

Machine Learning and Data Analytics

ME 5013- Fall 2019

Lecture 05

- Cost Function
- Gradient Descent



The University of Texas at San Antonio™

Adel Alaeddini, PhD

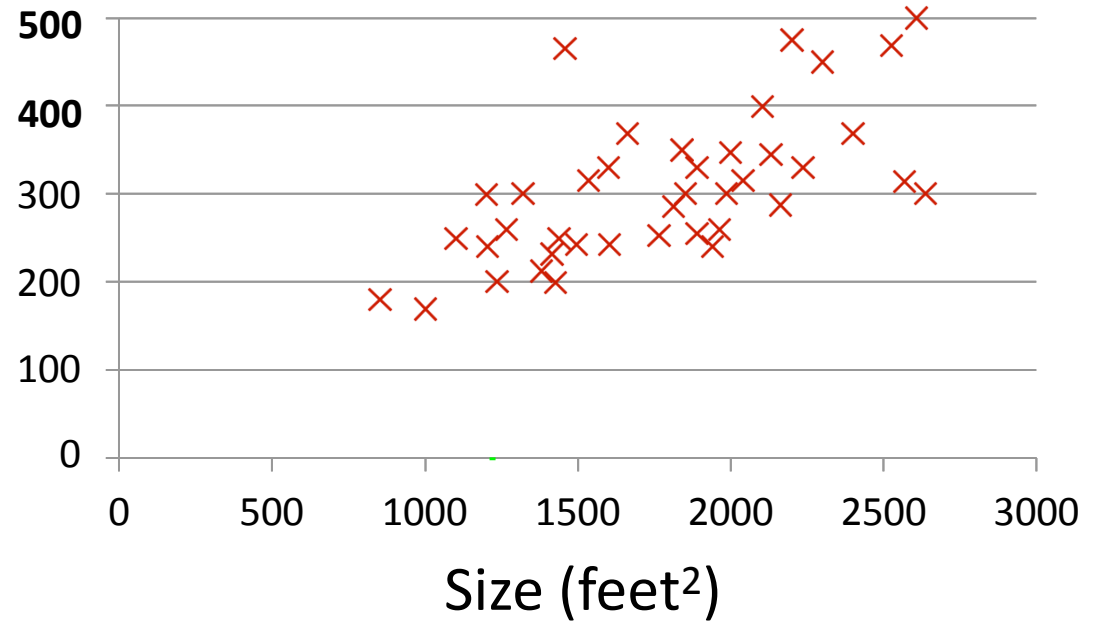
Associate Professor of Mechanical Engineering

Advanced Data Engineering Lab

adel.alaeddini@utsa.edu

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

Training Set

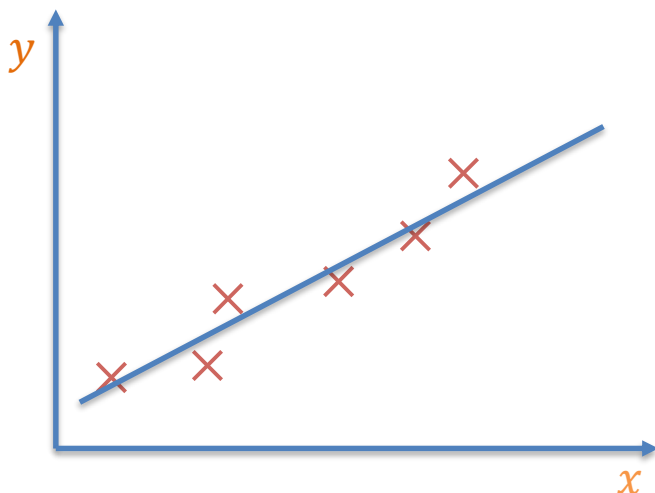
$$Y = \beta_0 + \beta_1 X + \epsilon$$

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

Hypothesis: $h_{\beta}(x) = \beta_0 + \beta_1 x$

β_i 's: Parameters

How to choose β_i 's ?



Idea: Choose β_0, β_1 so that $h_{\beta}(x)$ is close to y for our (x, y) training examples

Cost function I: $RSS = e_1^2 + e_2^2 + \dots + e_n^2$

$$RSS = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$$

$$RSS = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

Closed formed formula:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Cost function II:

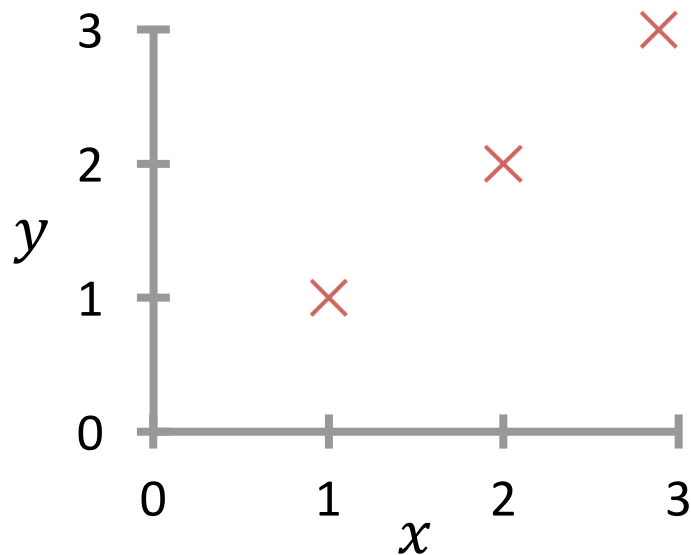
Hypothesis: $h_{\beta}(x) = \beta_0 + \beta_1 x$

$$\min_{\beta_0, \beta_1} J(\beta_0, \beta_1) = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

- Let $\beta_0 \equiv 0$

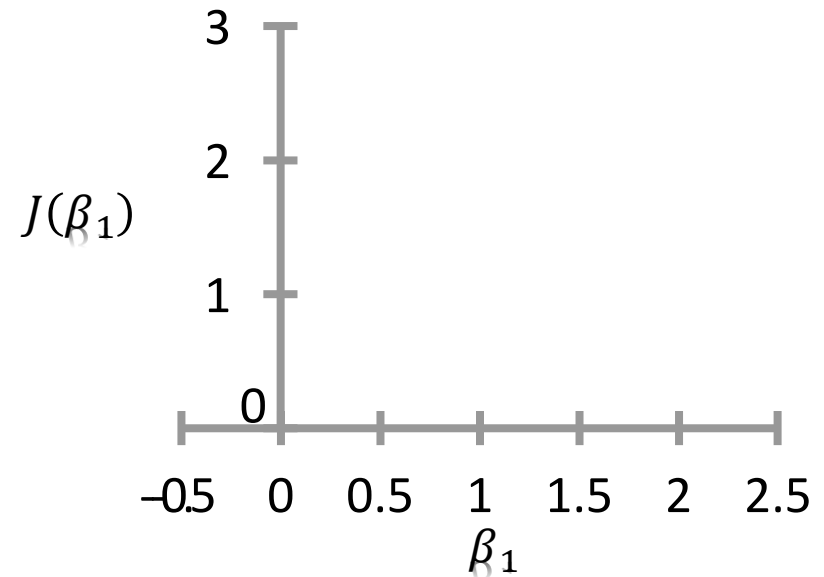
$$h_{\beta_1}(x)$$

For fixed β_1 this is a function of x



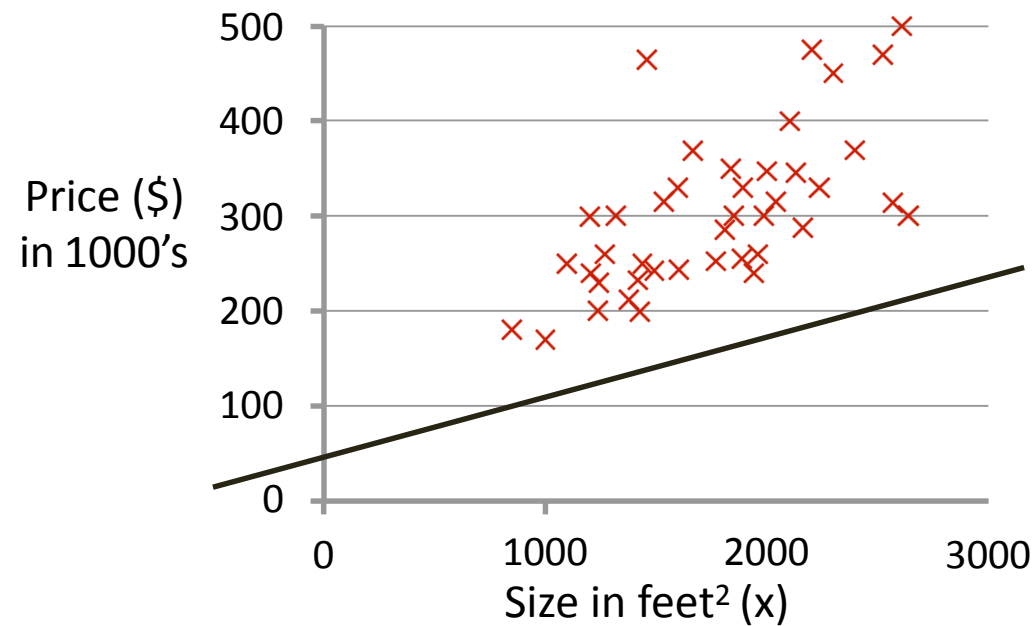
$$J(\beta_1)$$

(function of the parameter β_1)



$$h_{\beta_0, \beta_1}(x)$$

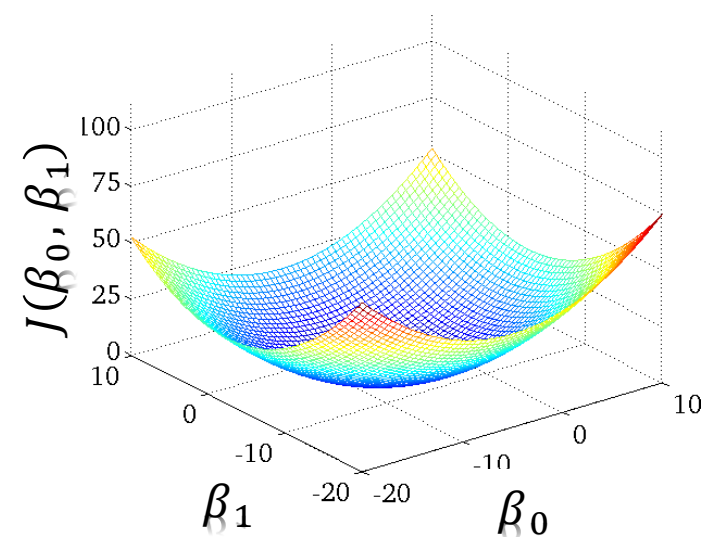
For fixed β_0, β_1 this is a function of x



$$h_{\beta_0, \beta_1}(x) = 50 + 0.06x$$

$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



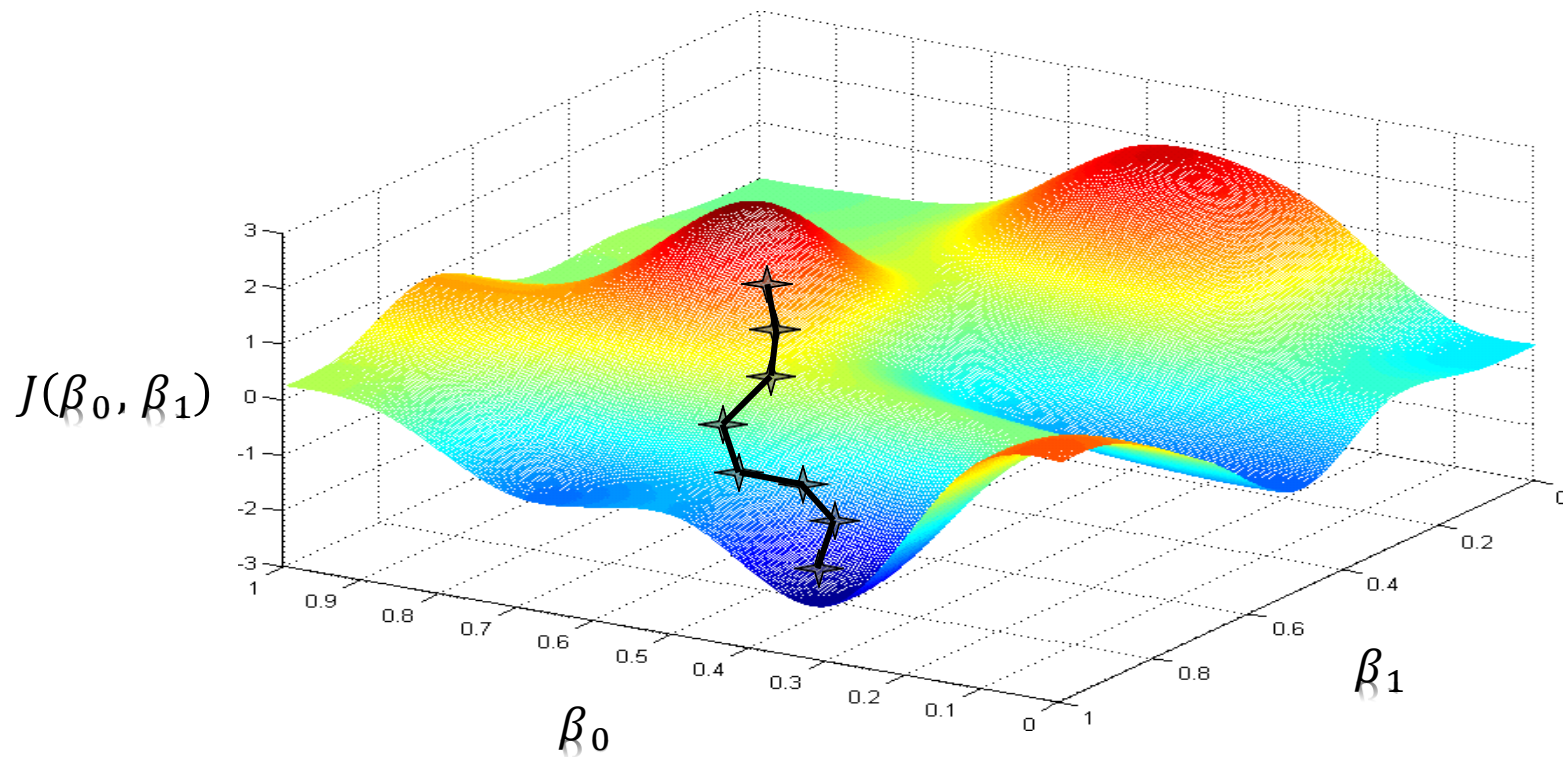
Have some function $J(\beta_0, \beta_1)$

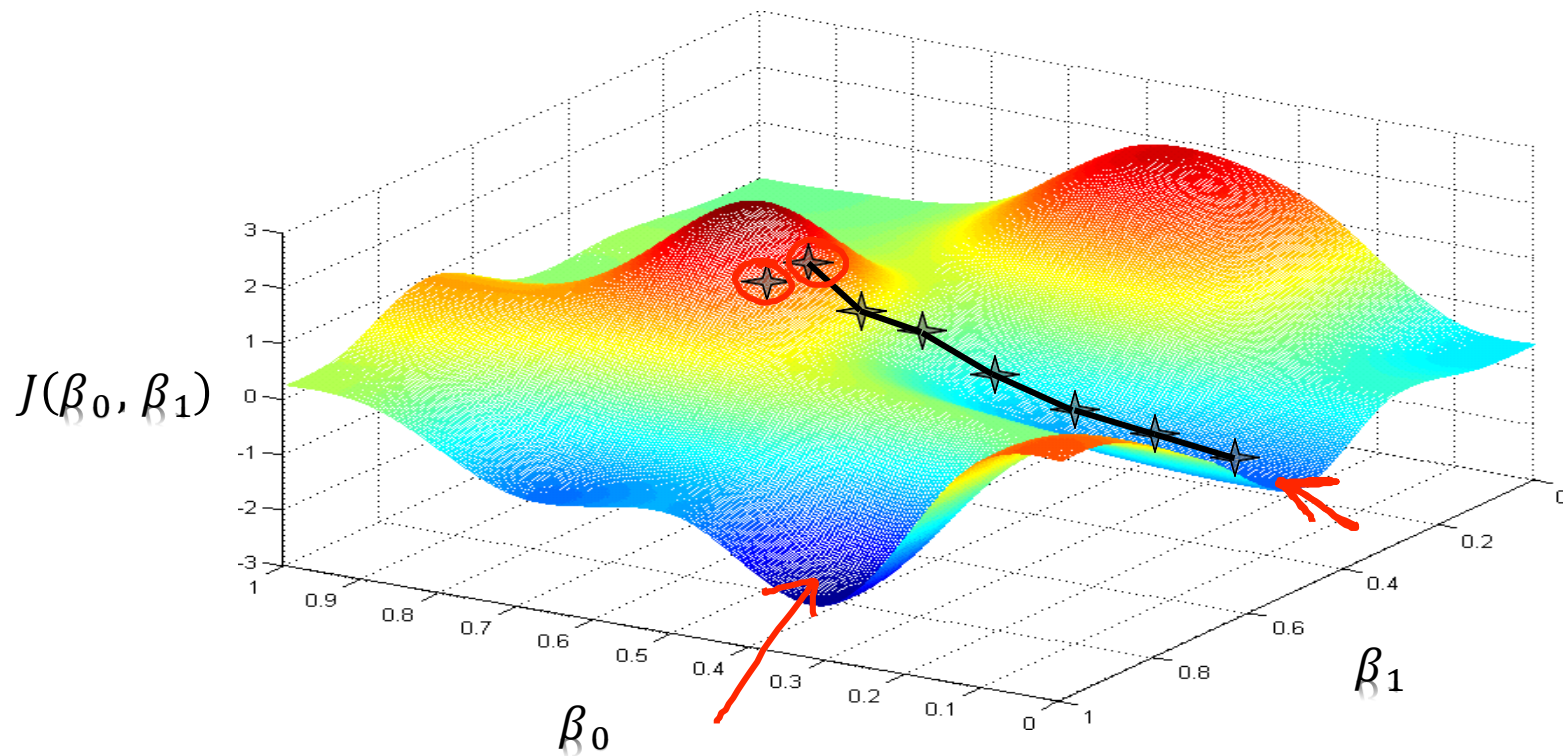
Want $\min_{\beta_0, \beta_1} J(\beta_0, \beta_1)$

Outline:

- Start with some β_0, β_1
- Keep changing β_0, β_1 to reduce $J(\beta_0, \beta_1)$

Until we hopefully end up at a minimum





Gradient descent algorithm

repeat until convergence {

$$\beta_j := \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\beta_0, \beta_1) \text{ (for } j = 0 \text{ and } j = 1)$$

}

Correct: Simultaneous update

$$\text{temp0} := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1)$$

$$\text{temp1} := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

$$\beta_0 := \text{temp0}$$

$$\beta_1 := \text{temp1}$$

Incorrect:

$$\text{temp0} := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1)$$

$$\beta_0 := \text{temp0}$$

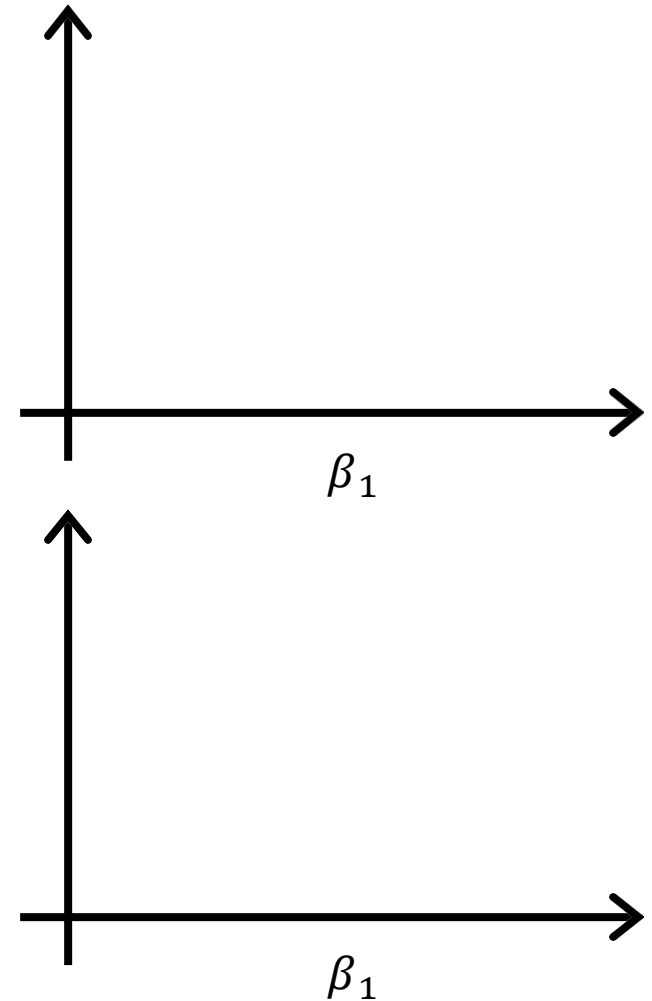
$$\text{temp1} := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

$$\beta_1 := \text{temp1}$$

$$\beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

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

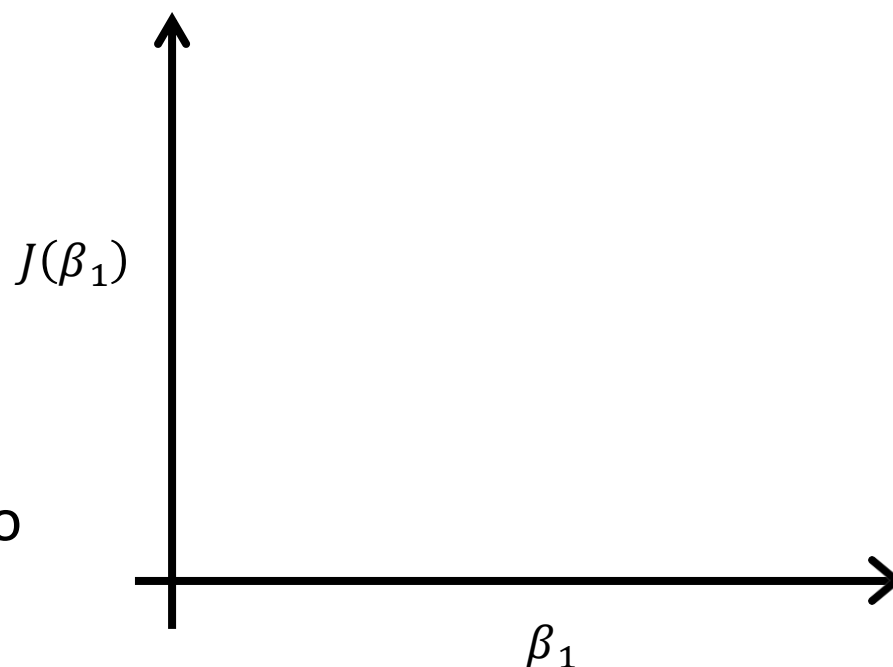
If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



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

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

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



Gradient descent algorithm

repeat until convergence {
 $\beta_j := \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\beta_0, \beta_1)$
 (for $j = 0$ and $j = 1$)
}

Hypothesis: $h_{\beta}(x) = \beta_0 + \beta_1 x$

$$\begin{aligned} & \min_{\beta_0, \beta_1} J(\beta_0, \beta_1) \\ &= \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \end{aligned}$$

repeat until convergence {

$$\beta_0 := \beta_0 + \alpha \frac{1}{n} \sum_{i=1}^n (y_n - \beta_0 - \beta_1 x_i)$$

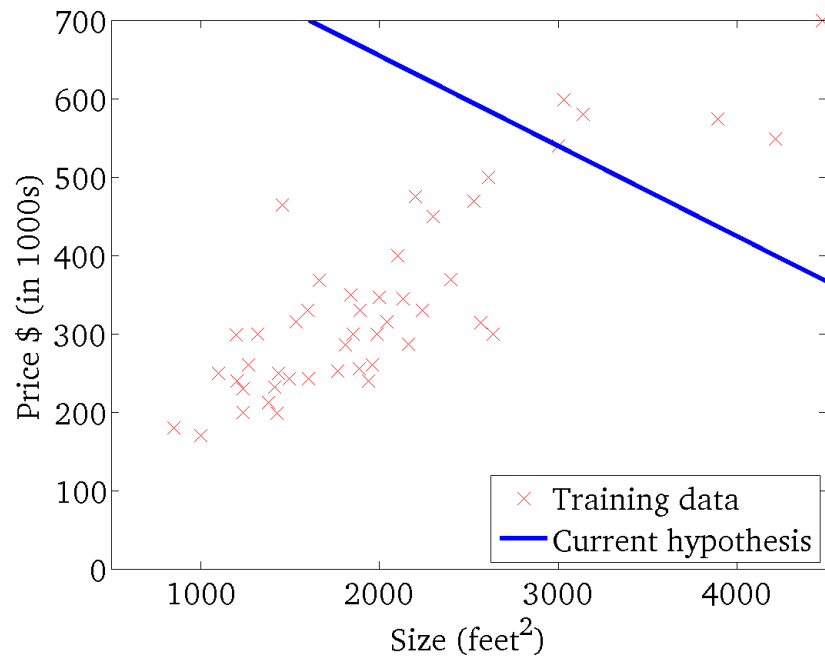
$$\beta_1 := \beta_1 + \alpha \frac{1}{n} \sum_{i=1}^n (y_n - \beta_0 - \beta_1 x_i) x_i$$

} update
 β_0 and β_1
 simultaneously

}

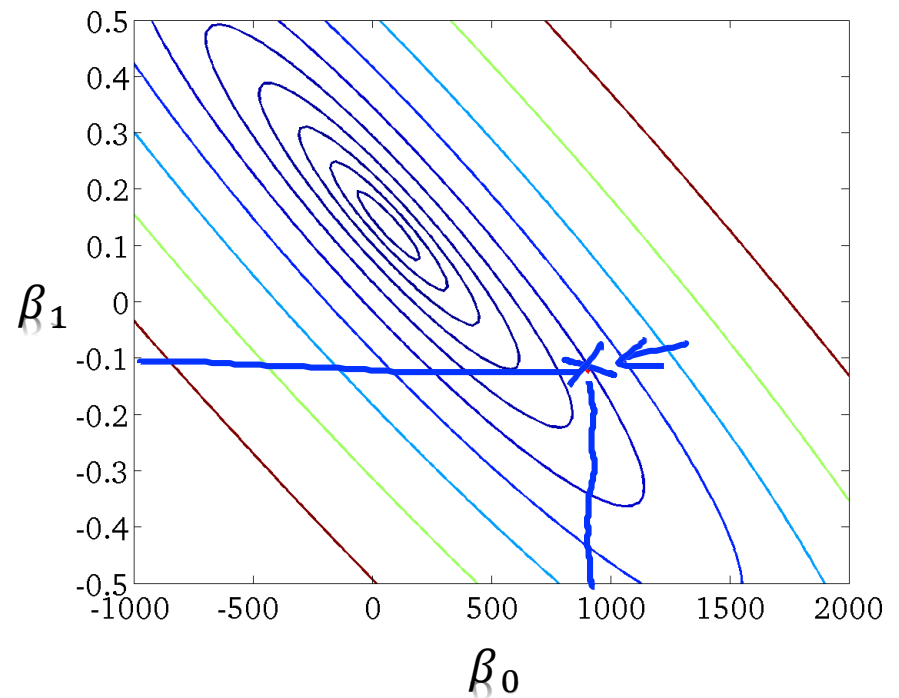
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



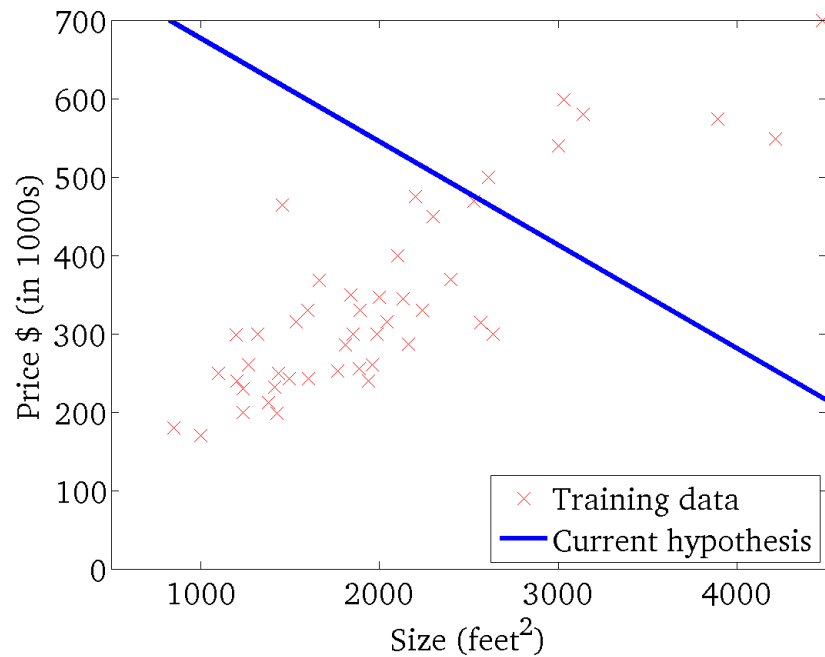
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



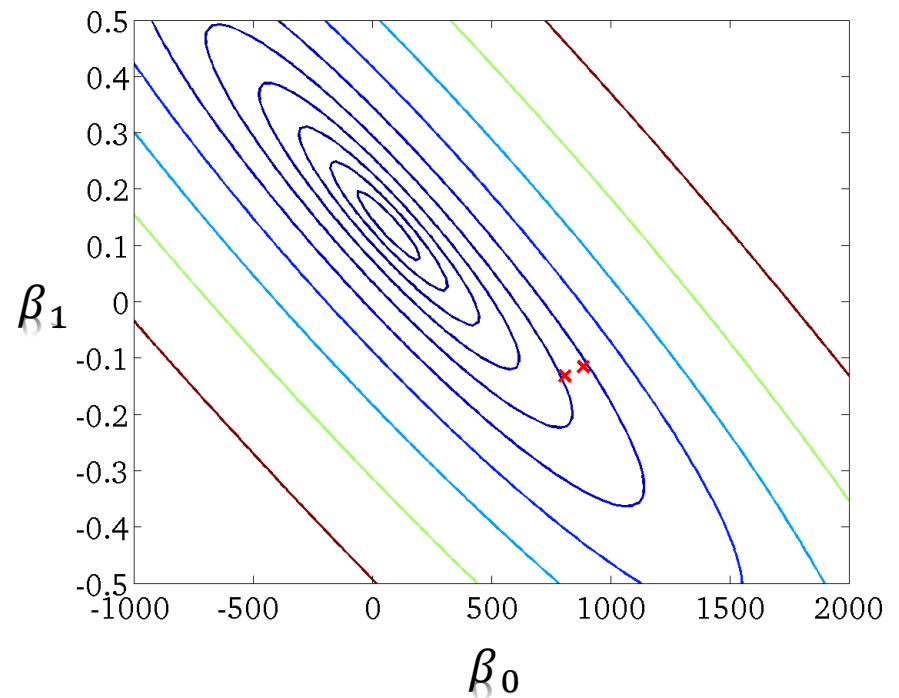
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



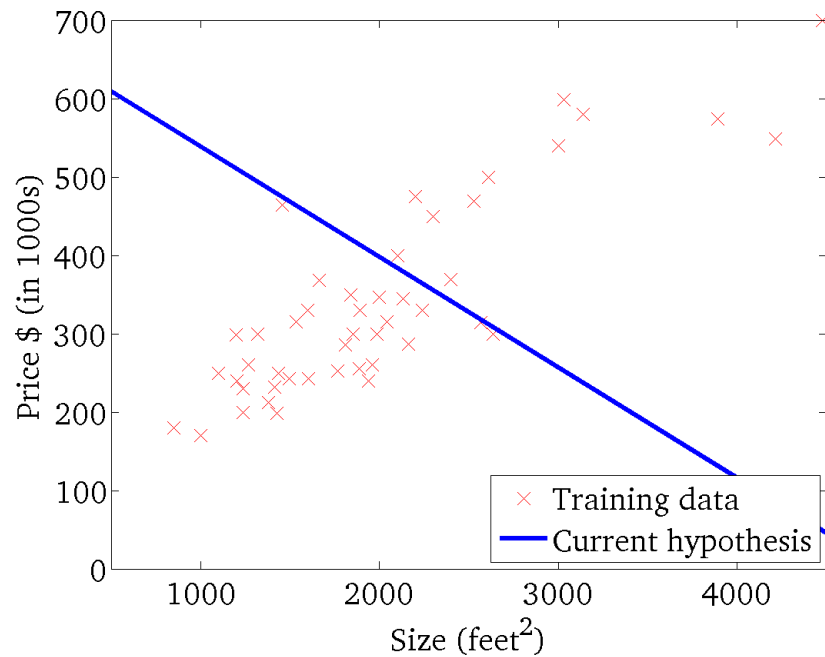
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



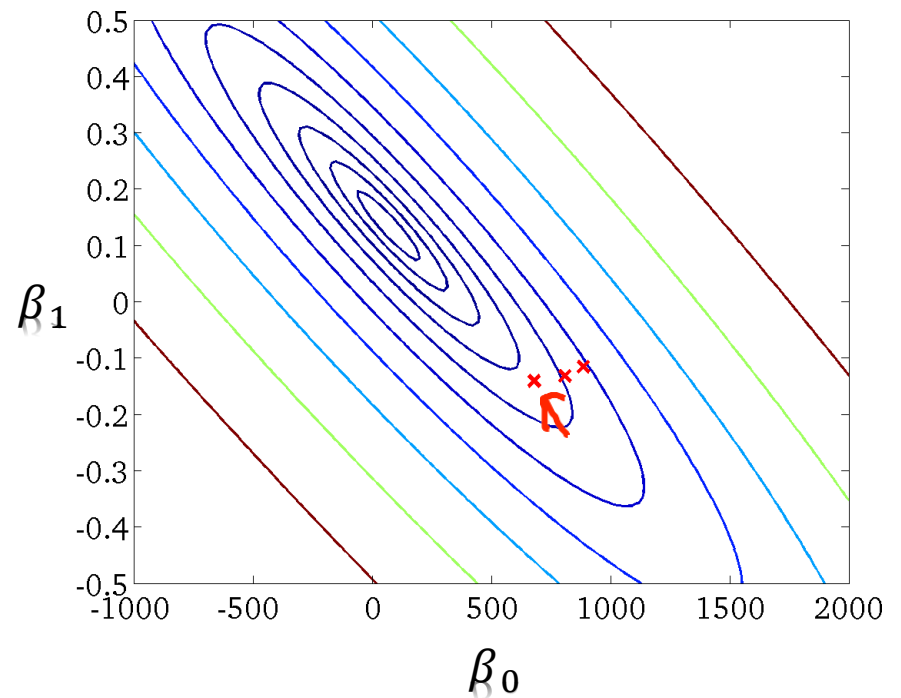
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



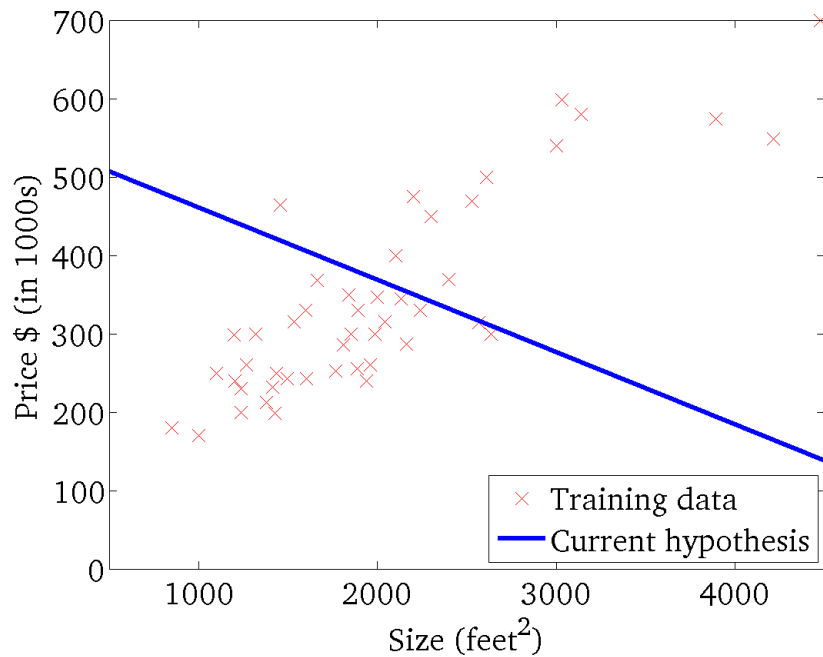
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



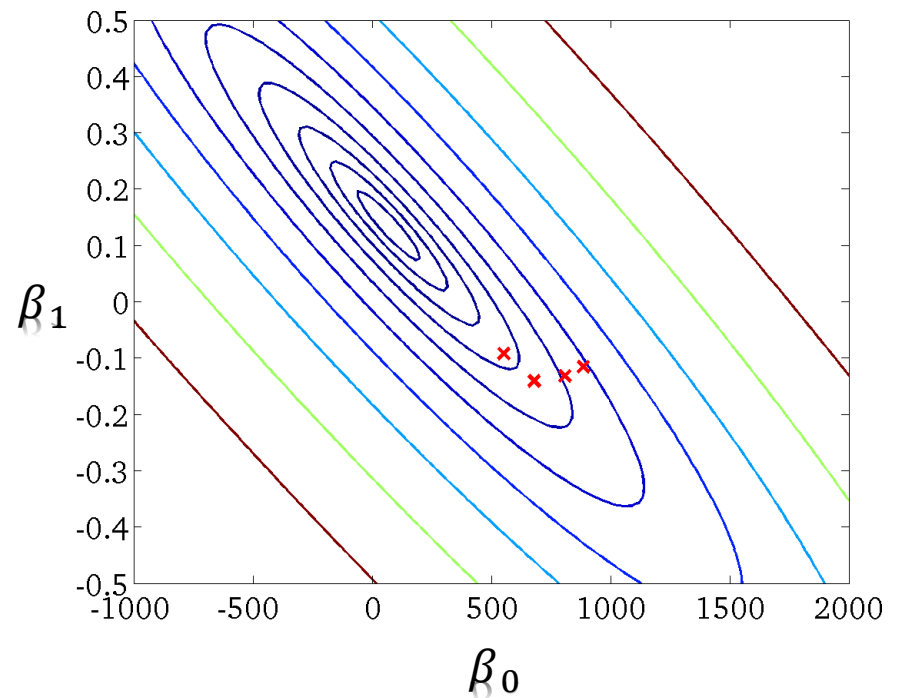
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



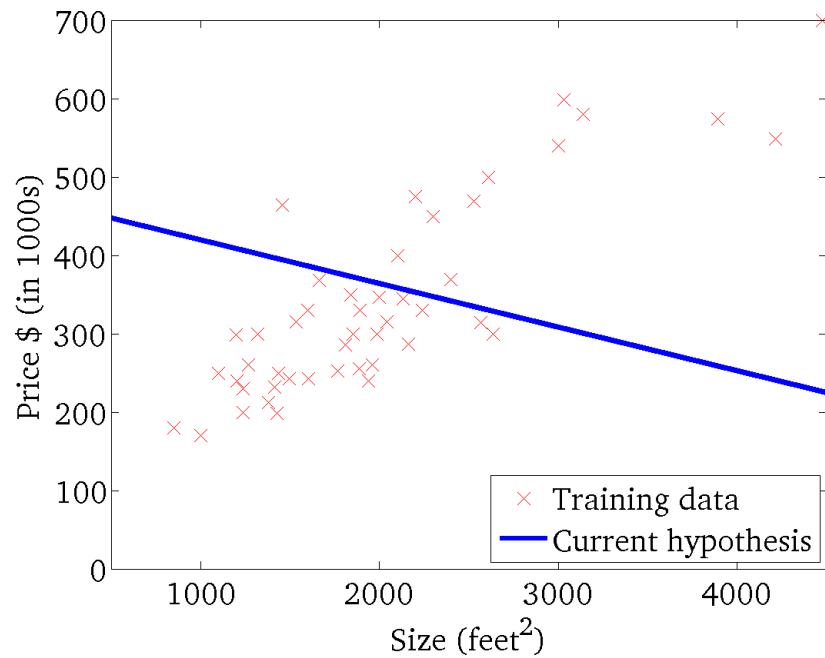
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



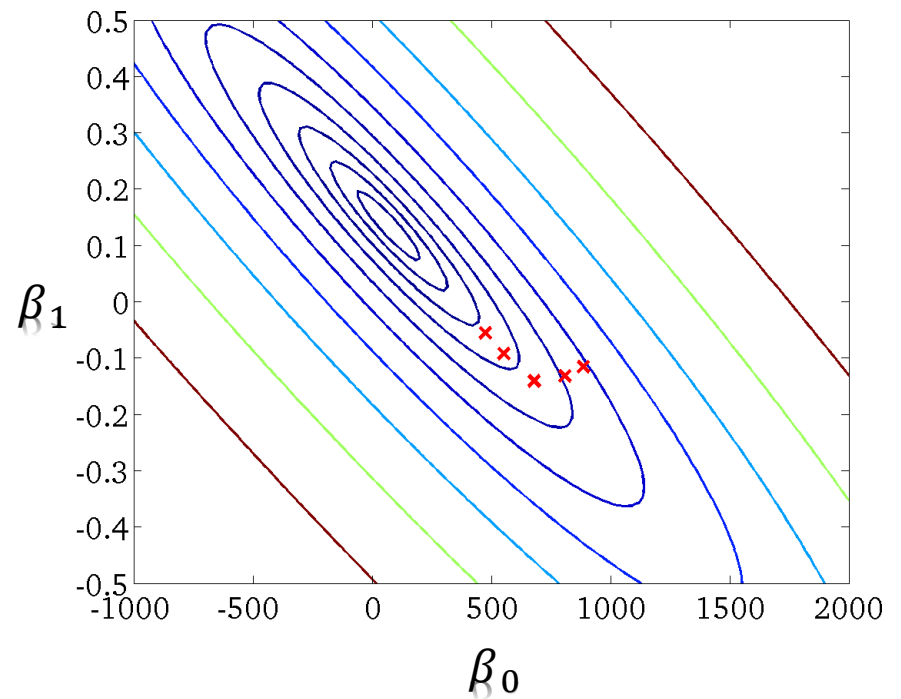
$$h_{\beta}(x)$$

For fixed β_0, β_1 this is a function of x



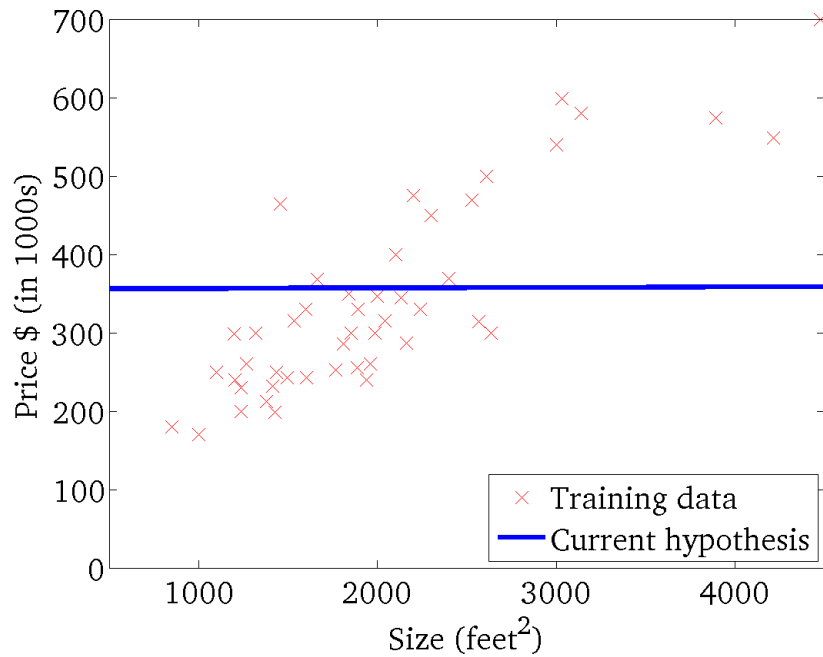
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



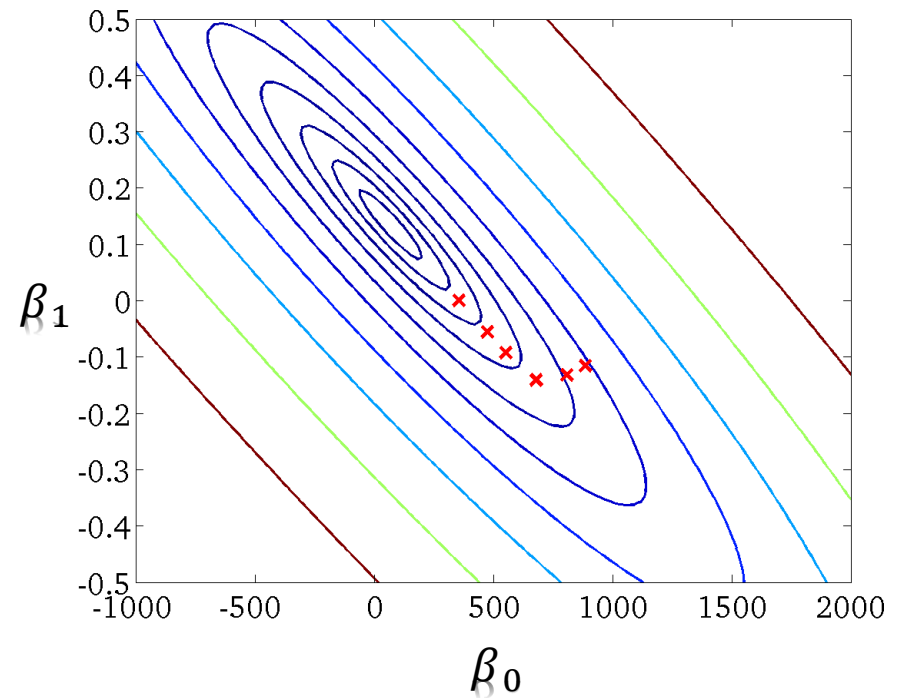
$$h_{\beta}(x)$$

For fixed β_0, β_1 this is a function of x



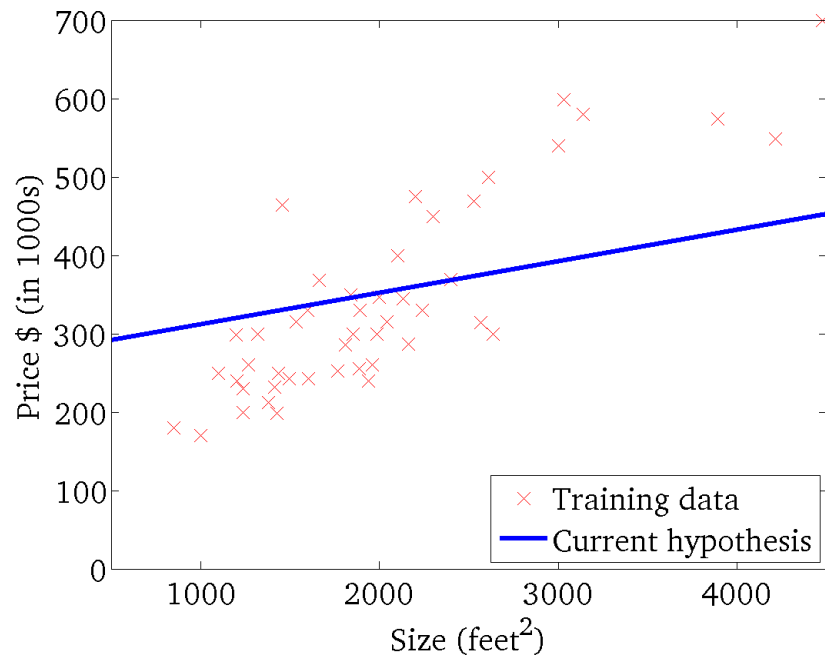
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



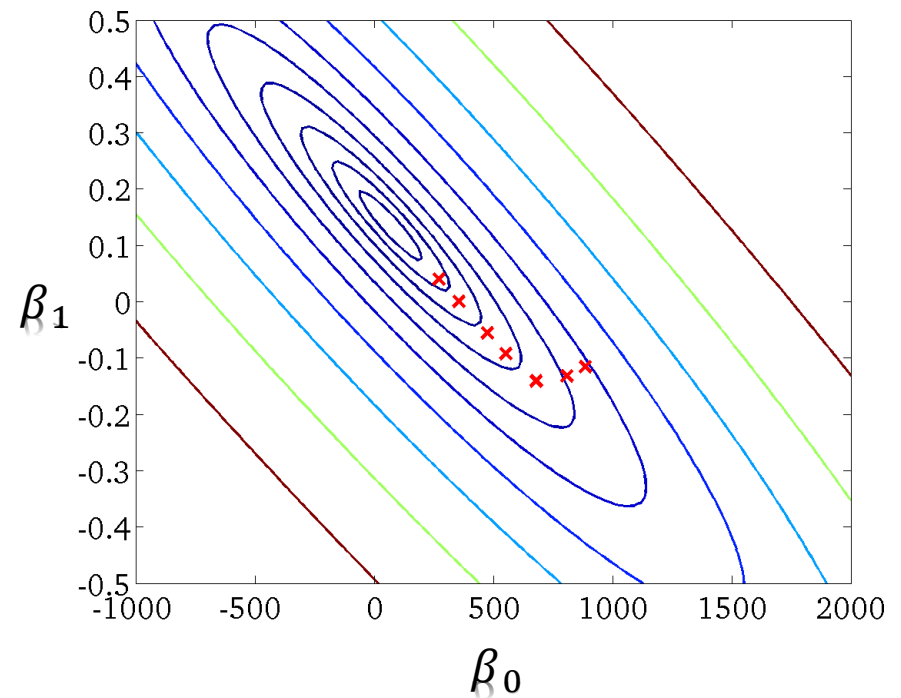
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



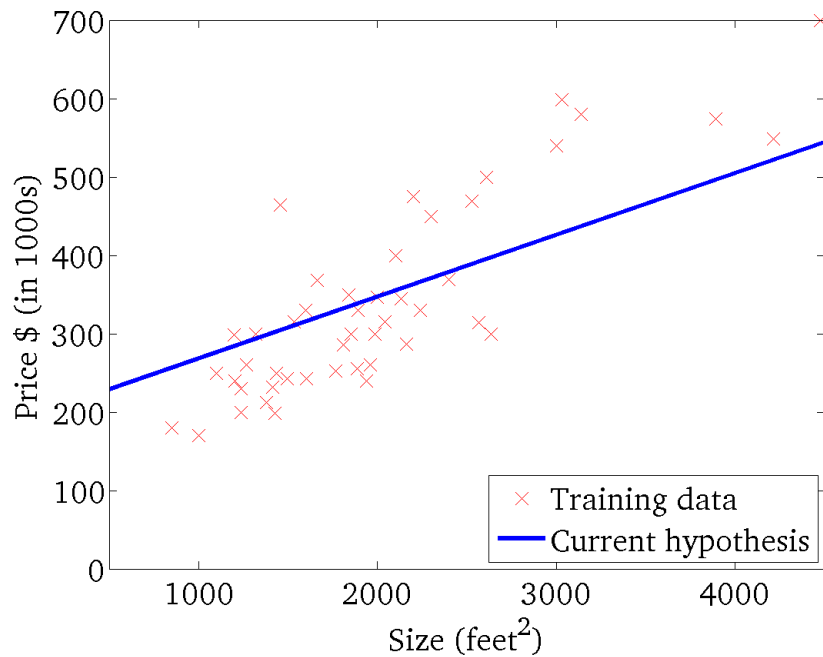
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



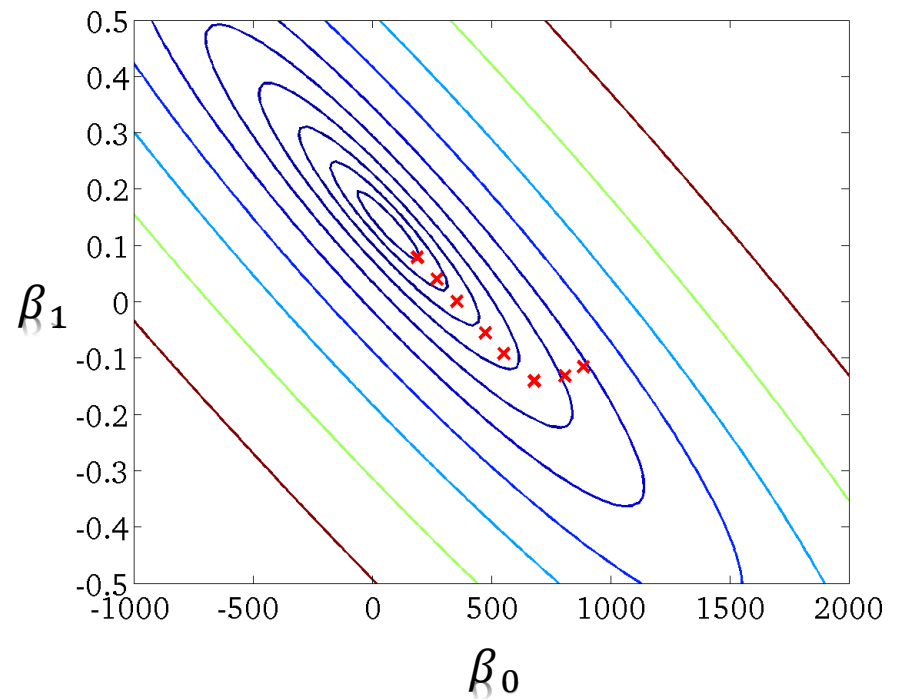
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



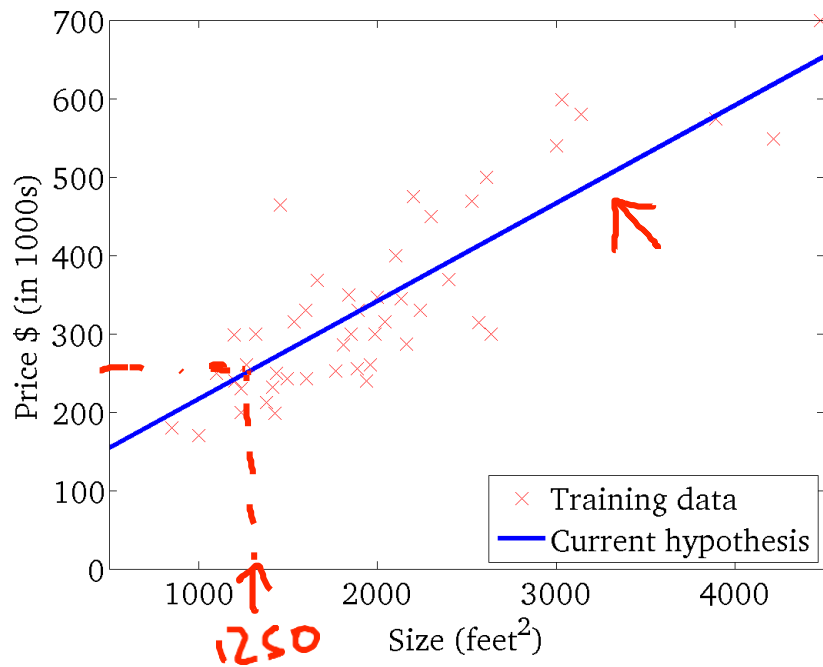
$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



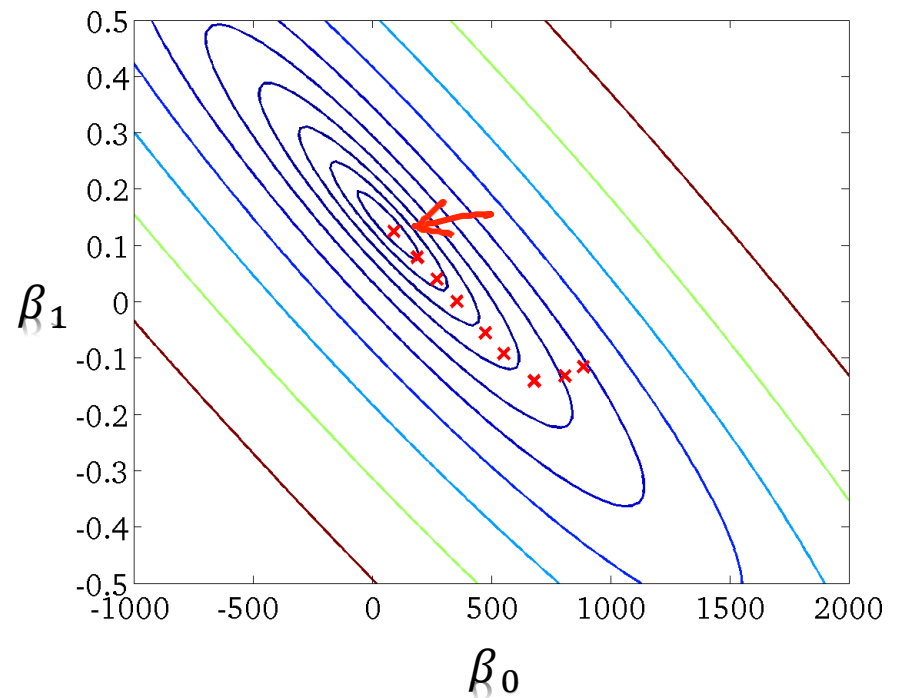
$$h_{\beta_0, \beta_1}(x)$$

For fixed β_0, β_1 this is a function of x



$$J(\beta_0, \beta_1)$$

(function of the parameters β_0, β_1)



“Batch” Gradient Descent

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