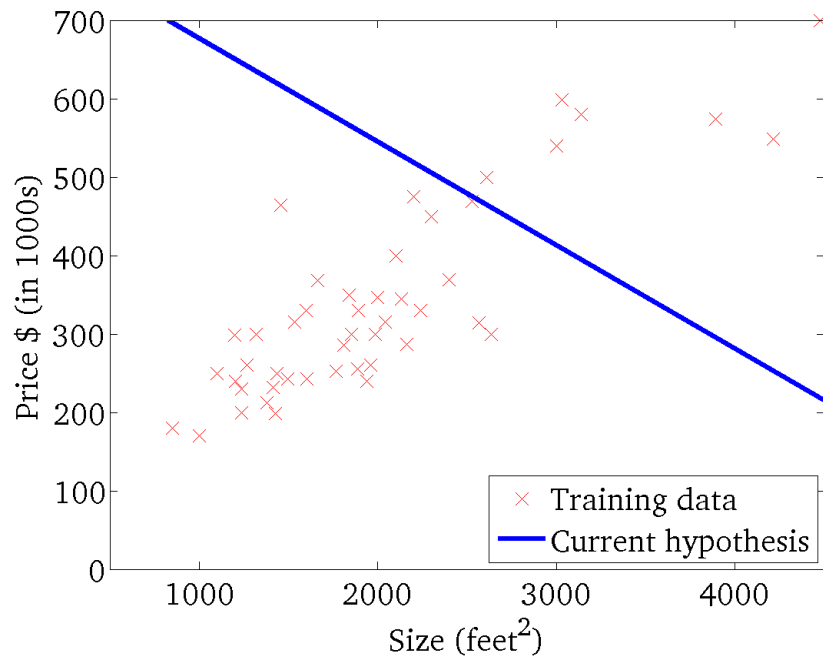
$$\underline{J(\theta_0, \theta_1)}$$

(function of the parameters θ_0, θ_1)



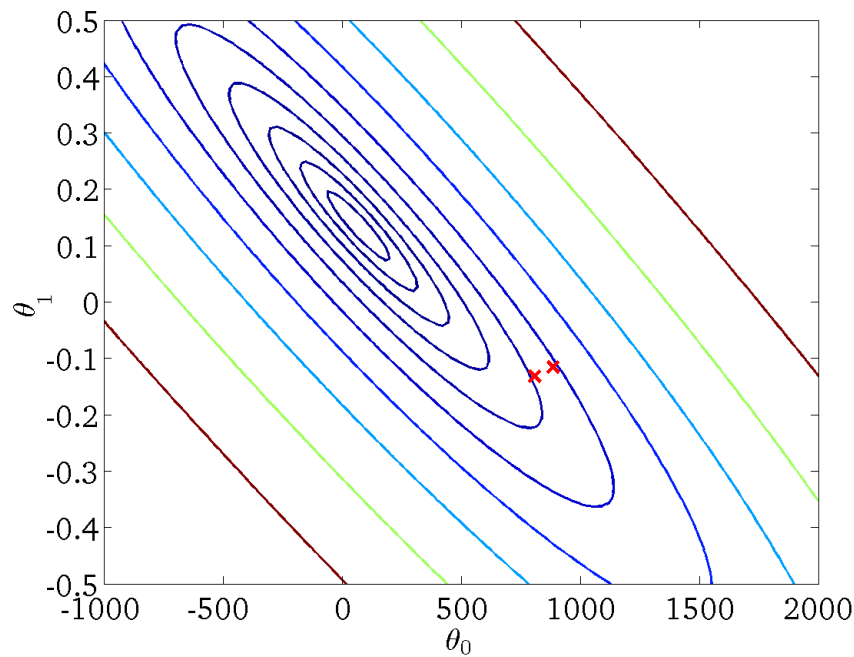
$$h_{\theta}(x)$$

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



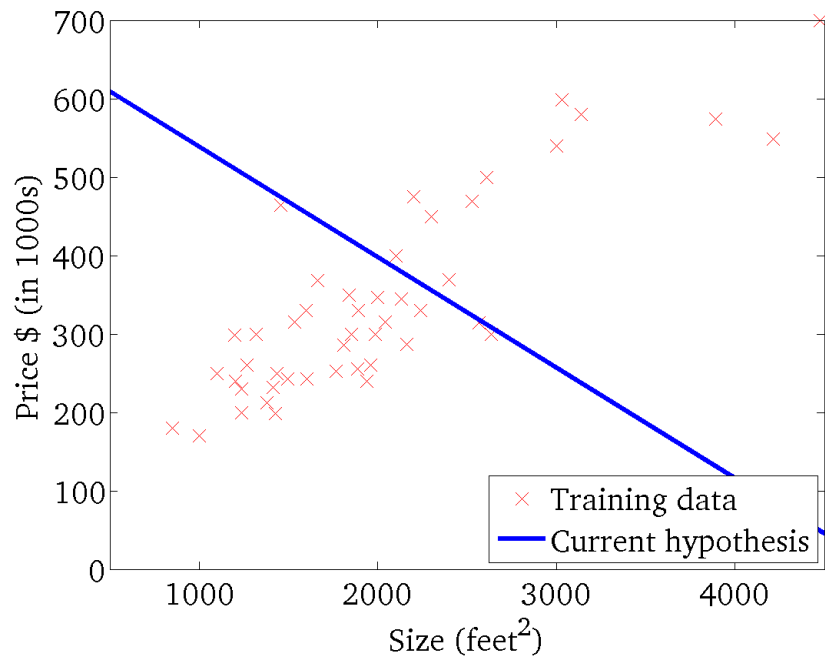
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



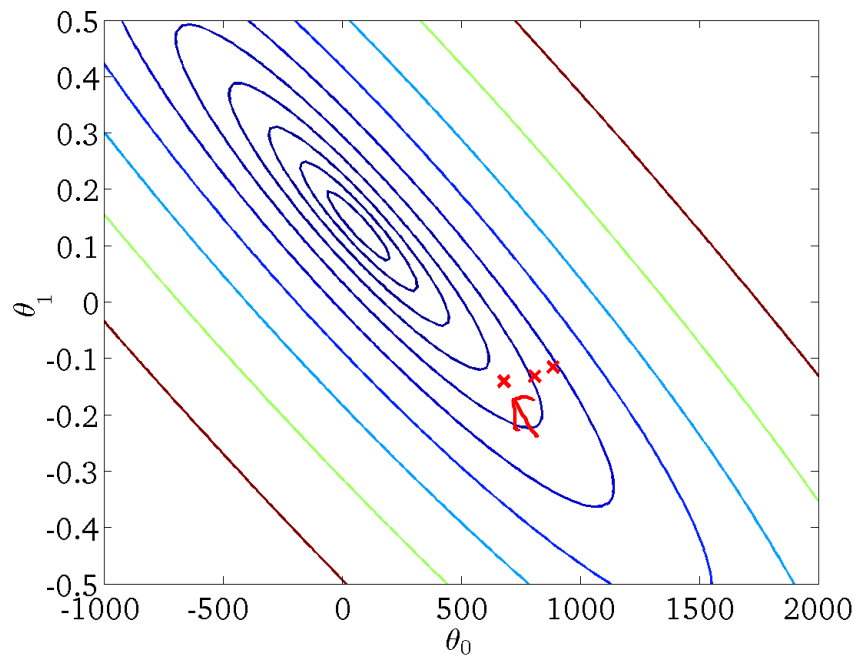
$$h_{\theta}(x)$$

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



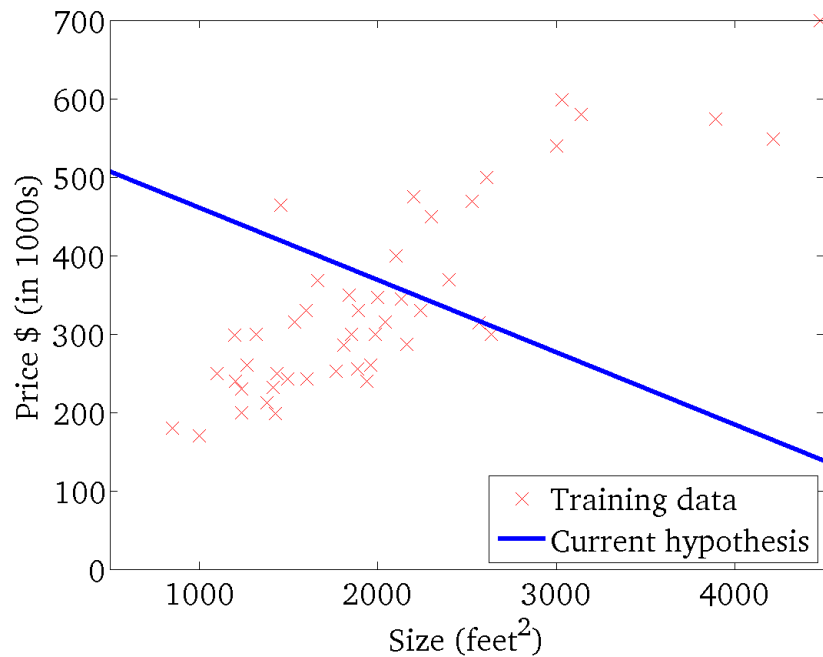
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



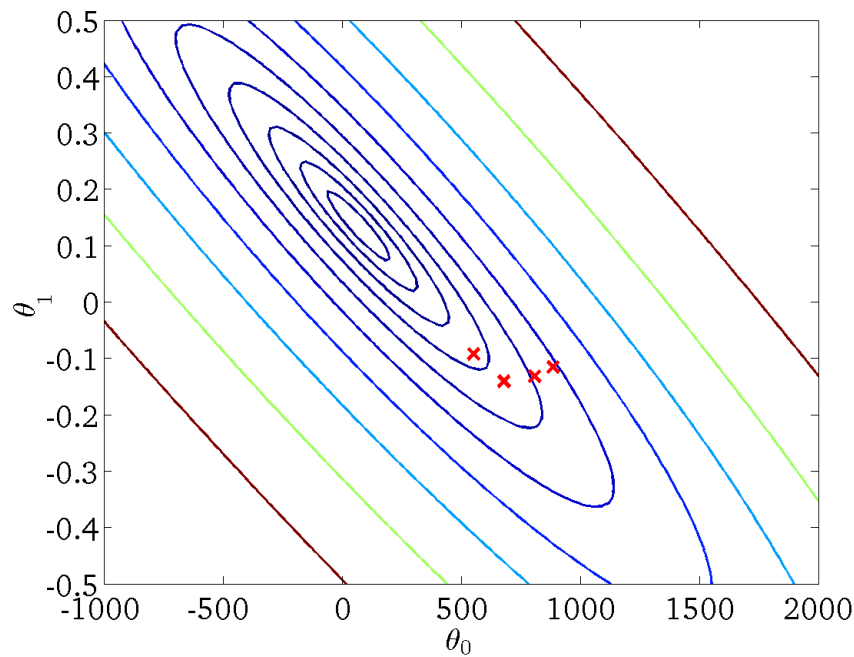
$$h_{\theta}(x)$$

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



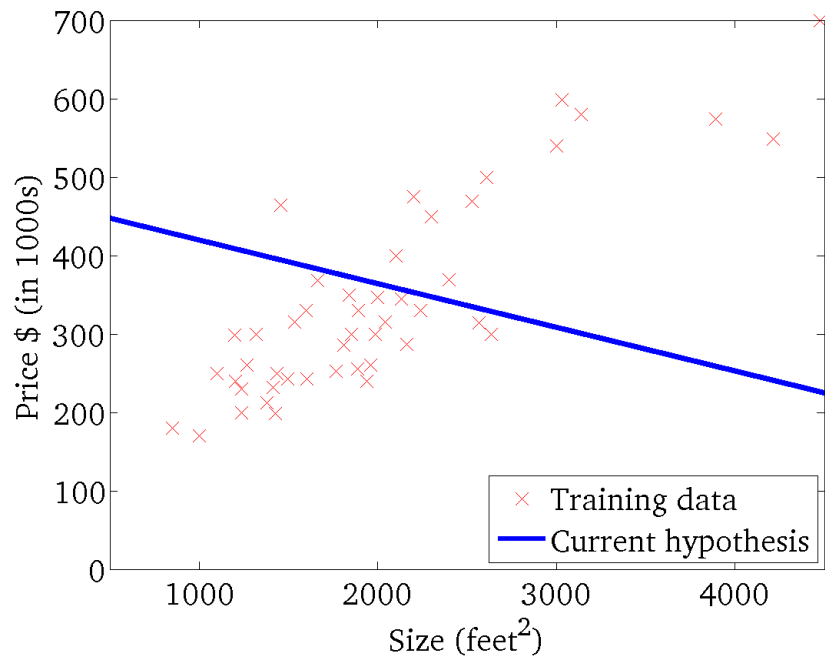
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



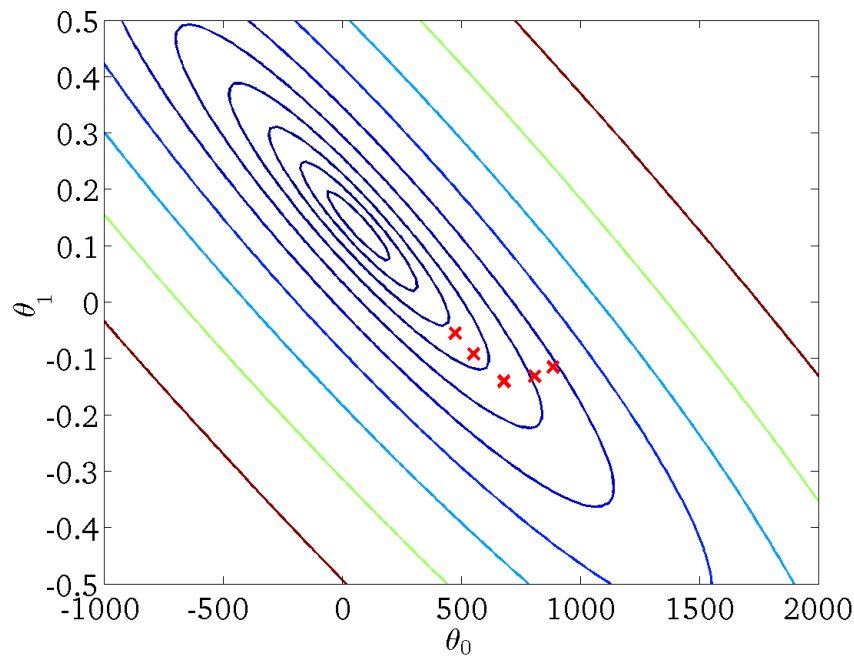
$$h_{\theta}(x)$$

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



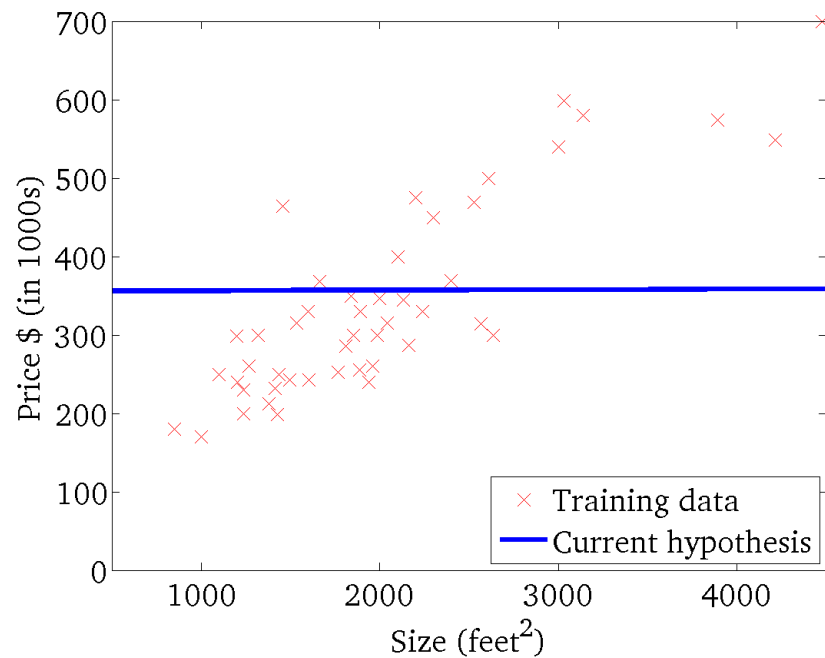
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



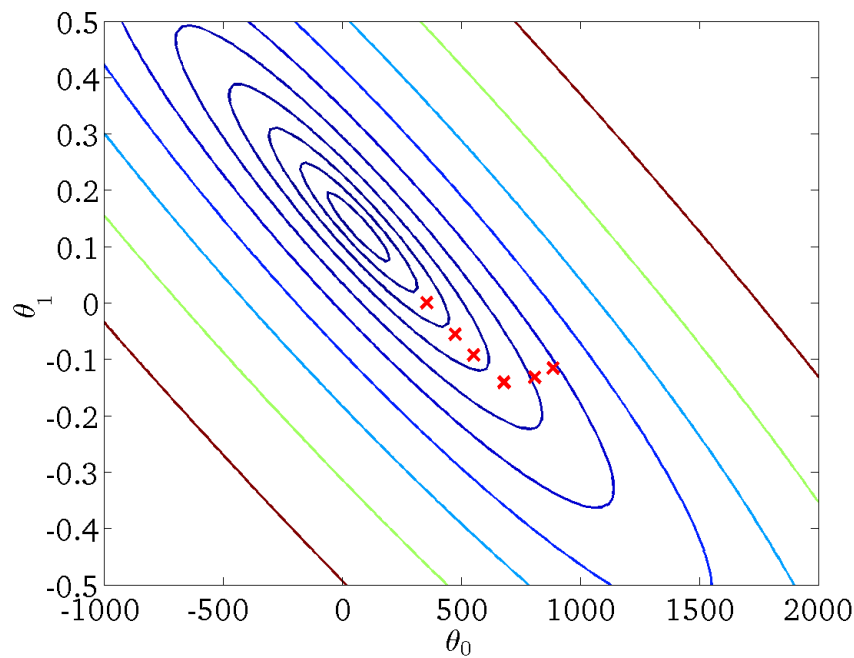
$$h_{\theta}(x)$$

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



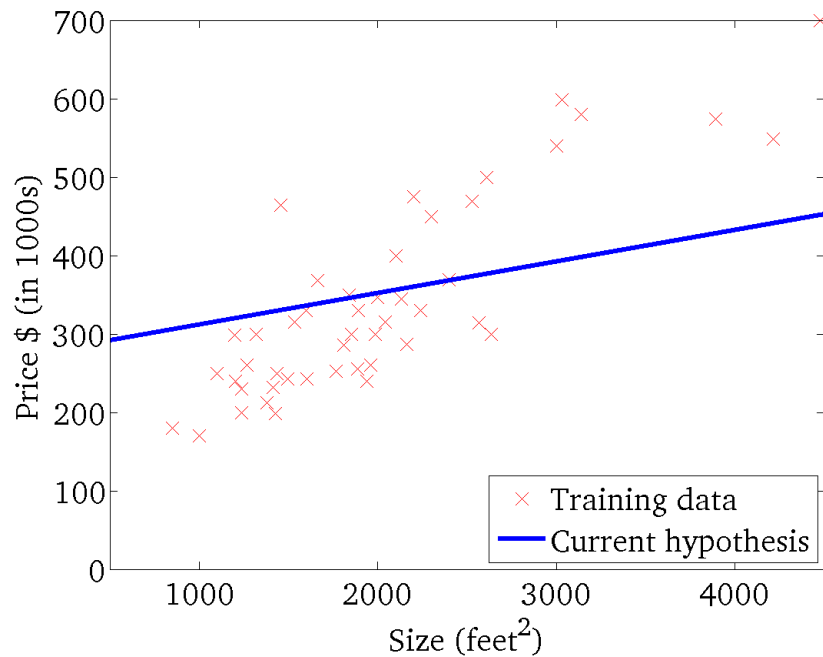
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



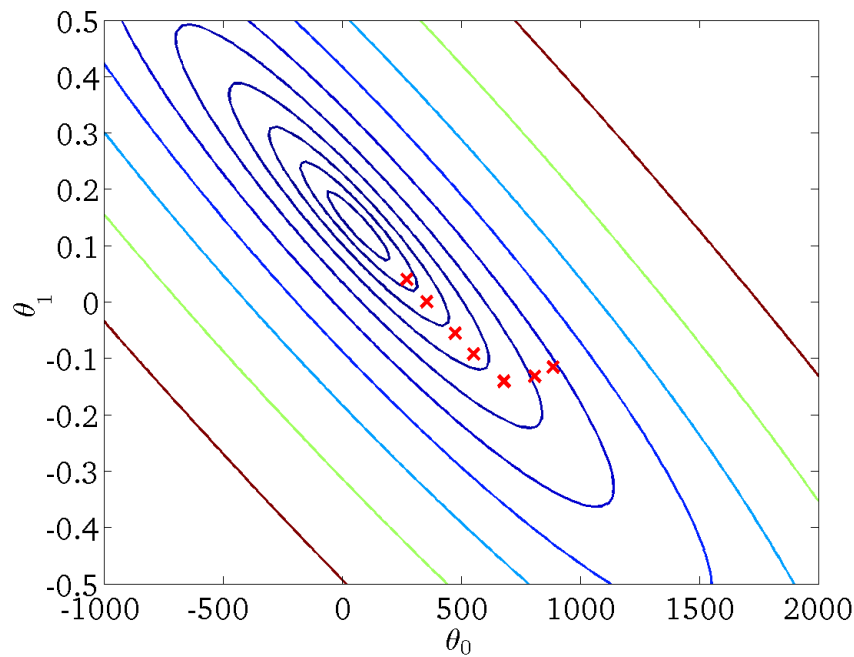
$$h_{\theta}(x)$$

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



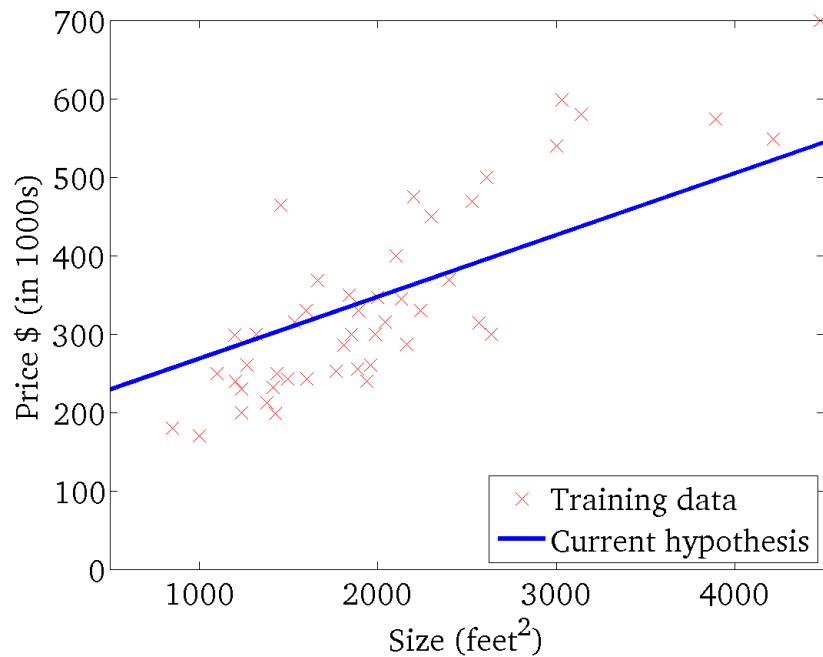
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



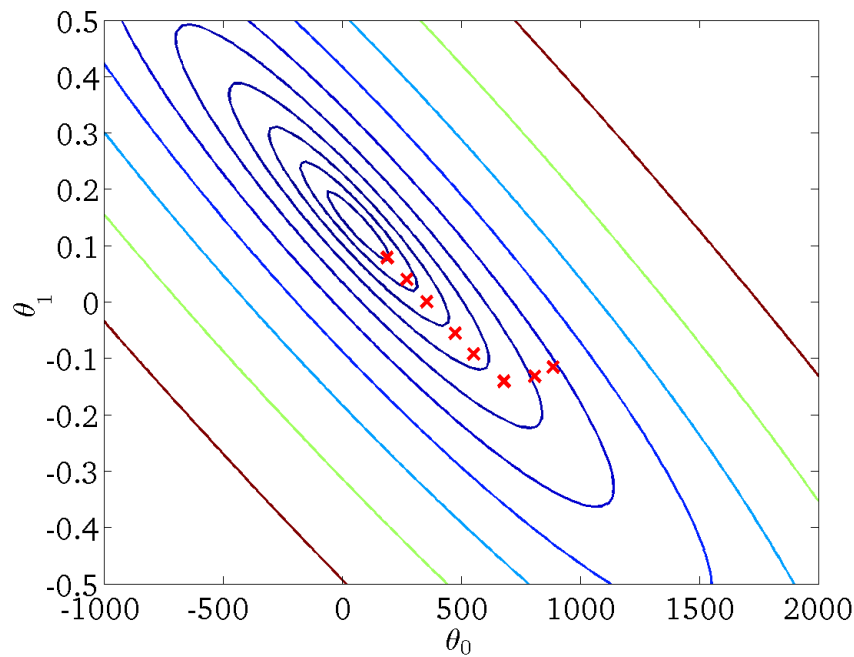
$$h_{\theta}(x)$$

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



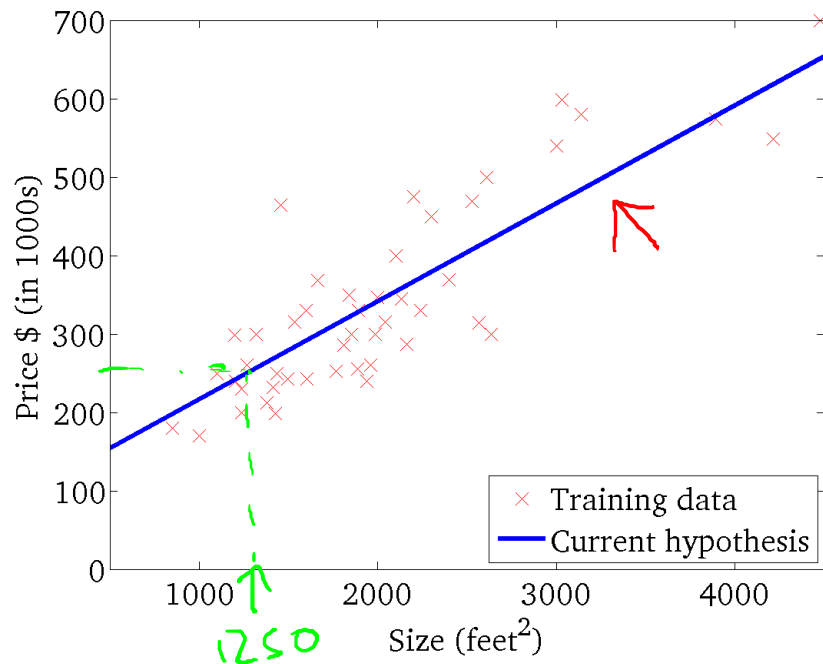
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



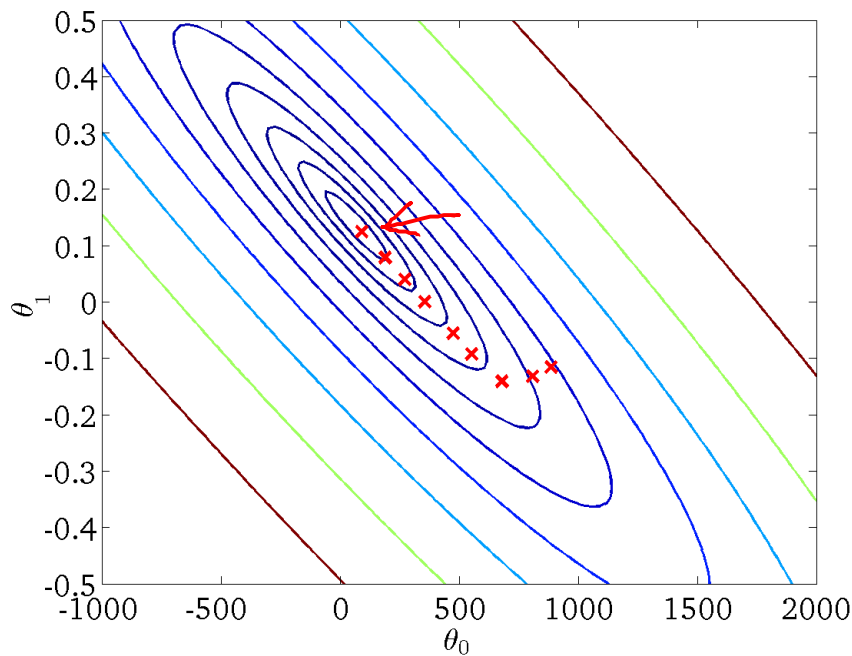
$$h_{\theta}(x)$$

(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)



“Batch” Gradient Descent

“Batch”: Each step of gradient descent uses all the training examples.

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

Two extensions:

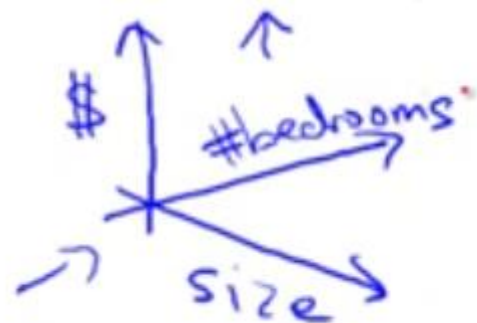
1. In $\min J(\theta_0, \theta_1)$, solve for θ_0, θ_1 exactly, without needing iterative algorithm (gradient descent).
2. Learn with larger number of features.

Size (feet ²)	Price (\$1000)
2104	460
1416	232
1534	315
852	178

Two extensions:

1. In $\min J(\theta_0, \theta_1)$, solve for θ_0, θ_1 exactly, without needing iterative algorithm (gradient descent).
2. Learn with larger number of features.


<u>Size (feet²)</u>	<u>Number of bedrooms</u>	<u>Number of floors</u>	<u>Age of home (years)</u>	<u>Price (\$1000)</u>
x_1	x_2	x_3	x_4	y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178




Linear Algebra

Notation and set of the things you can do with matrices and vectors.

Matrix:

$$X = \begin{bmatrix} 2104 & 5 & 1 & 45 \\ 1416 & 3 & 2 & 40 \\ 1534 & 3 & 2 & 30 \\ 852 & 2 & 1 & 36 \end{bmatrix}$$


Vector:

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 172 \end{bmatrix}$$


Topics:

- What are matrices and vectors
- Addition, subtraction, multiplication with matrices and vectors
- Matrix inverse, transpose