

# Gradient Descent

$\theta_0, \theta_1$

- Now we know how to evaluate a model, using a cost function, how do we make the model **learn** the optimal parameters ?
- In other words, how do we **minimize** the cost function **without testing all** the different possible models ?
- The algorithm used to do this is called **Gradient Descent**, and is **essential** to most machine learning algorithms, not just linear regression !
- In DL libraries this type of algorithm is called an **optimizer** and other variants exist.

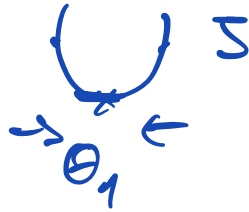
# Gradient Descent

- We have some function  $J(\theta_1, \theta_2)$  which we want to minimize...

Outline :

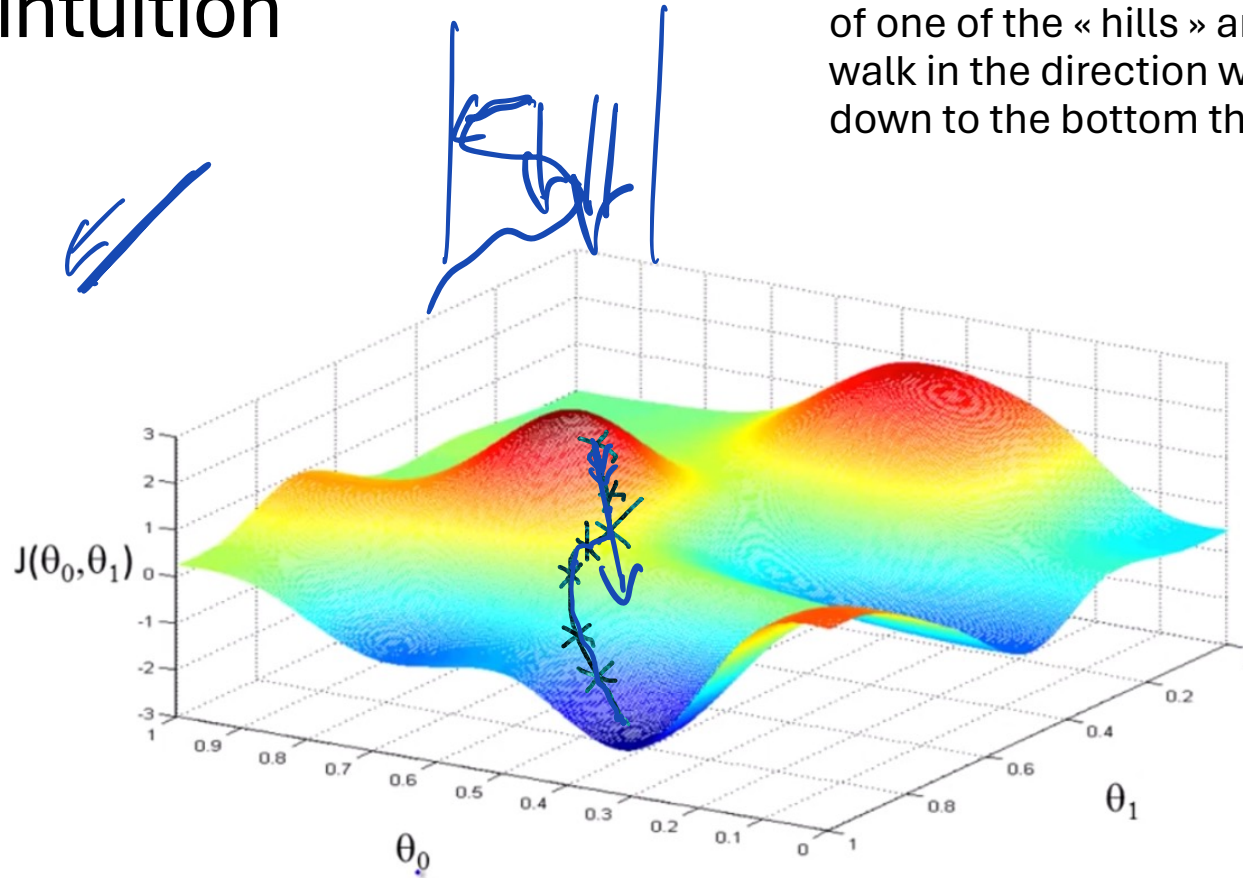


- Start with some initial guess, some random values for  $\theta_1, \theta_2$
- Keep **updating**  $\theta_1, \theta_2$  a little bit to reduce  $J(\theta_1, \theta_2)$  until we end up at a **minimum** (global or local)



# GD intuition

- This is your cost function in 3D
- Imagine you start somewhere near the top of one of the « hills » and your goal is to walk in the direction which will take you down to the bottom the fastest.



# GD formula

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

$$\theta_{j\_new} : \theta_j - \alpha$$

- This is the update formula for each of the parameters
- $:=$  signifies assignment
- $\alpha$  is a number called the **learning rate**. If  $\alpha$  is very **large**, then it corresponds to an **aggressive** learning procedure and big steps being taken « downhill » and vice versa.
- $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$  is a derivative term, which requires a bit of calculus

# GD Intuition

- Why does this update make sense ?
- Why are we putting those 2 terms together ?
- Let's try and get a basic understanding of derivatives before we go any further.

# Derivatives

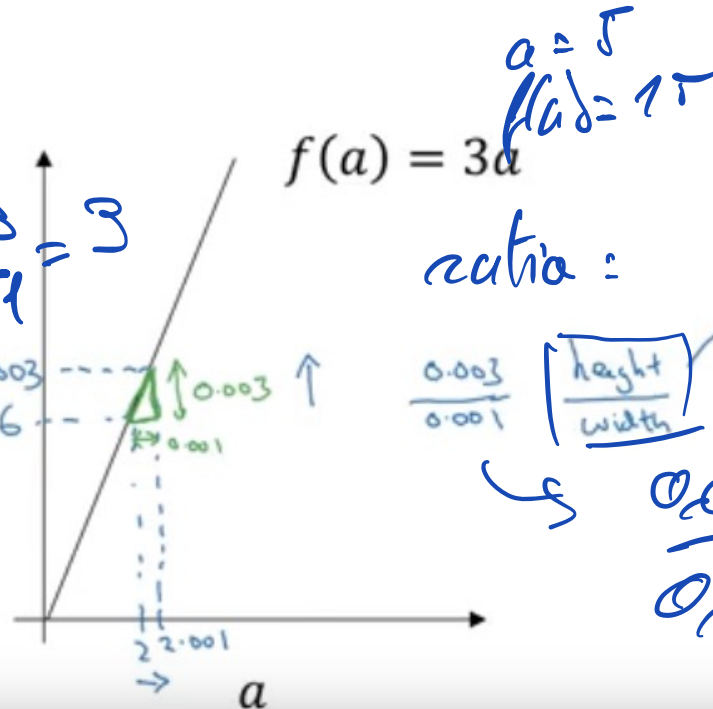
- The derivative describes **how the output of a function varies** with regard to a very very very **tiny positive nudge** to the **input**, to the point where we consider *almost* no variation in input....
- Informally, the derivative tells you how a function behaves at a particular « instant », i.e. for a given input value.
- The derivative is commonly referred to as « **instantaneous rate of change** »

# Derivatives

- Here is a linear function as an example. What happens when we shift the input by a 'small' value like 0.001
- when  $a=2$  ?
- when  $a=5$  ?

$$\frac{15.003 - 15}{5.001 - 5} = \frac{0.003}{0.001} = 3$$

(change in  $a$ )



$$\begin{aligned} a &= 2 & f(a) &= 6 \\ a &= 2.001 & f(a) &= 6.003 \end{aligned}$$

slope (derivative) of  $f(a)$  at  $a=2$  is 3

$$\frac{0.003}{0.001} = 3$$

# Derivatives

- With this function, we expect a small positive nudge in the input to make the output increase by 3 times the value of that nudge.

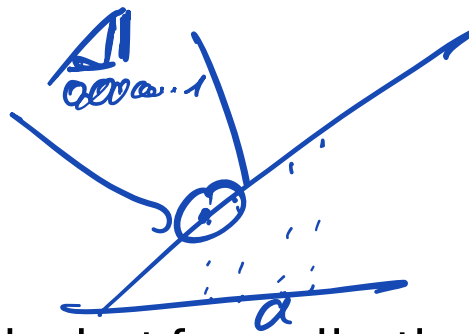
$$f(5.001) = 15.003$$

- In other words the **ratio** between the change in output and the change in input is 3 :

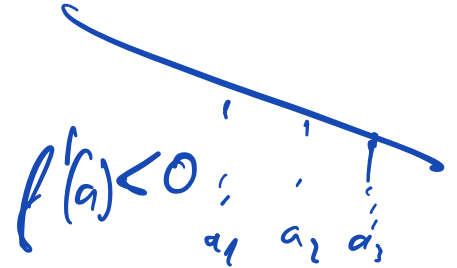
$$\frac{\text{change in } f(a)}{\text{change in } a} = \frac{df(a)}{da} = \frac{0.003}{0.001} = 3$$



# Derivatives



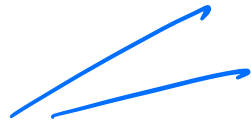
$$f'(a) > 0$$



- This is just an example, but formally, the derivative considers this ratio when the input is increased by a **much tinier** amount !
- $\text{Nudge} < 0.000000000000.....1 \Rightarrow$  nudge gets **as close to 0 as possible**
- In this previous example, whatever input value we pick, the derivative will be the same.
- This makes sense since the **function is a line** and the output increases at a constant rate
- Question : What if the derivative was negative everywhere ? What would the function look like ?

$$f'(a) \quad f'(a) = 3$$

$$f'(a) = 0$$

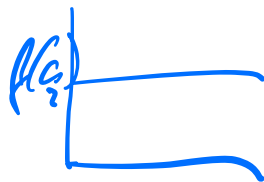


$$\frac{dy}{dx} > 0$$

$$f'(a) = -5$$

$$f(a) = 2a + 2$$

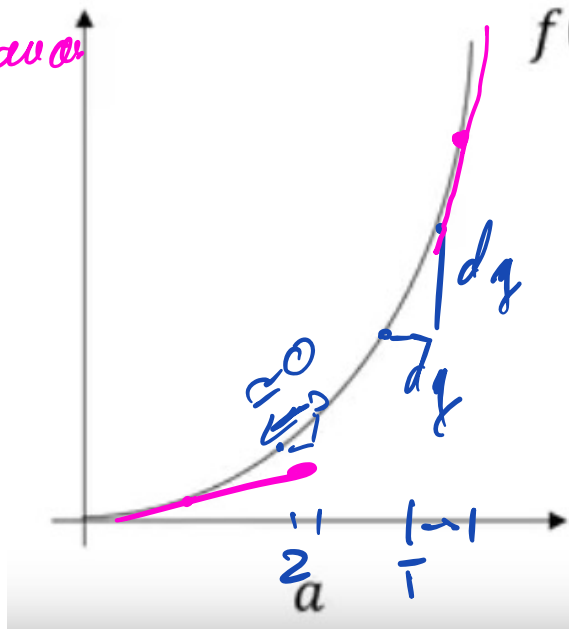
# Derivatives



$$\frac{dy}{dx}$$

- What if our function isn't a line ?

derivative for:  
rate of change



$$f(a) = a^2$$

$$a = 2$$

$$a = 2.001$$

$$f(a) = 4$$

$$f(a) \approx \underline{4.004}$$

$$(4.004004)$$

$$f'(a=2) = \frac{0.004}{0.001} = 4$$

$$f'(a=2) = \frac{f(a=2.001) - f(a=2)}{2.001 - 2}$$

$$= \frac{4.004004 - 4}{0.001}$$

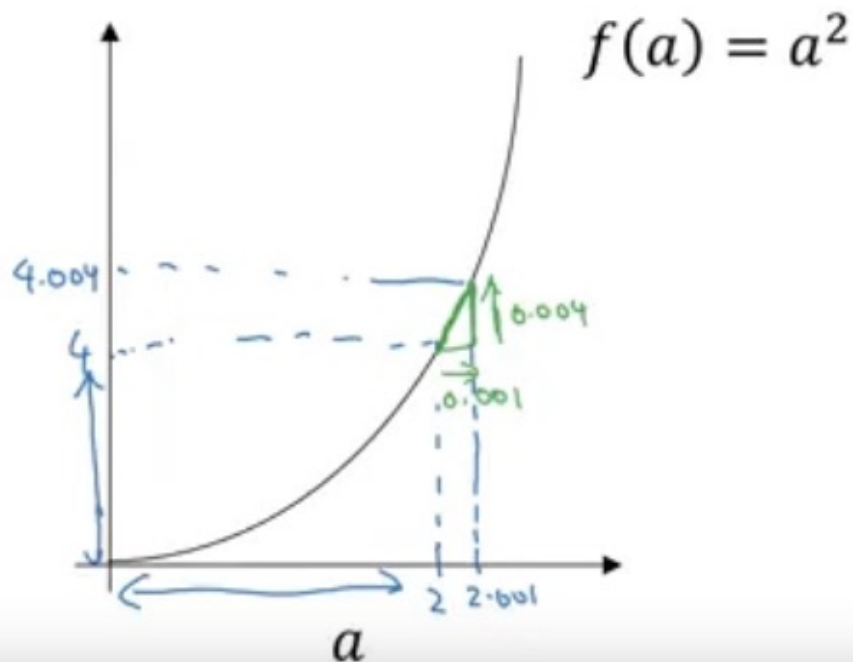
$$= \frac{0.004004}{0.001} = 4.004004$$

$$f'(2) = 4$$

$$f'(5) = 10$$

$$= \frac{0,01}{0,001} = \boxed{10}$$

- The derivative at  $a=2$  is ...



$$a=2 \quad f(a)=4$$

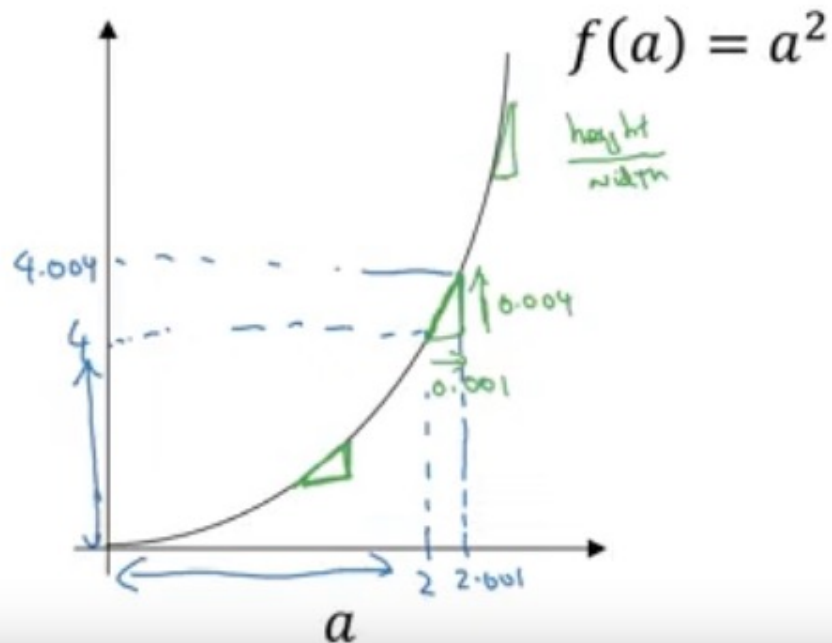
$$a=2,001 \quad f(a) \approx 4,004$$

$$(4,004001)$$

slope (derivative) of  $f(a)$  at  $a=2$  is  $4$ .

$$\frac{d}{da} f(a) = 4 \quad \text{when } a=2.$$

- The derivative at  $a=5$  is ...



$$\begin{aligned}
 a &= 2 & f(a) &= 4 \\
 a &= 2.001 & f(a) &\approx 4.004 \\
 & & & (4.004004) \\
 \text{slope (derivative) of } f(a) \text{ at } & & & \\
 a=2 & \text{ is } 4. \\
 \frac{d}{da} f(a) &= \underline{4} \text{ when } a=\underline{2}. \\
 a &= 5 & f(a) &= 25 \\
 a &= 5.001 & f(a) &\approx 25.010 \\
 \frac{d}{da} f(a) &= \underline{10} \text{ when } a=\underline{5}.
 \end{aligned}$$

- Rules exist to compute derivatives
- For example, the function

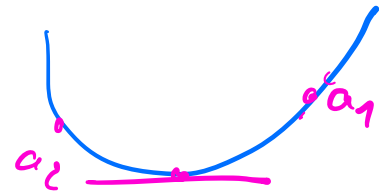
$$f(a) = a^2 \quad \boxed{na^{n-1}} \quad f(a^2) = 2a^1 = 2a$$

$$f'(a) = \frac{d}{da} f(a) = \boxed{2a}$$

(The notations are called Lagrange and Leibniz notations and are both common)

- If we look at the derivatives/slopes/ratios we calculated previously, this does indeed seem to work !
- Note: the derivative is equal to the **slope of the tangent line** on the graph **at our input value**.

$f'(a) > 0$  :  $f$  is increasing at that point  $a$



$f'(a_2) < 0$  : decreasing  
 $f'(a_3) = 0$  : constant  
 $f'(a_4) > 0$  : increasing

# Derivatives: (optional)



$$f'(x) = \frac{df}{dx} = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x} \quad \text{diff. } x$$

Example 2:  $f(x) = x^2$   $f'(x) = 2x$

$$\begin{aligned}
 \frac{df}{dx} &= \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x} \\
 &= \lim_{\Delta x \rightarrow 0} \frac{x^2 + (2x\Delta x + (\Delta x)^2) - x^2}{\Delta x} \\
 &= \lim_{\Delta x \rightarrow 0} \frac{2x\Delta x + (\Delta x)^2}{\Delta x} \\
 &= \lim_{\Delta x \rightarrow 0} (2x + \Delta x)
 \end{aligned}$$

- As  $\Delta x$  approaches 0, the derivative
- Approaches  $2x$ .

# GD Intuition

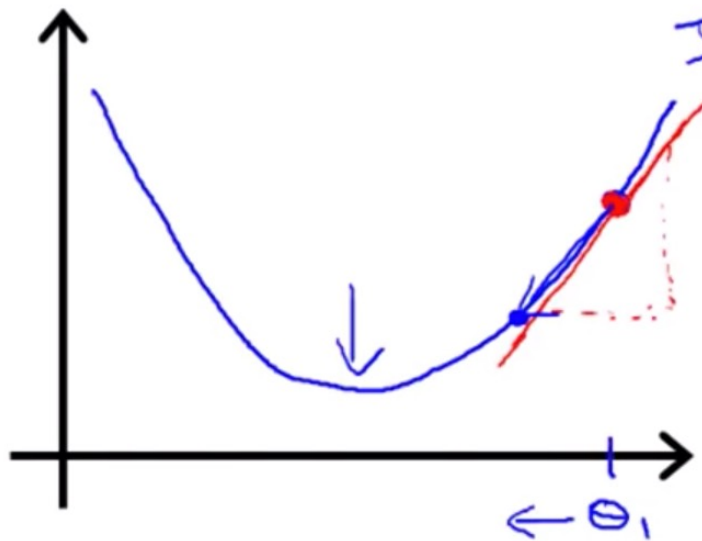
- Now we have a basic understanding of derivatives, let's apply this understanding to the gradient descent algorithm by using a **simpler example**, with a **cost function of only 1 single parameter**.
- We use  $J(\theta_1)$  instead of  $J(\theta_0, \theta_1)$
- Let's look at a couple scenarios to see how Gradient Descent updates our parameter  $\theta_1$ .

$$\theta_1 \leftarrow \theta_1 - \frac{dJ(\theta_1)}{d\theta_1}$$

When the derivative is **positive**...

- Remember, our cost function looks like a parabola.
- When  $\theta_1$  is too high, we want our optimizer to **reduce** this parameter and bring it closer to the « **sweet spot** », where the cost is minimized.
- Let's see if it does the right thing :

minimize  $J(\theta_1)$



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

$$\theta_1 := \theta_1 - \alpha$$

$$\boxed{\frac{\alpha}{d\theta_1} J(\theta_1)} \geq 0$$

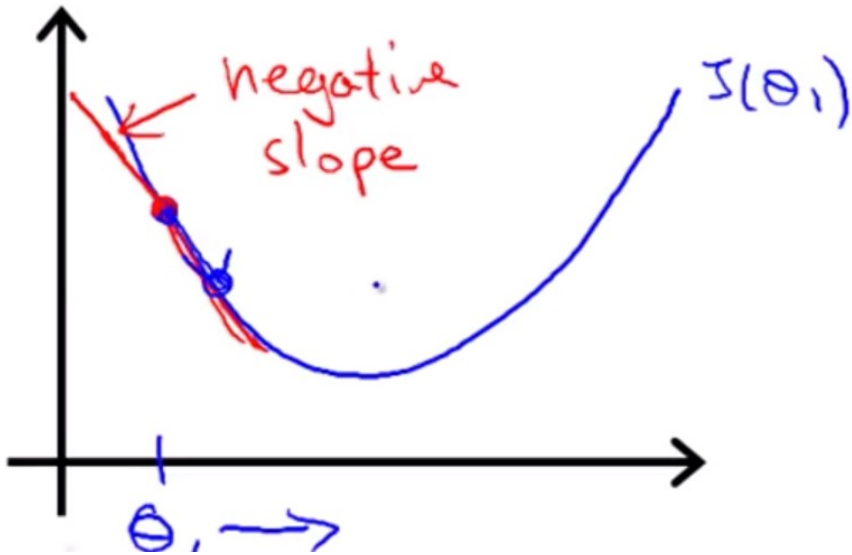
$$\theta_1 := \theta_1 - \underline{\alpha} \cdot (\text{positive number})$$

$$\theta_1 = 5 - 1.2 = 3.8$$



When the derivative is negative...

- When  $\theta_1$  is too low, let's see if Gradient Descent **increases** it and brings it closer to the « sweet spot », where the cost is minimized :

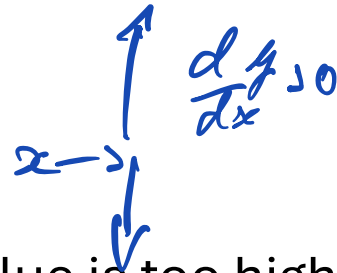


$$\frac{\frac{d}{d\theta_1} J(\theta_1)}{\leq 0}$$

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

$\uparrow$                        $\uparrow$   
 $- \times - 2$   
 $+ \sqrt{2}$

# Recap



Cost is going up

- When the parameter value is too high, the (derivative is positive) and the update rule decreases the value for the parameter.

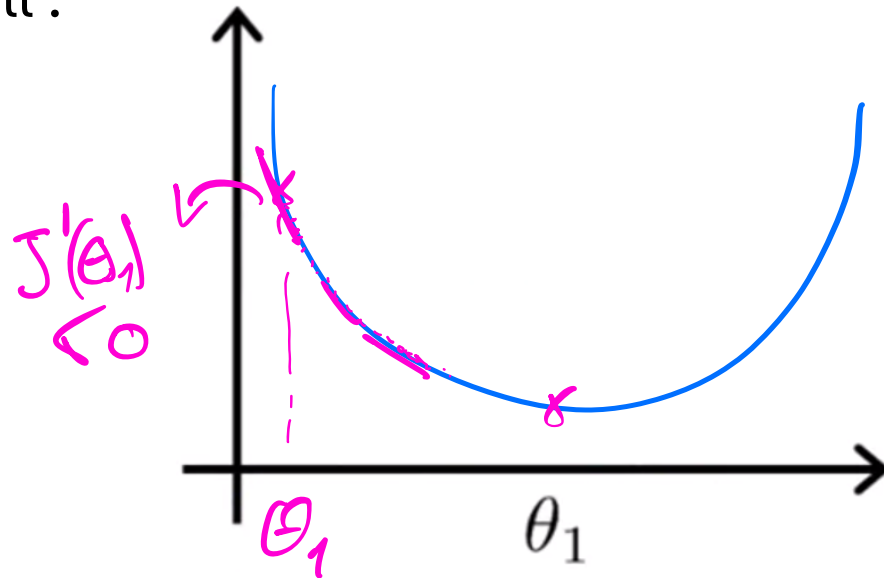
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) > 0$$

- Conversely, when the parameter value is too low, the parameter value will be increased by the update rule.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) < 0$$

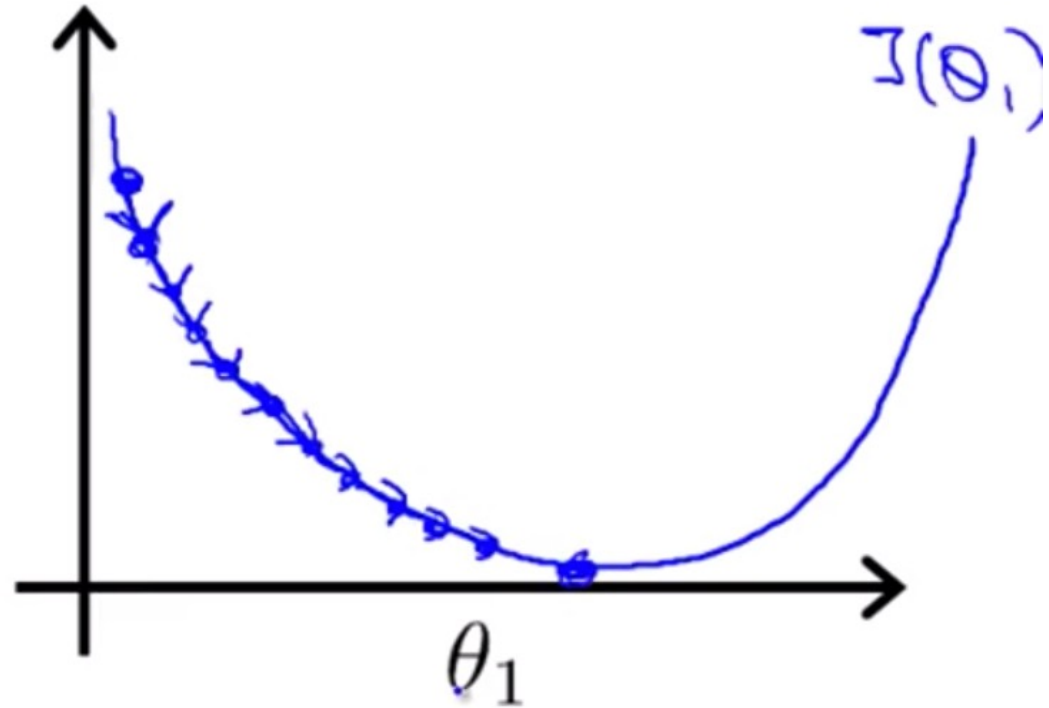
# Okay so now what about $\alpha$ ?

- Remember the update rule :  $\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$
- How does  $\alpha$  influence the update of our parameter  $\theta_1$  ?
- If  $\alpha$  is too small :



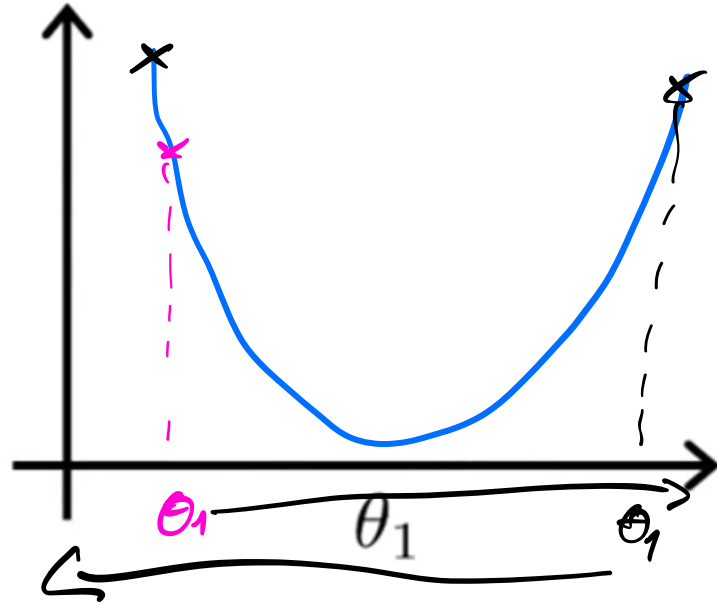
# If $\alpha$ is too small

- Many small steps will be taken, which makes Gradient Descent very slow



# If $\alpha$ is too large...

- Gradient descent may « overshoot », go past the minimum. It may even never converge (never find the minimum) and keep jumping around.

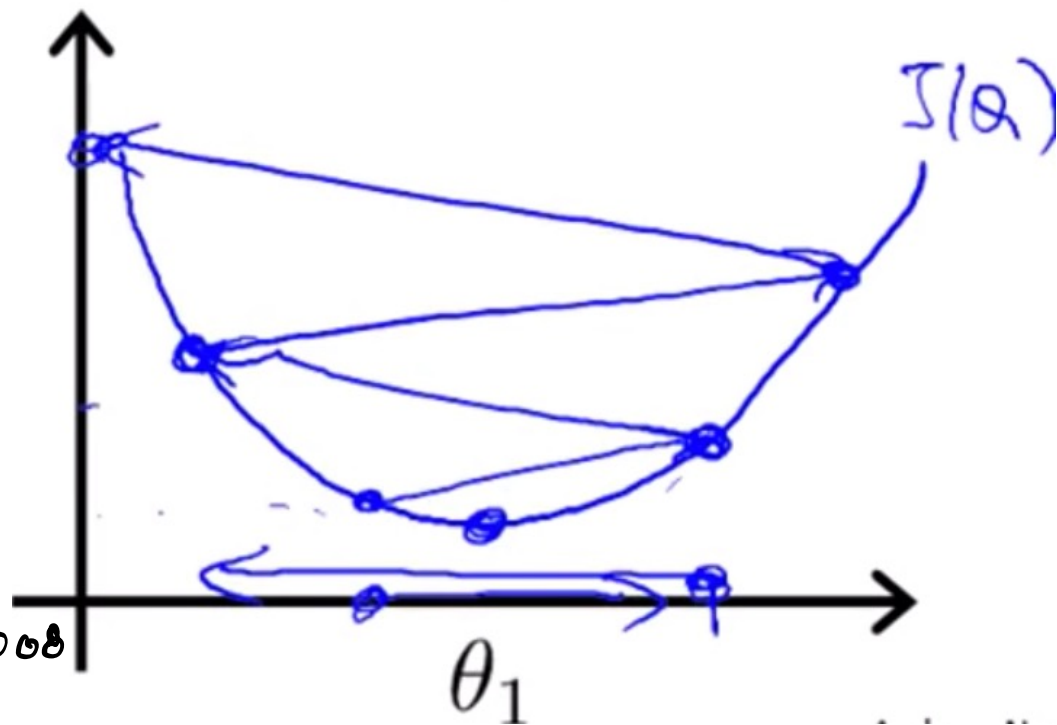


If alpha too large

$\alpha$  : hyperparameter

$0,1 \times 10^{-1}$   $0,01 \times 10^{-1}$   $0,001$

$0,002$   $0,0008$



Scheduler : linear  $0,9$   
corine

GD - momentum  
adam

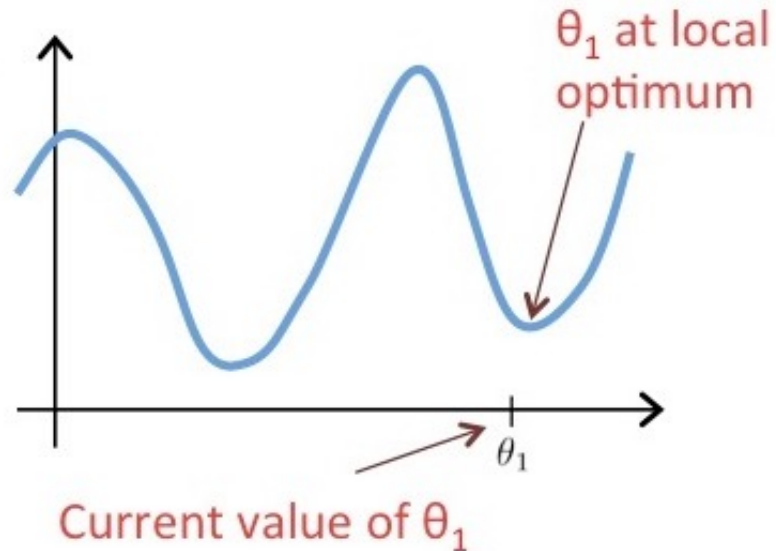
global  
vs

## Question

Suppose  $\theta_1$  is at a local optimum of  $J(\theta_1)$ , such as shown in the figure.

What will one step of gradient descent  $\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$  do?

1. Change  $\theta_1$  in a random direction?
2. Move  $\theta_1$  in the direction of the global minimum of  $J(\theta_1)$ ?
3. Leave  $\theta_1$  unchanged?
4. Decrease  $\theta_1$ ?



# Recap

- To update our parameter with the Gradient Descent algorithm, we perform 2 essential steps :
  1. Compute the derivative of the parameter with respect to the value we want to minimize (ie. our cost: a score to express how good our model is doing)
  2. Take an optimization step/update the parameter. This update will be proportional to the derivative and the learning rate.  
Large derivative (steep tangent line) + large learning rate = big update



# Piecing everything together

- This is all we need :
  - A **hypothesis** function (our model)
  - A **cost function** (to tell us how well/bad our model is doing)
  - **Gradient Descent or variant** (to update our parameters and get closer to a better model)

## 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  $j = 0$ )  
}

$|\theta_{j_t} - \theta_{j_{t+1}}| \leq 0,0001$  then stop.

## Linear Regression Model

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

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

# Derivatives vs. Partial derivatives

- Except, instead of having a cost function with a single input, we are back to 2 inputs, our 2 parameters  $\theta_0$  and  $\theta_1$ .
- When we have functions with multiple inputs (known as multivariate functions), computing 1 single derivative is no longer enough!
- The function's « **instantaneous rate of change** » for a given combination of parameters is now determined by 2 values :
  - How does a tiny change in  $\theta_0$  change  $J(\theta_0, \theta_1)$  ?
  - How does a tiny change in  $\theta_1$  change  $J(\theta_0, \theta_1)$  ?

=> Packed together into a vector, these 2 derivatives make up what is referred to as the **gradient**
- Each derivative is a **partial derivative**. (you need both together to get the whole picture !)

# Derivatives vs. Partial derivatives

- *Partial Derivative* :

This comes down to calculating the derivative at each input value, treating the other input as a constant

- We pretend for a second that the other input value has basically no effect on the function
  - when looking at  $\theta_0$  , we treat  $\theta_1$  as a constant
  - when looking at  $\theta_1$  , we treat  $\theta_0$  as a constant

# Partial derivatives visually

- To help illustrate things and relate them to our simple Gradient Descent intuition:



# Gradient Descent

- Each partial derivative tells us how the function behaves (increases/decreases, quickly/slowly, stays constant...) with respect to a single input
- We can then use this information to know if we should increase or decrease each input to get closer to our minimum cost value !
- **Gradient** : the partial derivatives packed together in a vector
- **Descent** : we want to find the cost function's minimum, using the gradient as a source of information to tell us if the cost is increasing/decreasing with respect to each input.

# Update rule

- So we need to figure out the partial derivatives for each parameter !
- the partial derivative of  $J(\theta_0, \theta_1)$  with respect to  $\theta_1$
- the partial derivative of  $J(\theta_0, \theta_1)$  with respect to  $\theta_2$

# Partial derivatives of $J(\theta_1, \theta_2)$

- You can treat these results as being **given**, in order not to go into the details of the derivation.

- General formula

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

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

See here for the [MSE derivative](#)

# Partial derivatives of $J(\theta_1, \theta_2)$

- Here are the partial derivatives obtained (take these at face value for now):

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

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

- These formulas allow us to compute the partial derivatives for each of the parameters, which we can then plug into our Gradient Descent algorithm.



# Gradient Descent

- We now have formulas to update our parameters !

repeat until convergence {

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

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

}

# Quick recap to put things into perspective

- We have :
- a **model**, which is a line :

$$h(x) = \theta_0 + \theta_1 x$$

- a **cost function**, to tell us how good/bad our model fits the data:

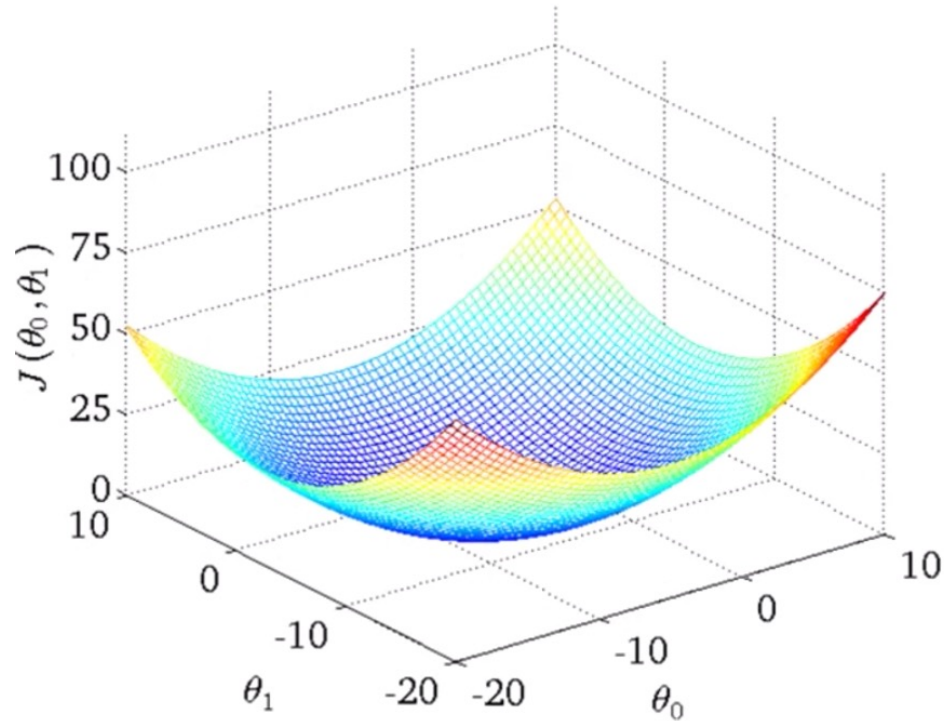
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

- **Gradient Descent**, a method to update our parameters so as to minimize the cost function:

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

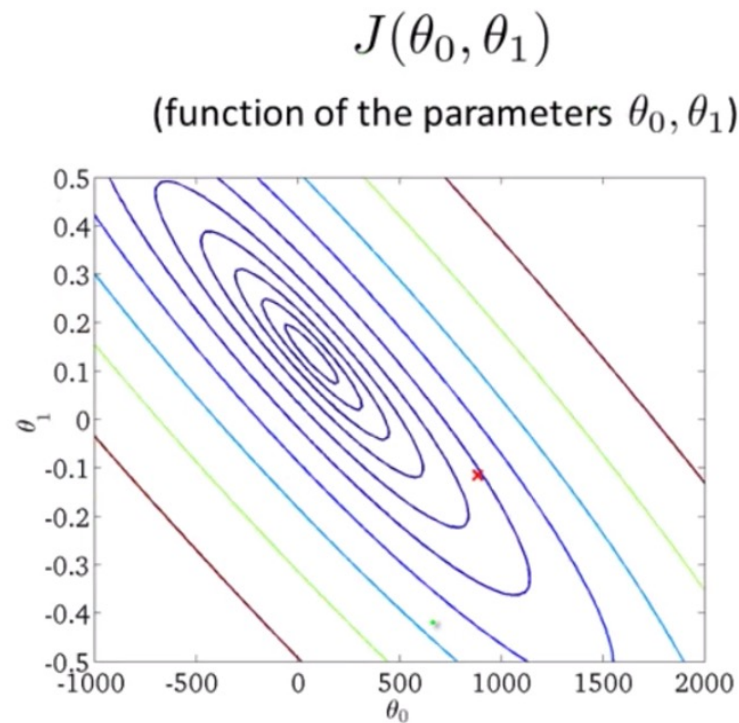
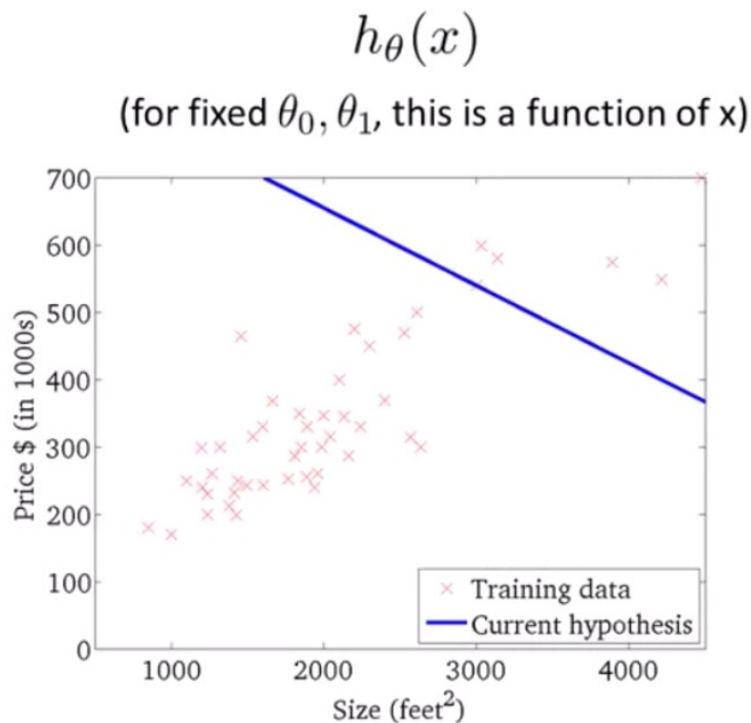
# Update examples

- For linear regression, the cost function will always be bowl-shaped



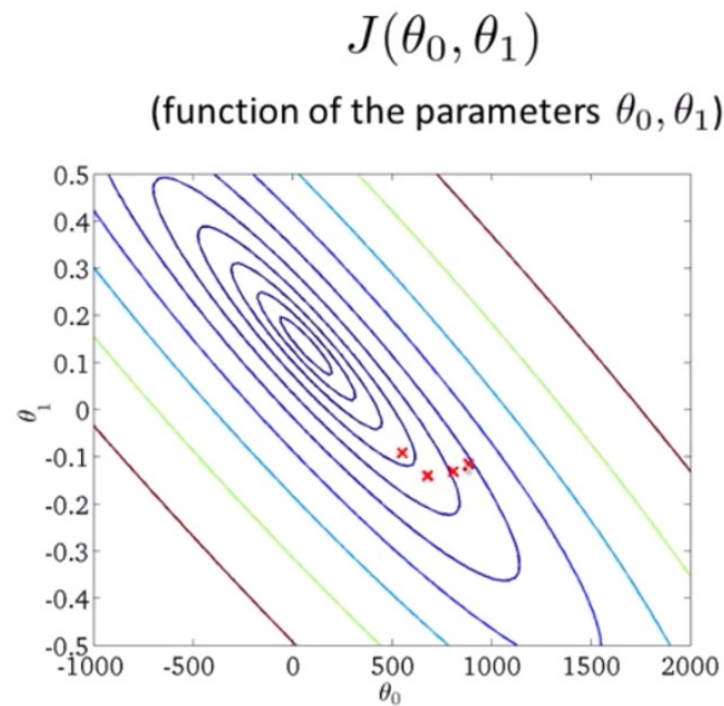
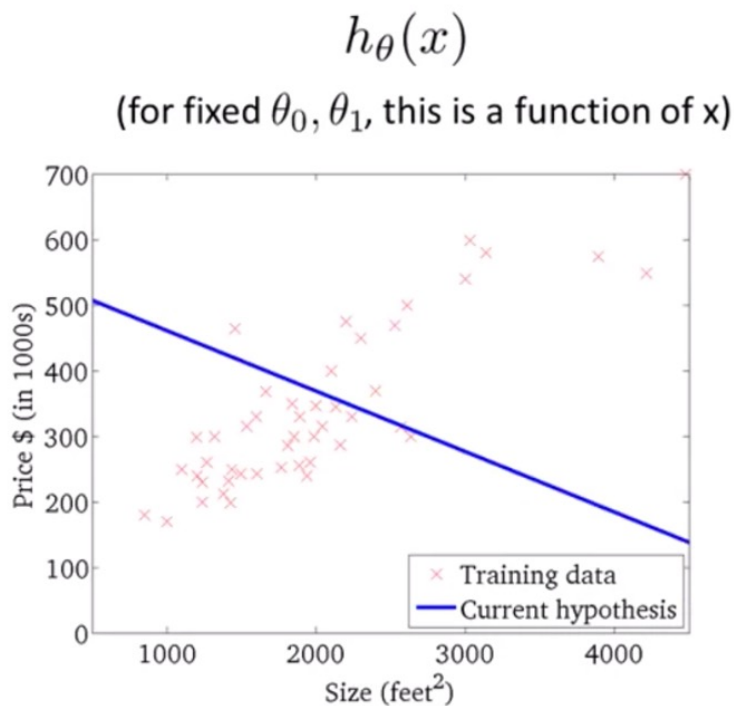
# Update examples

Say we initialize our parameters randomly, this is the model and cost :



# Update examples

As we take Gradient Descent steps, the model (line) seems to be fitting the data better

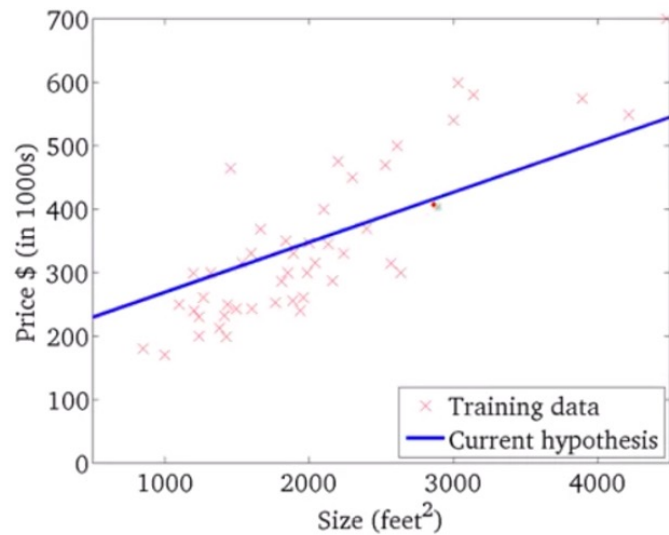


# Update examples

Until we reach the global minimum

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

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

(function of the parameters  $\theta_0, \theta_1$ )

