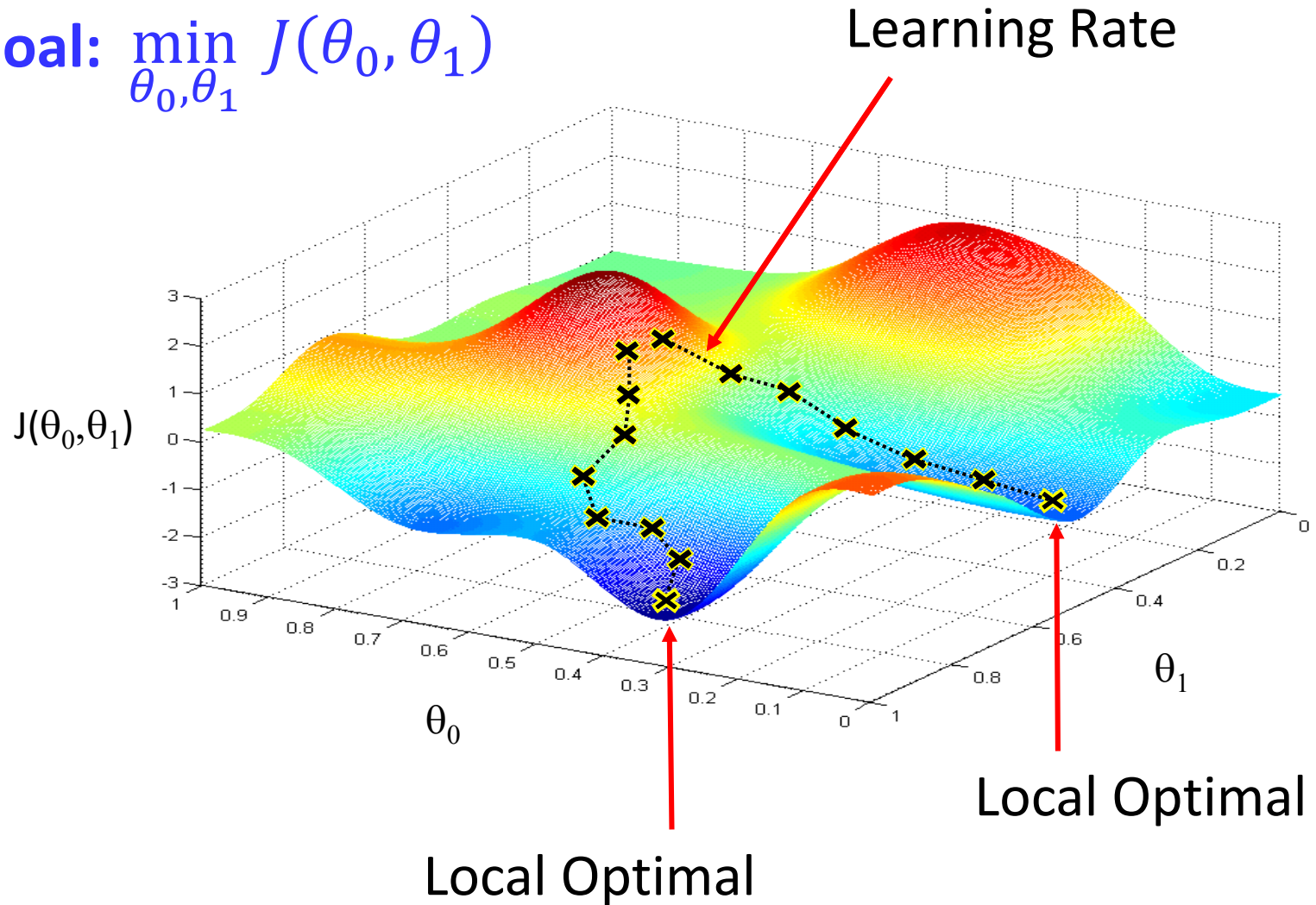


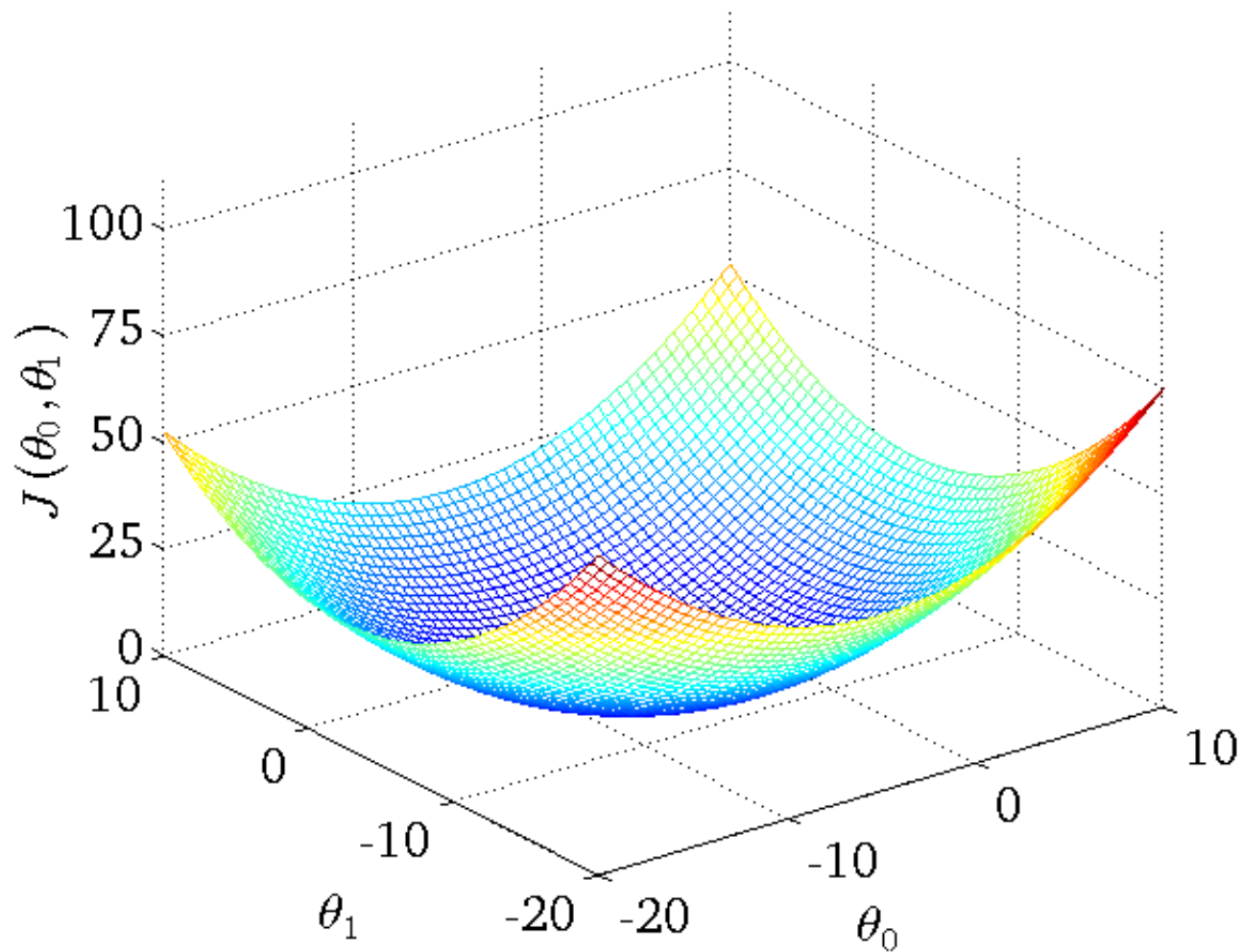
Gradient Descent

Gradient Descent

Goal: $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$



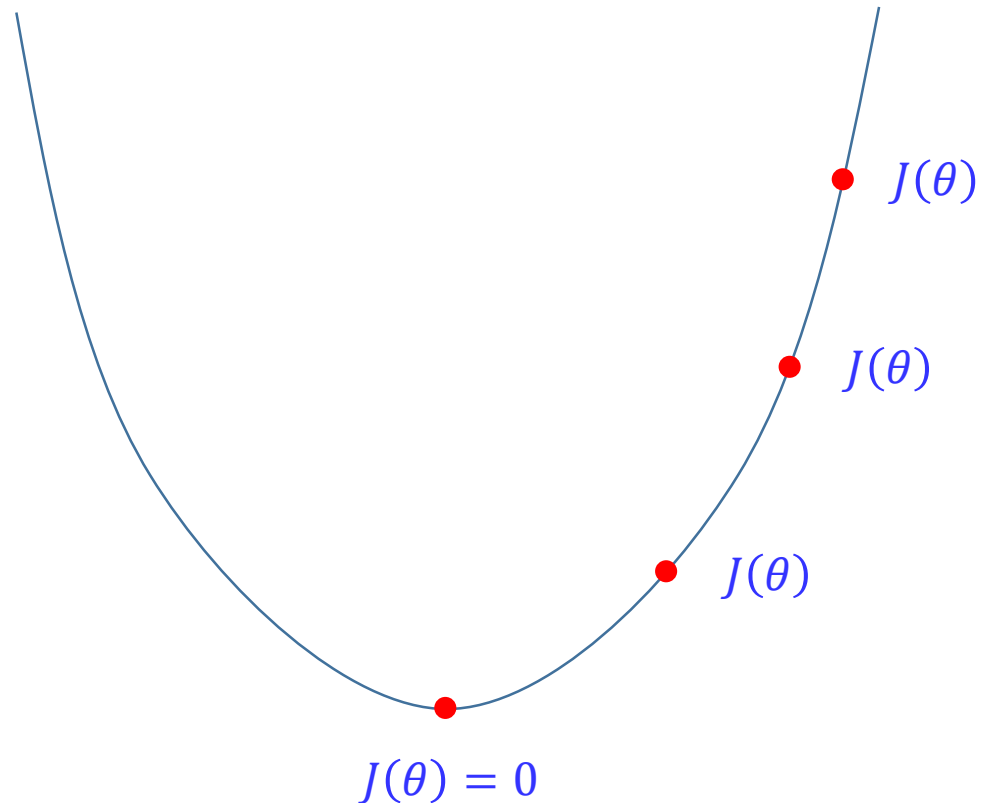
Cost Function With 2 Parameters



Cost Function With 1 Parameter

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

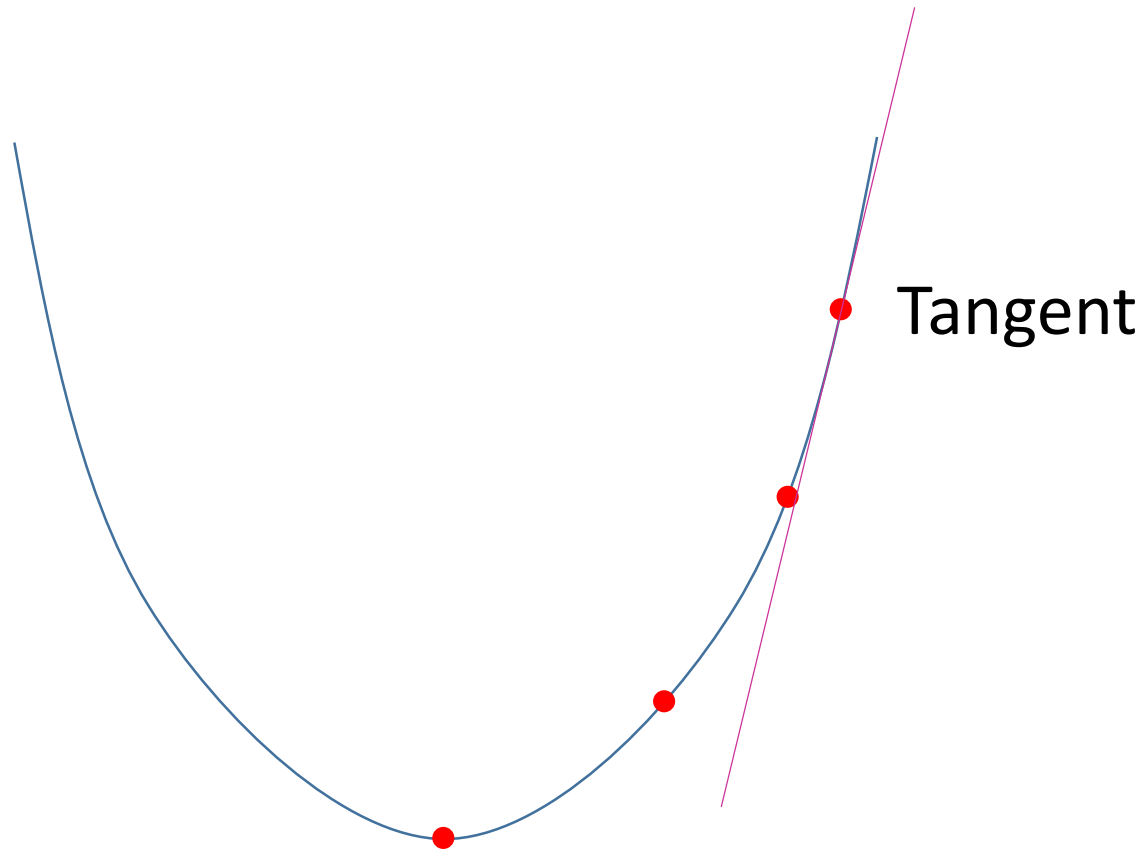
Will be using only one
Parameter (θ) to simplify
the explanation.



Gradient Descent

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

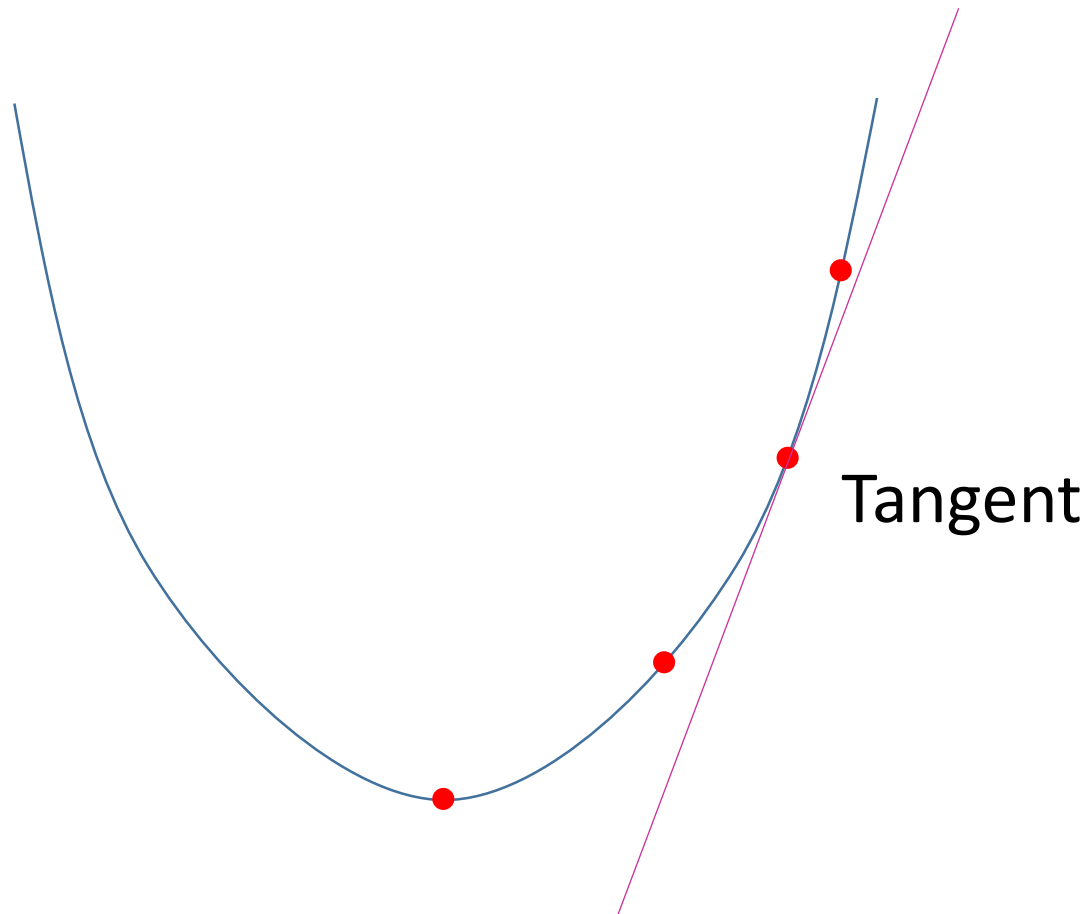
Will be using only one
Parameter (θ) to simplify
the explanation.



Gradient Descent

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

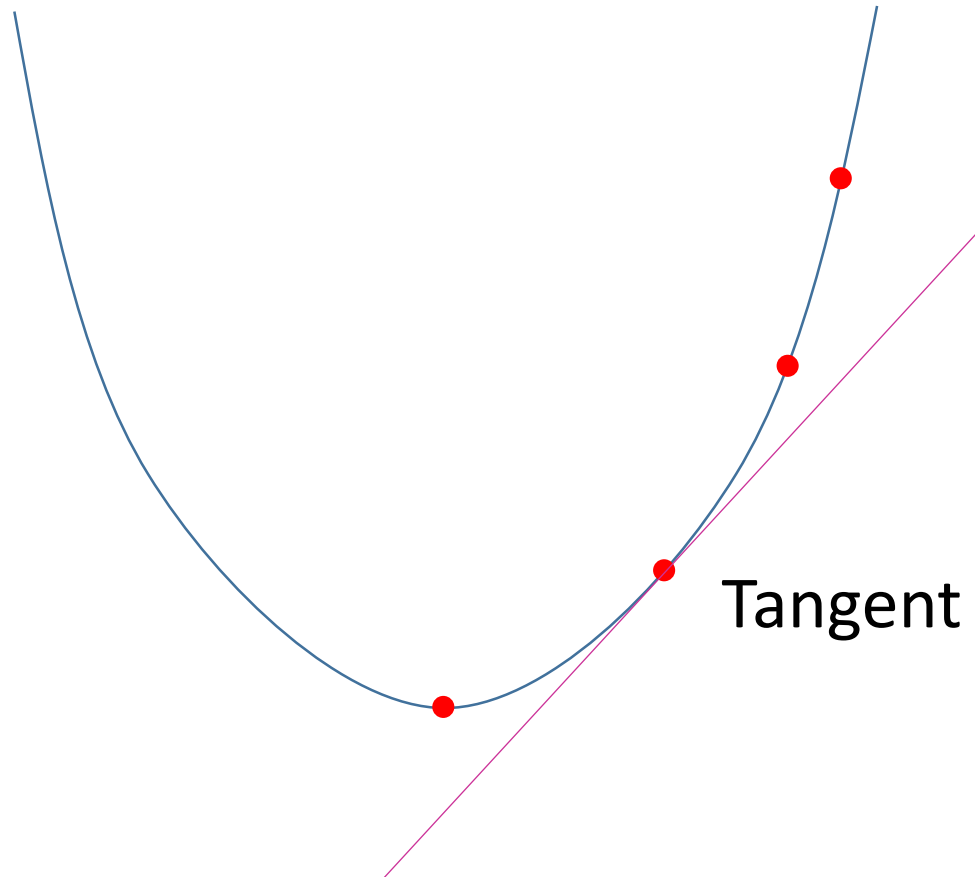
Will be using only one
Parameter (θ) to simplify
the explanation.



Gradient Descent

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

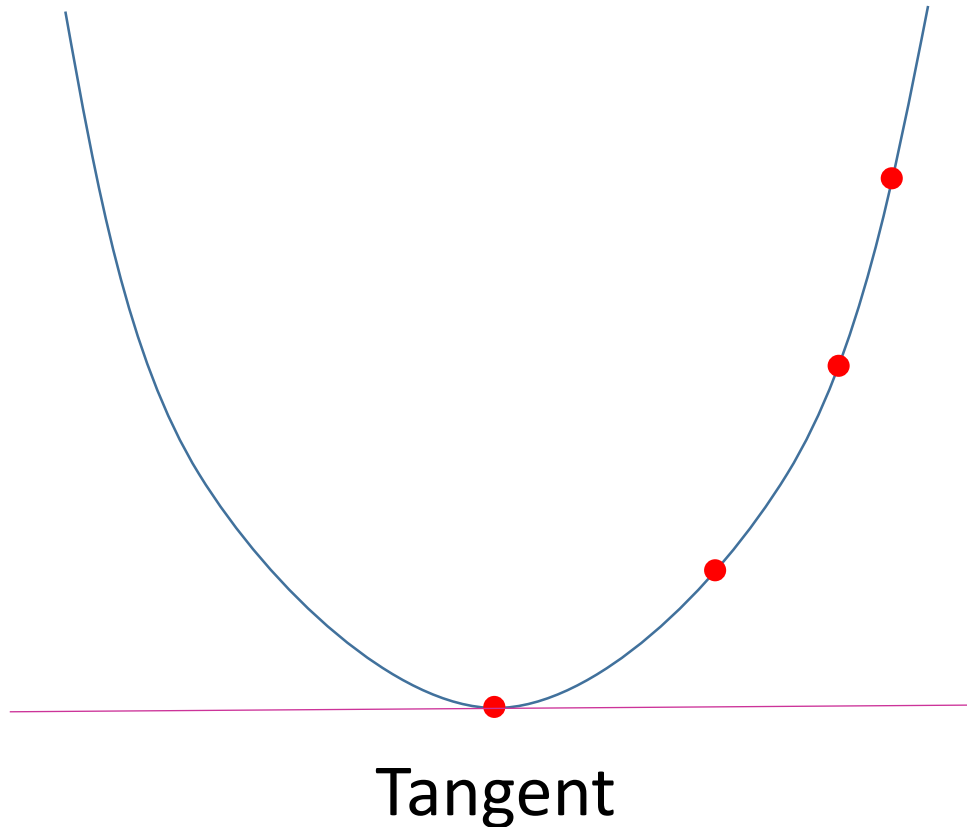
Will be using only one
Parameter (θ) to simplify
the explanation.



Gradient Descent

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

Will be using only one
Parameter (θ) to simplify
the explanation.



Gradient Descent

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

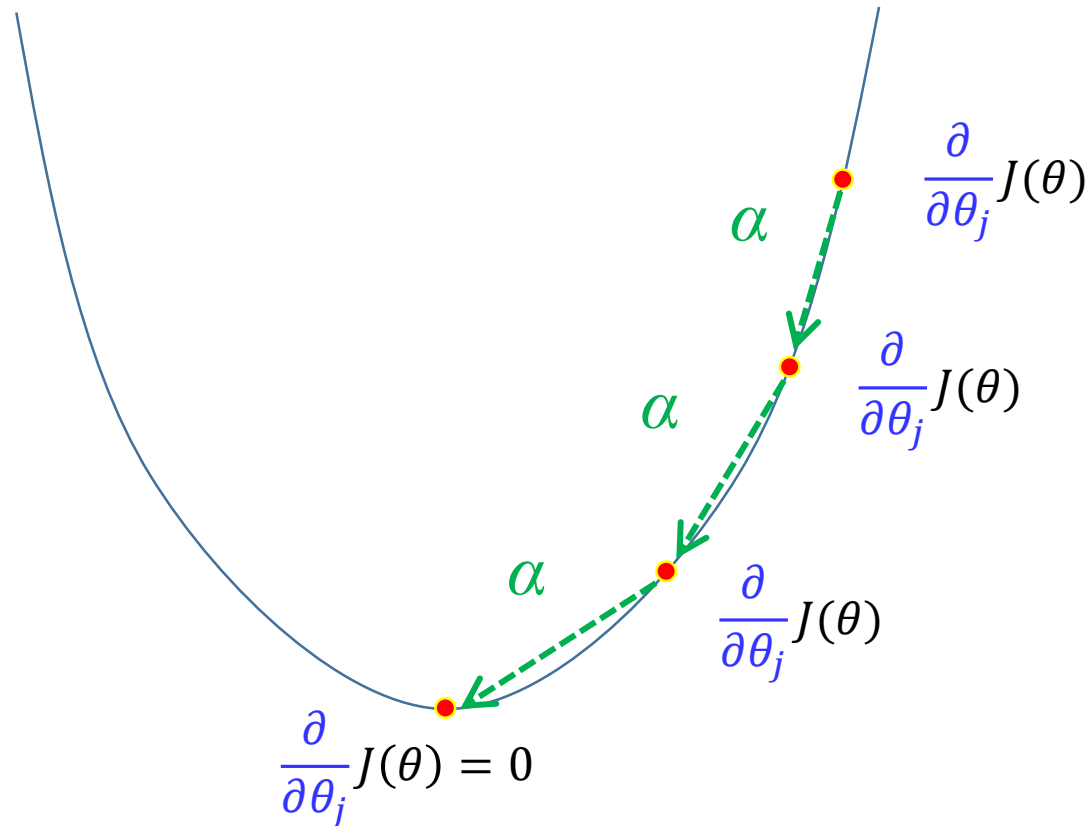
Tangent $\rightarrow \frac{\partial}{\partial \theta_j} J(\theta)$

Derivative

Partial (∂) / Total (d)

Gradient Descent

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



Gradient Descent Algorithm

*Resuming the
explanation with
2 Parameters*

$j = 0, j = 1$

repeat until convergence {

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

}

Gradient Descent Algorithm

$$j = 0, j = 1$$

repeat until convergence {

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

}

Update **MUST** be **SIMULTANEOUS**



Gradient Descent Algorithm Implementation

Incorrect

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

Correct

```
repeat until convergence {  
   $t0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
   $t1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
   $\theta_0 := t0$   
   $\theta_1 := t1$   
}
```

Gradient Descent Algorithm

$$j = 0, j = 1$$

repeat until convergence {

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

}

Gradient Descent Algorithm

$$j = 0, j = 1$$

repeat until convergence {

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

}



Update **MUST** be
SIMULTANEOUS

Gradient Descent Algorithm

$$j = 0, j = 1$$

repeat until convergence {

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

}

Update MUST be
SIMULTANEOUS

Learning
Rate

Gradient Descent Algorithm

$$j = 0, j = 1$$

repeat until convergence {

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

}

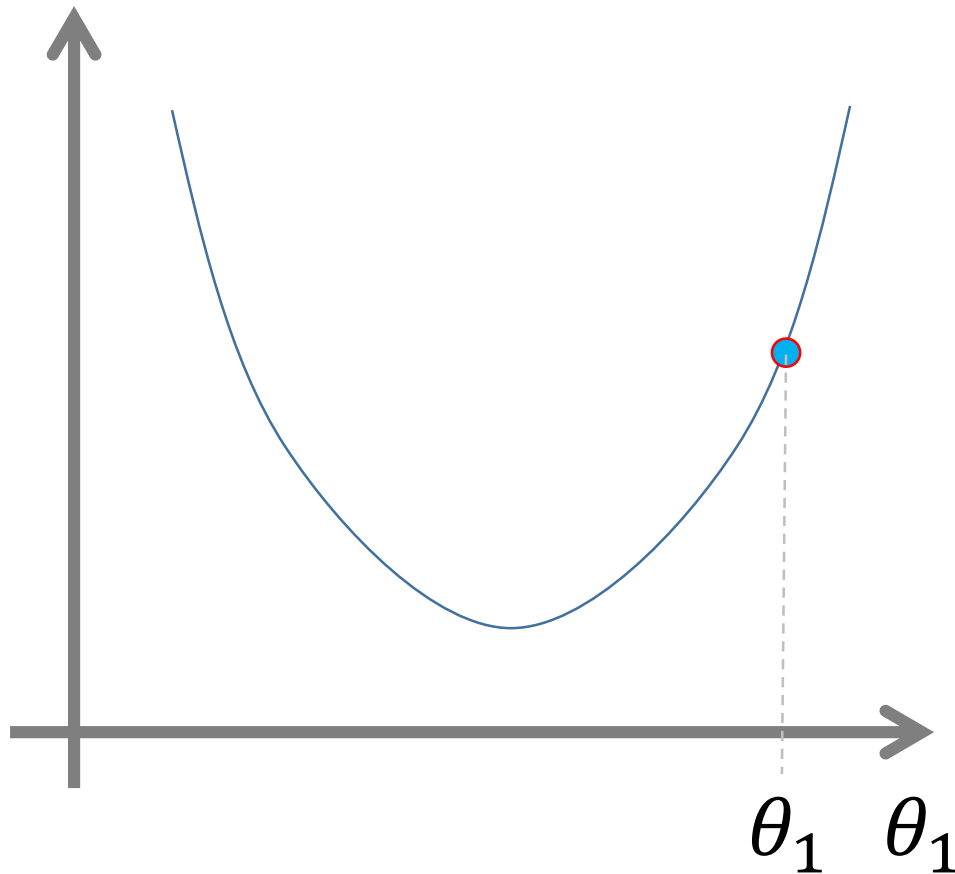
Update MUST be
SIMULTANEOUS

Learning
Rate

Tangent
Value

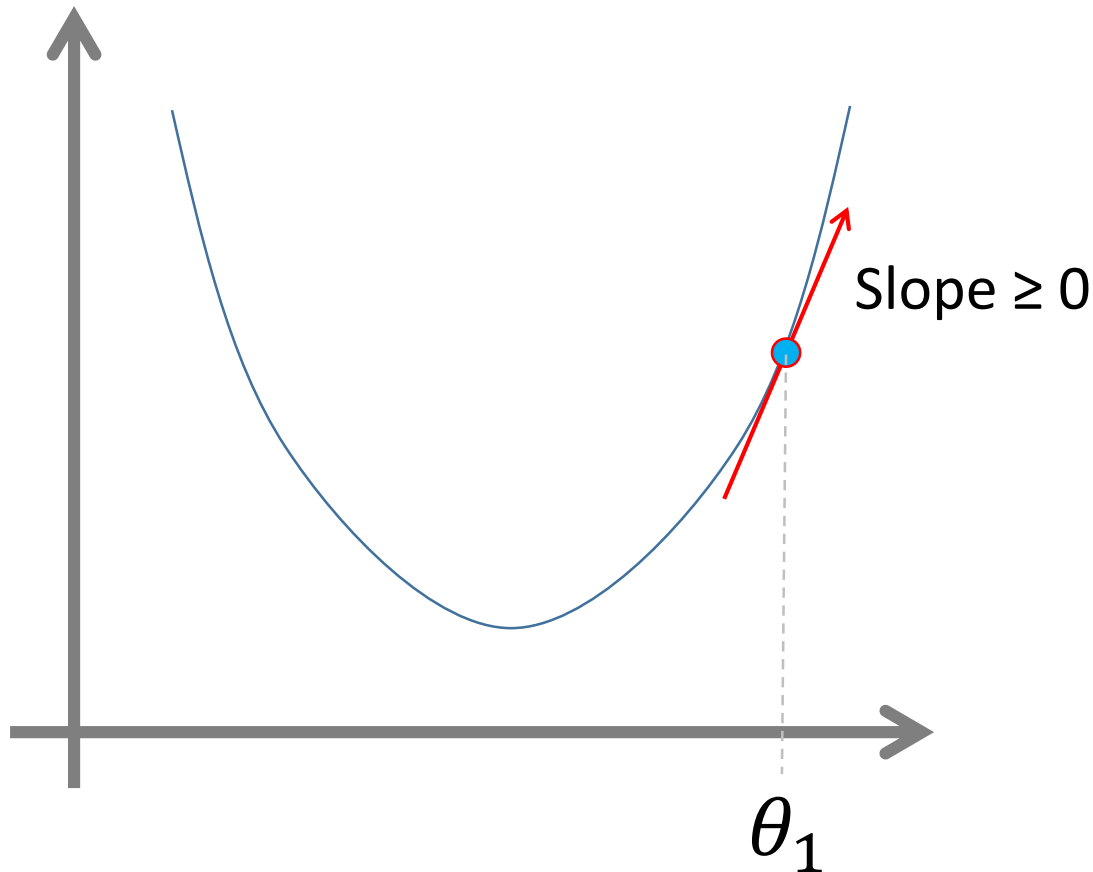
Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



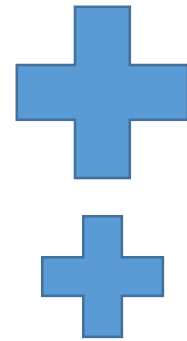
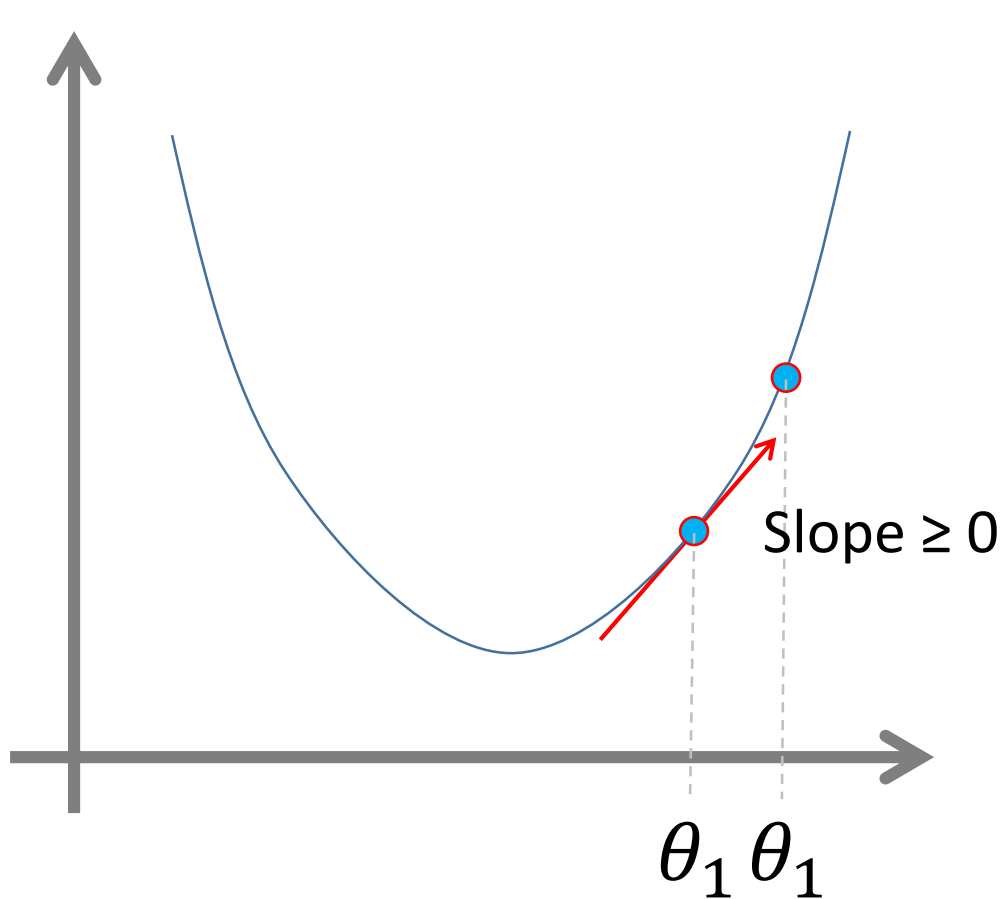
Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



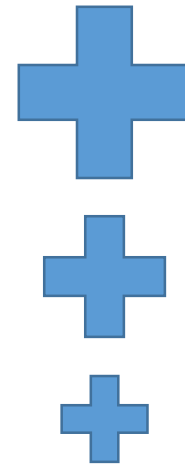
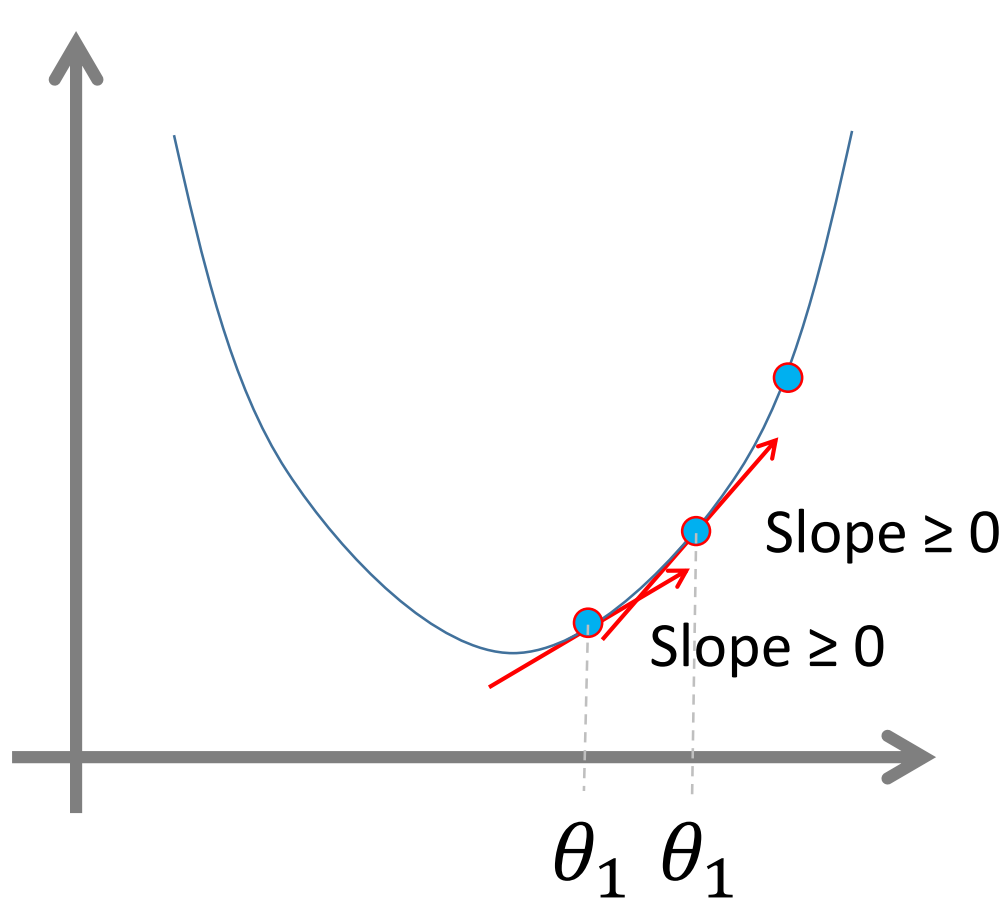
Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



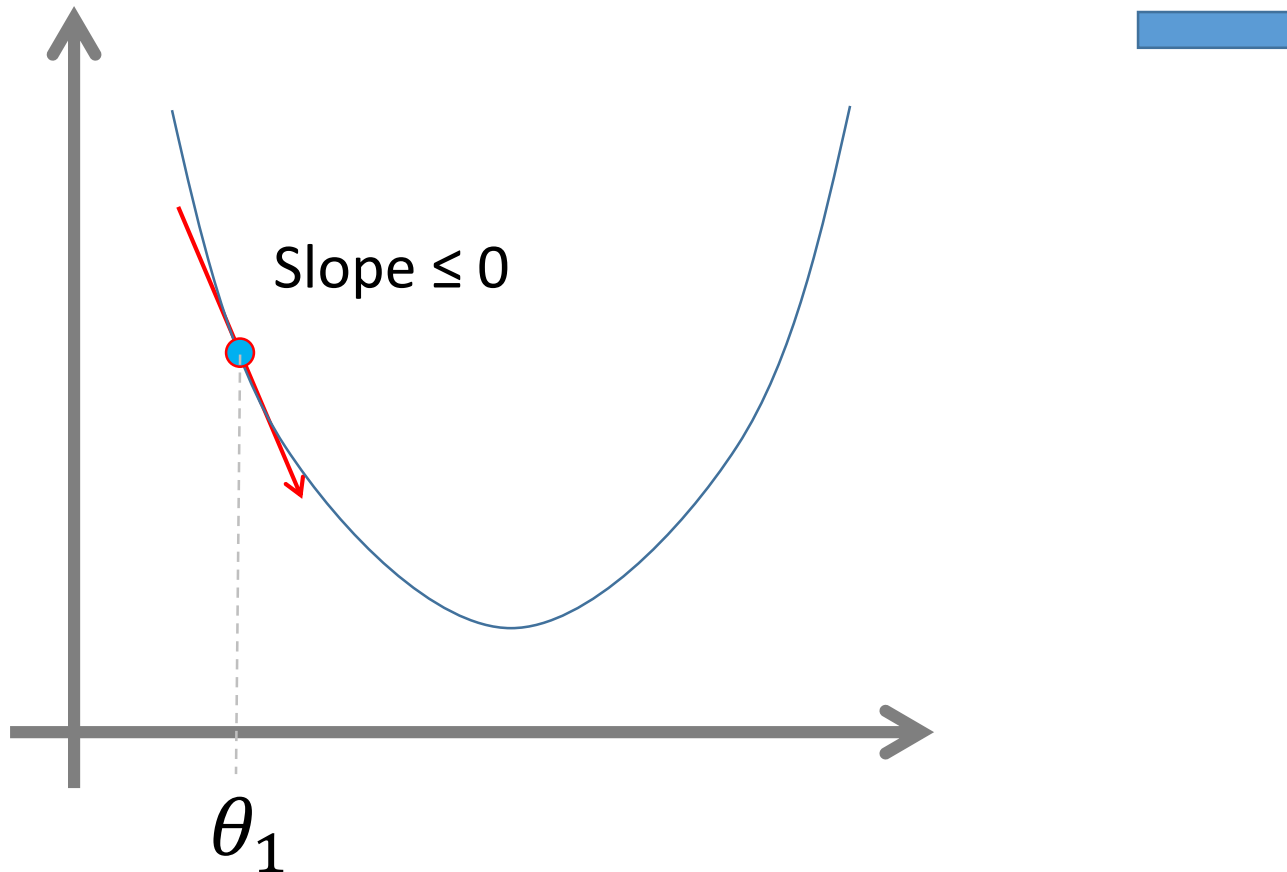
Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



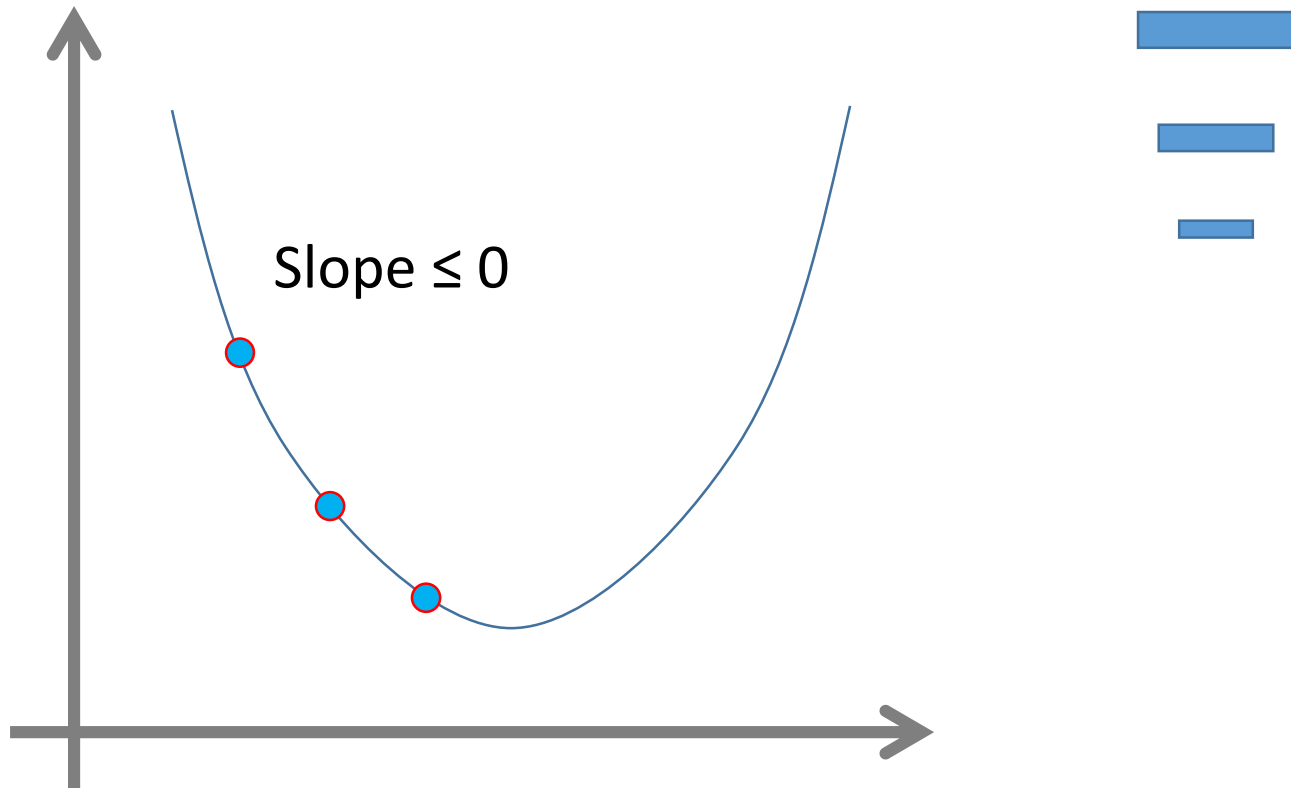
Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



Exercise

- As we approach a local minimum, will Gradient Descent automatically take smaller steps?
- Do we need to adjust the Learning Rate (α) manually over time?

Solution

- As we approach a local minimum, will Gradient Descent automatically take smaller steps?

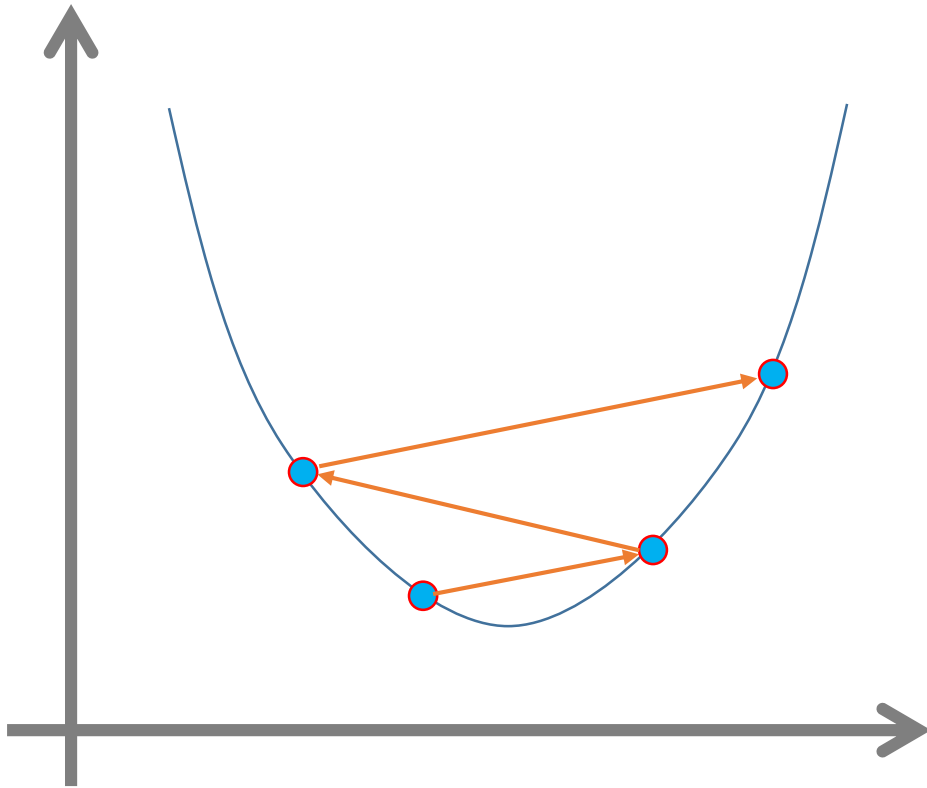
Yes

- Do we need to adjust the Learning Rate (α) manually over time?

No

Gradient Descent

$$\text{Let } \theta_0 := 0; \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



If α is too high?

Gradient Descent
can overshoot the
minimum; it may
fail to converge, or
even diverge.

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 Algorithm $\rightarrow \min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

$$j = 0, j = 1$$

repeat until convergence {

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

}

Derivative

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

Derivative

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

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

Derivative

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

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

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

Derivative

$$\frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 ; j \in \{0, 1\}$$



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

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

List of Derivative Rules

*[https://www.math.ucdavis.edu/~kouba/
Math17BHWDIRECTORY/Derivatives.pdf](https://www.math.ucdavis.edu/~kouba/Math17BHWDIRECTORY/Derivatives.pdf)*

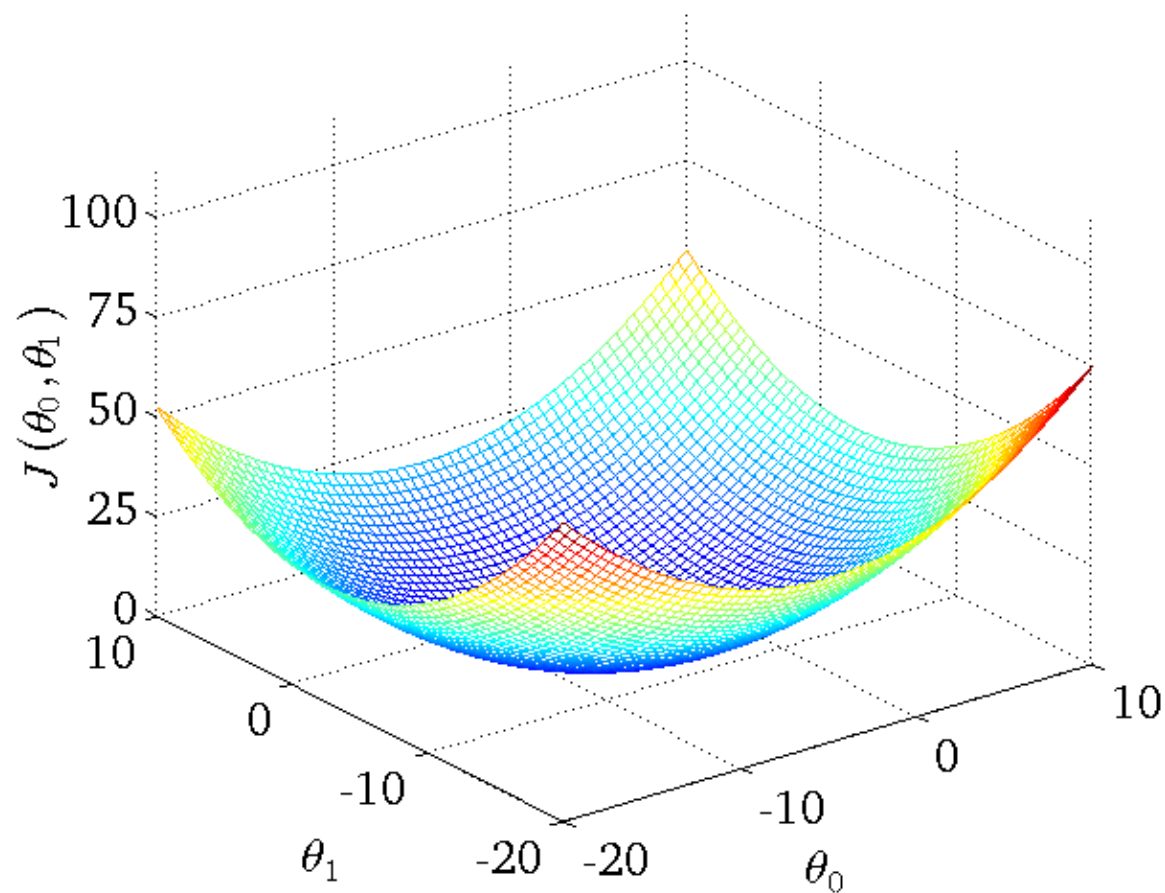
Power Rule: $f(x) = x^n$ then $f'(x) = nx^{n-1}$

Batch Gradient Descent

Batch = use **ALL** Training Examples
in each step of Gradient Descent

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

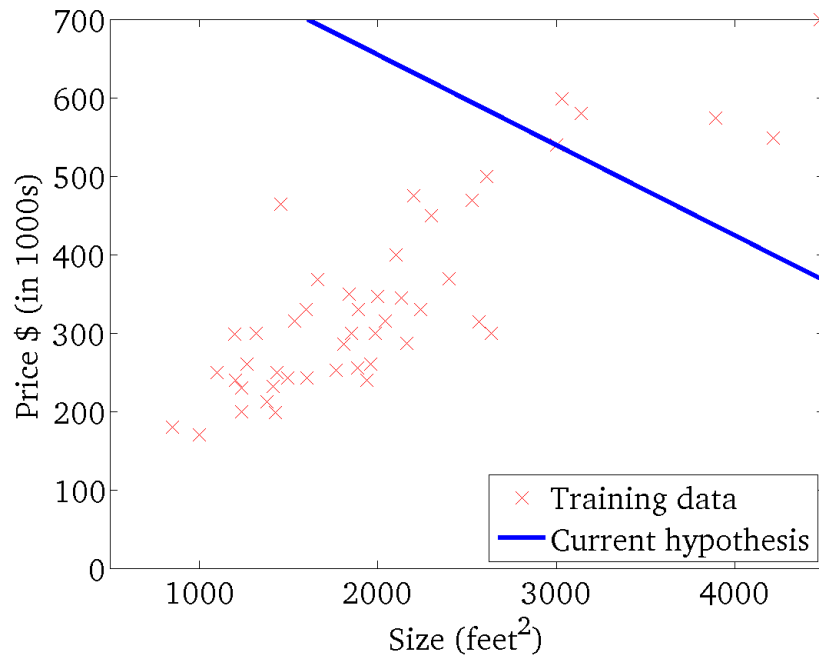
Batch Gradient Descent



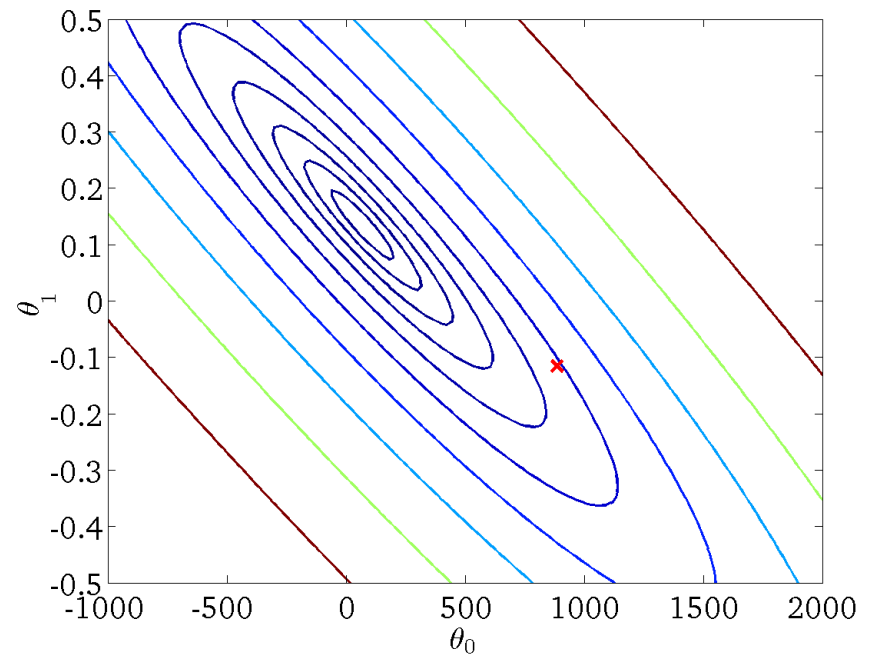
Batch Gradient Descent

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

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



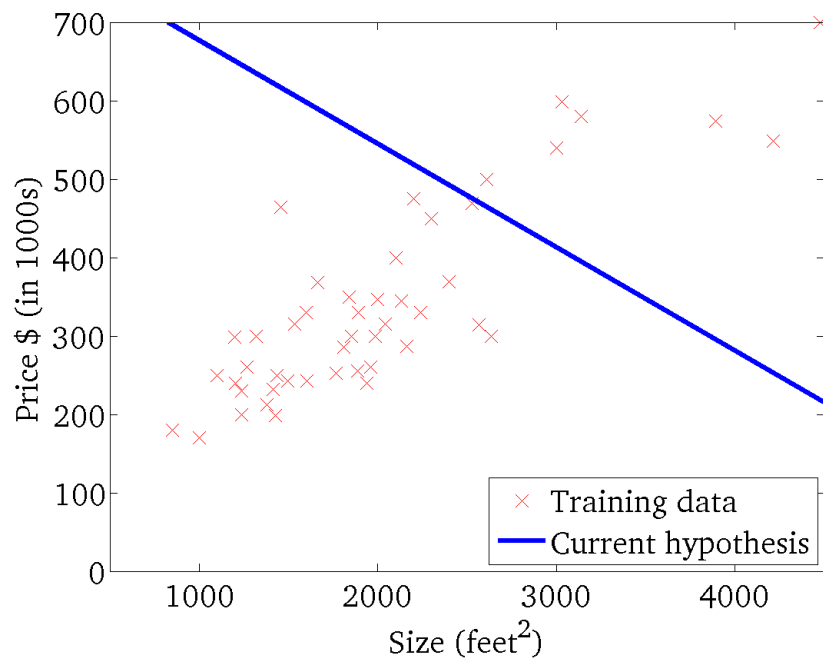
(function of the parameters θ_0, θ_1)



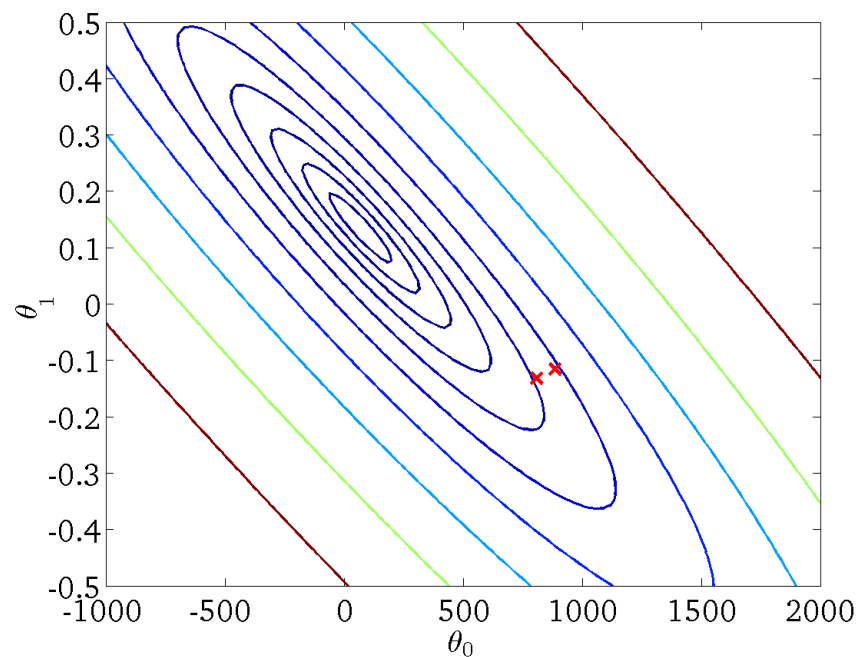
Batch Gradient Descent

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

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



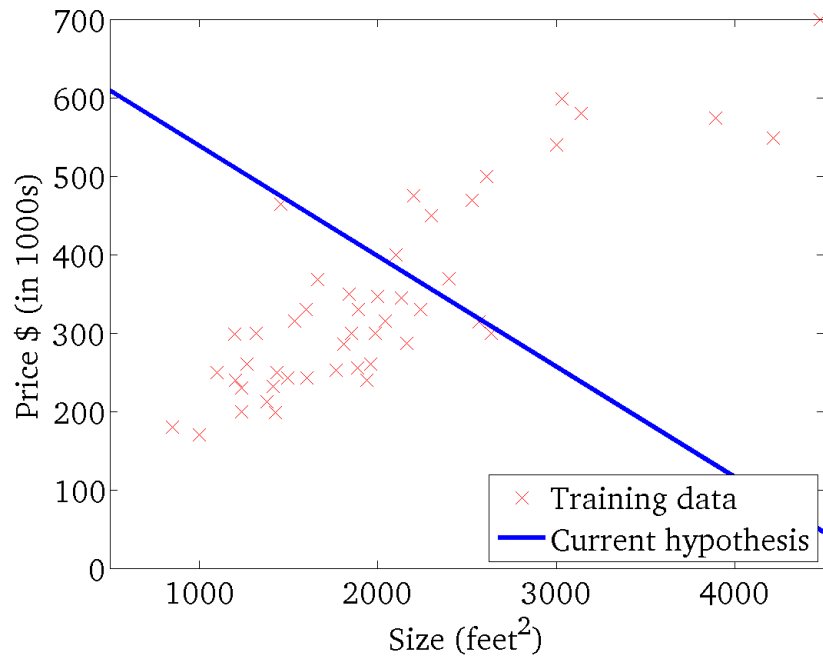
(function of the parameters θ_0, θ_1)



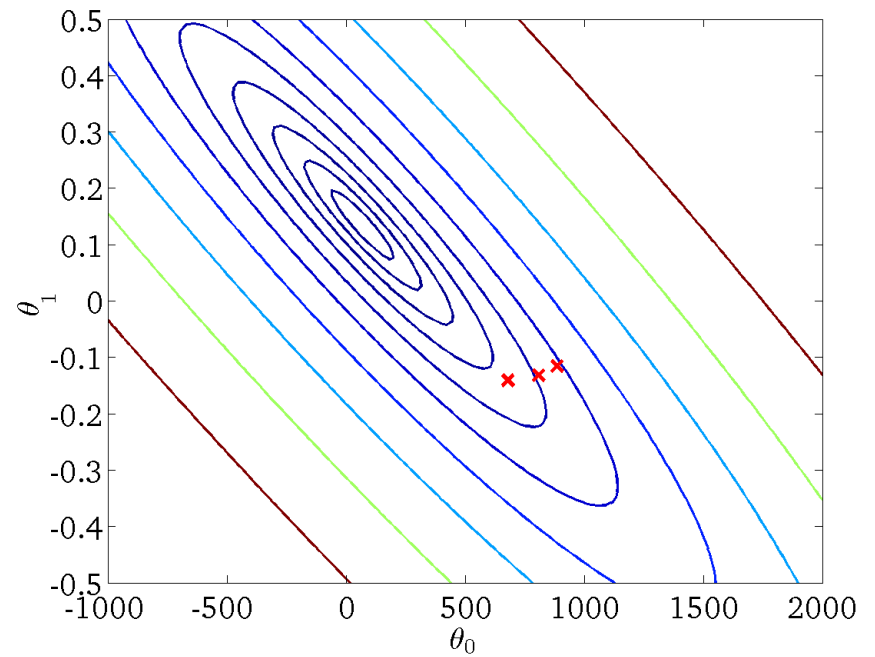
Batch Gradient Descent

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

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



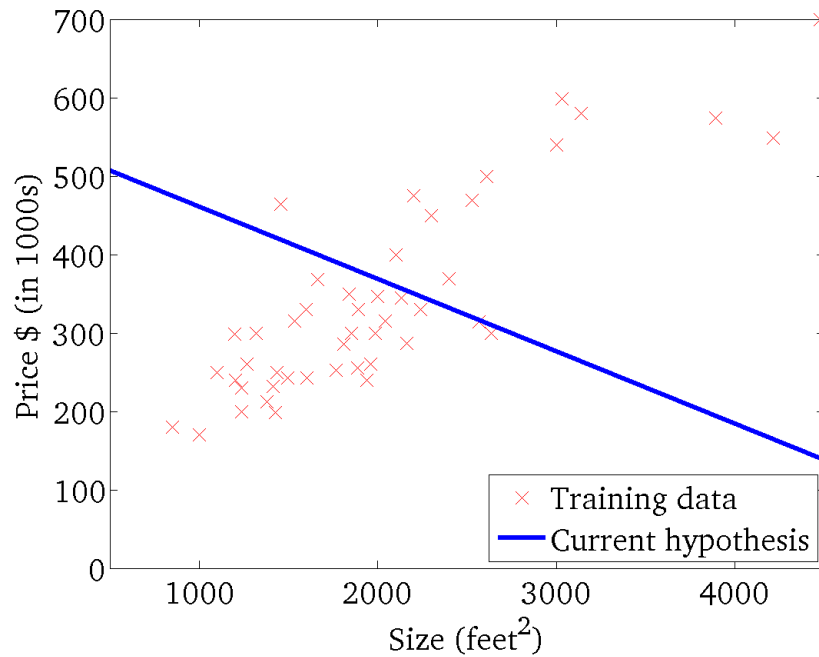
(function of the parameters θ_0, θ_1)



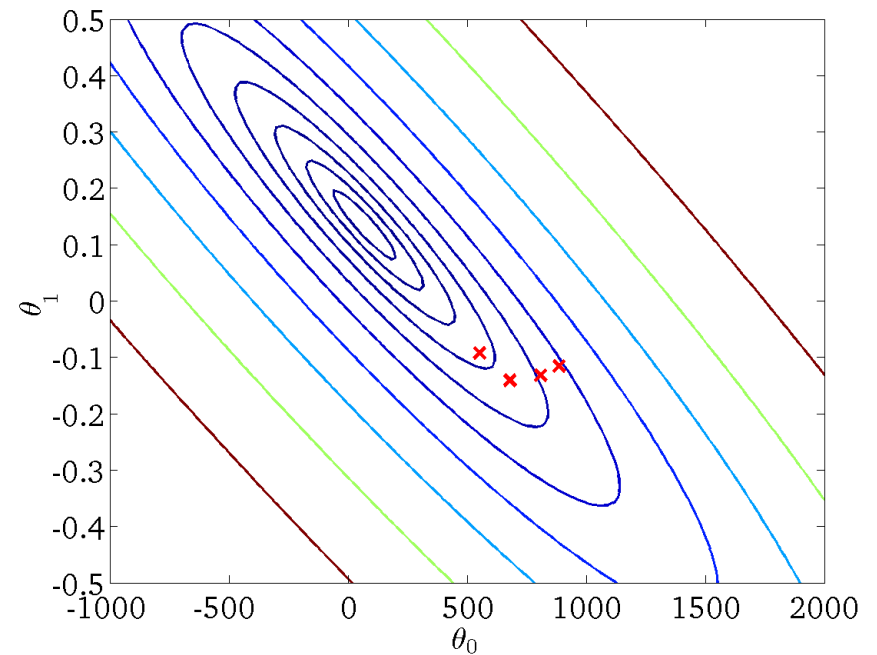
Batch Gradient Descent

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

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



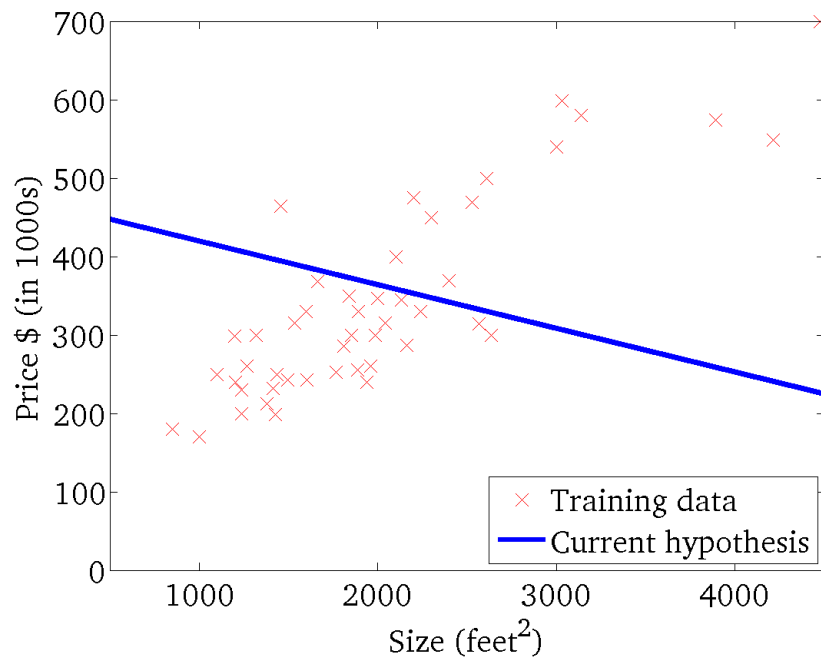
(function of the parameters θ_0, θ_1)



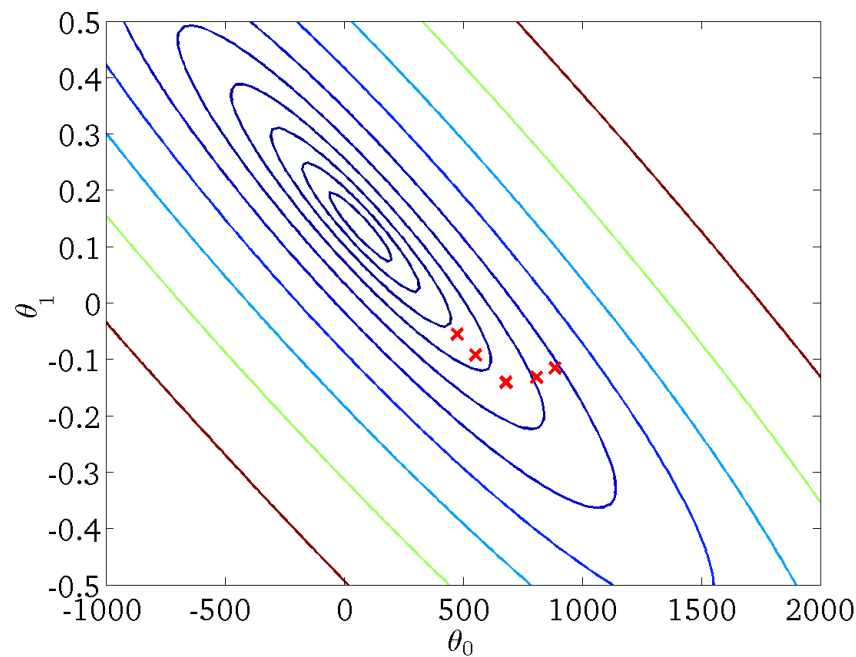
Batch Gradient Descent

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

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



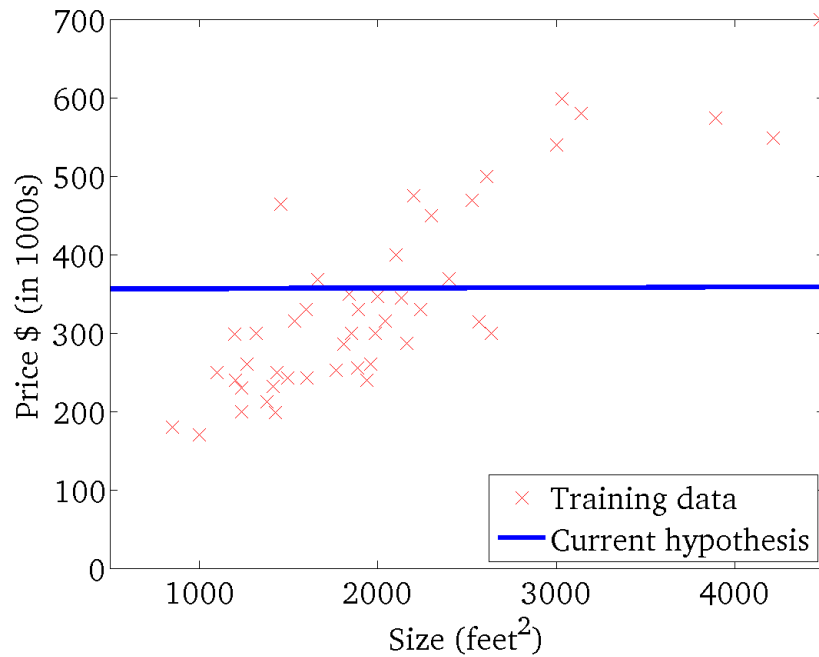
(function of the parameters θ_0, θ_1)



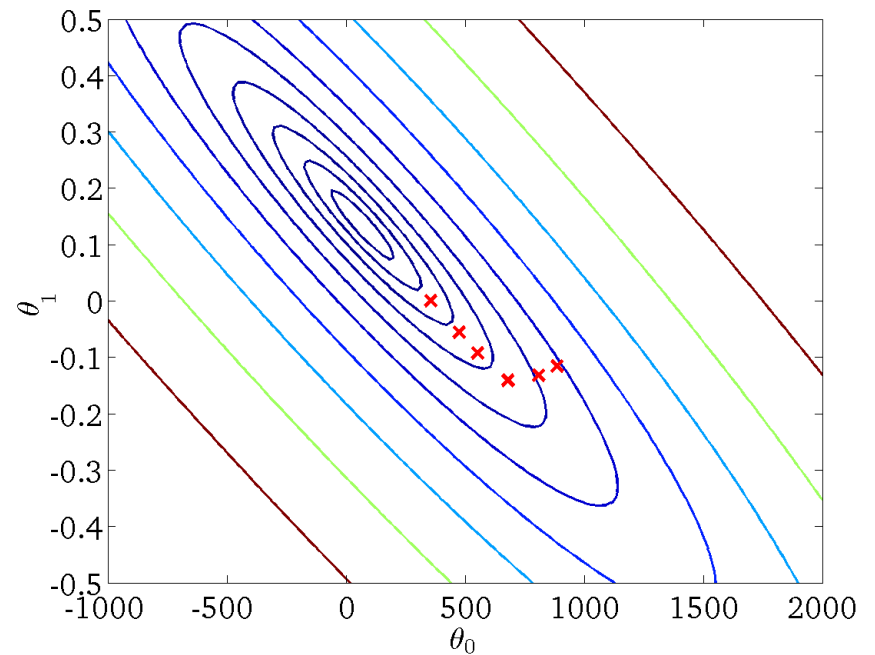
Batch Gradient Descent

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

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



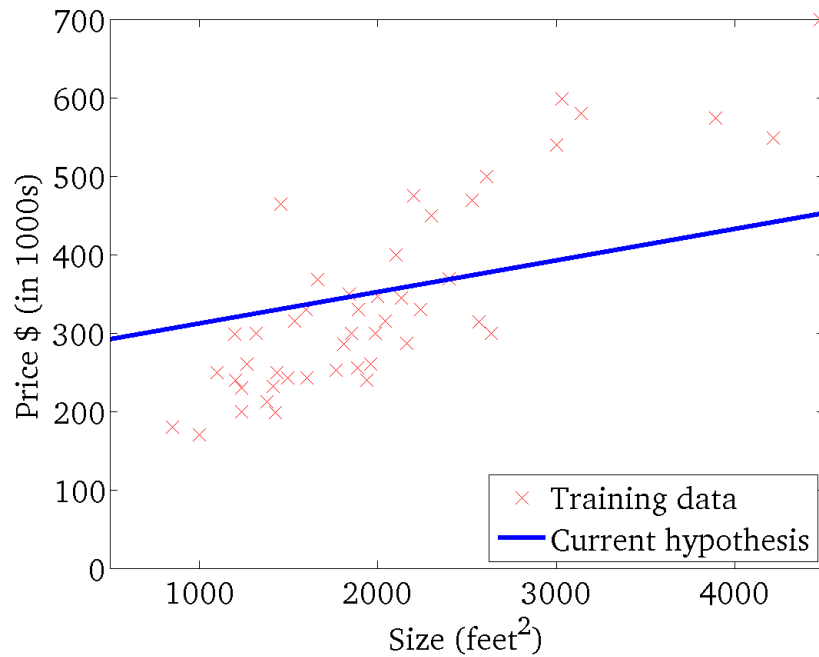
(function of the parameters θ_0, θ_1)



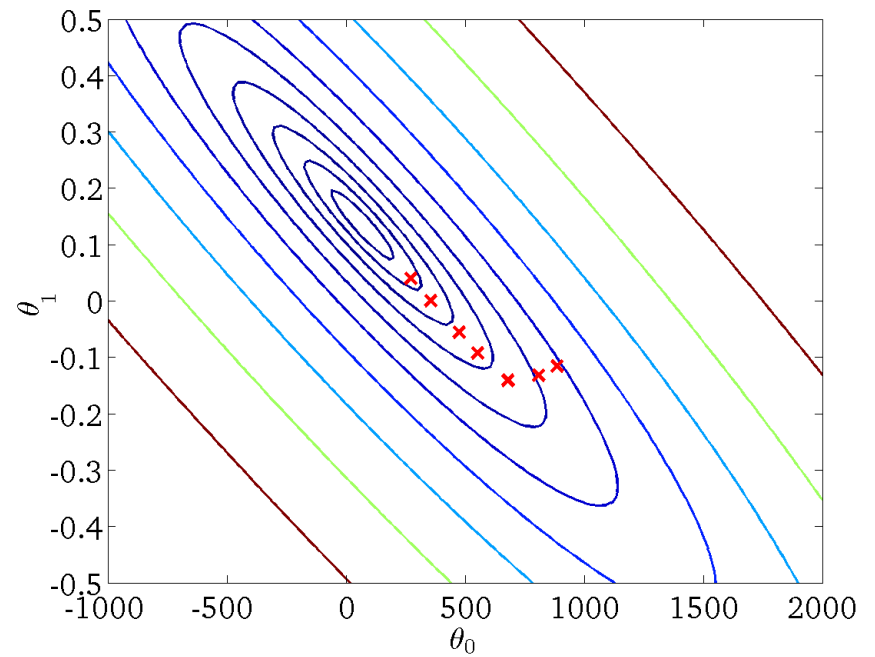
Batch Gradient Descent

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

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



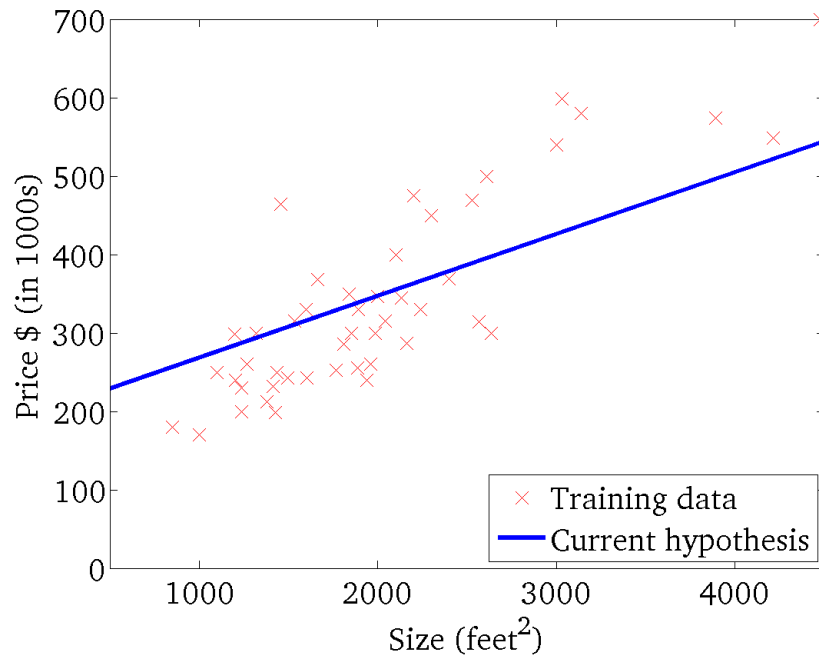
(function of the parameters θ_0, θ_1)



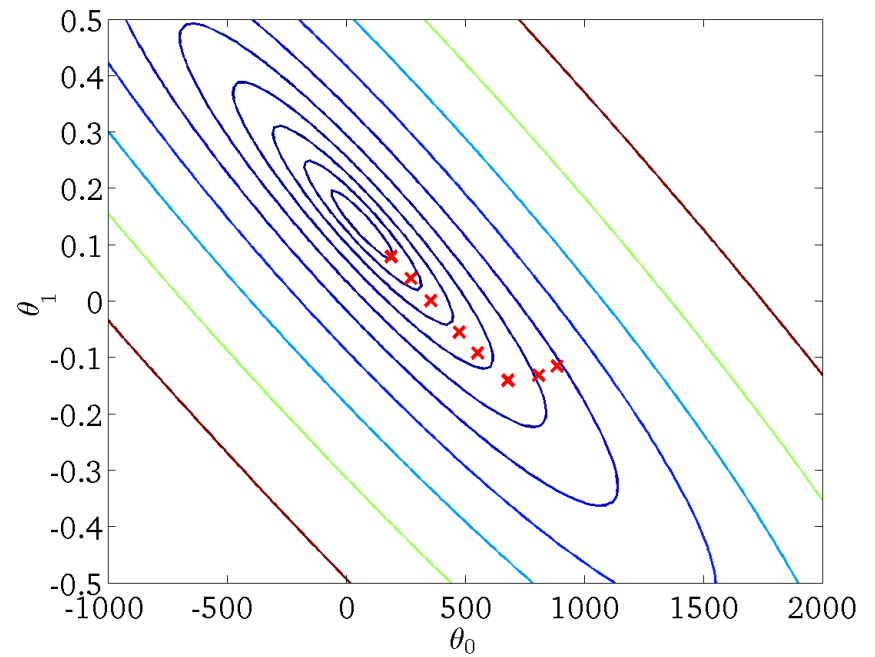
Batch Gradient Descent

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

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



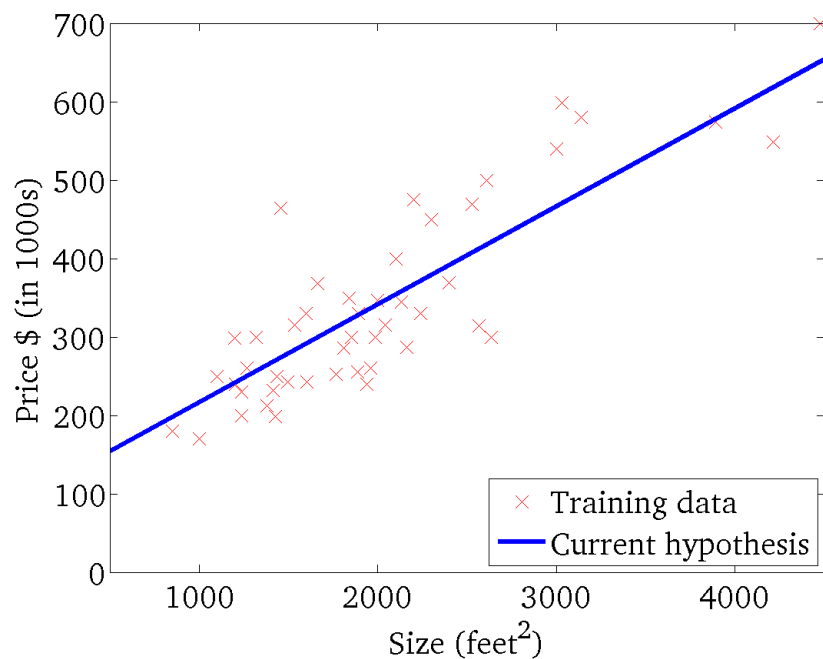
(function of the parameters θ_0, θ_1)



Batch Gradient Descent

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

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



(function of the parameters θ_0, θ_1)

