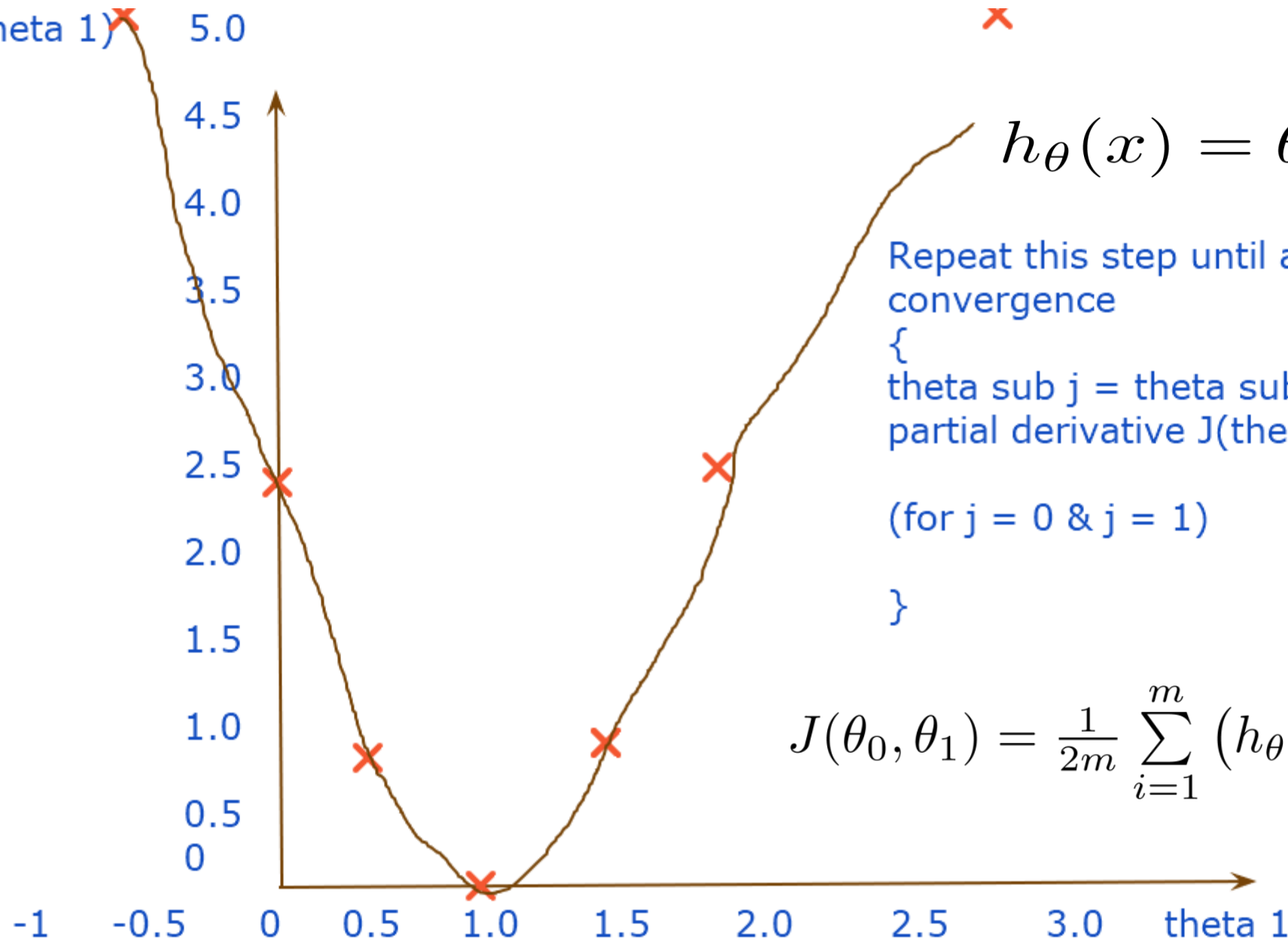J (Theta 1)

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

Repeat this step until achieve convergence
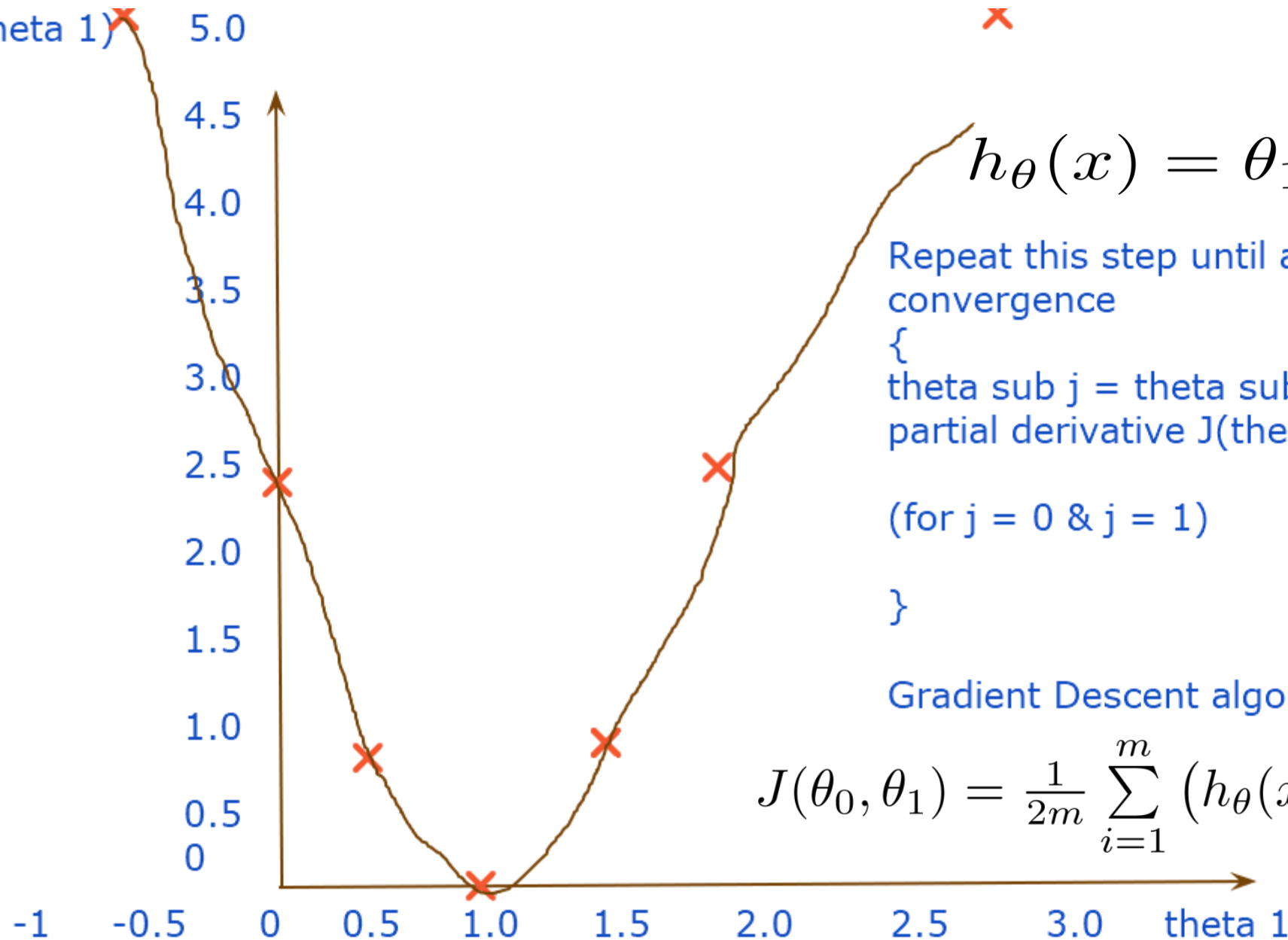{
theta sub j = theta sub j - alpha partial derivative J(theta's)/j

(for j = 0 & j = 1)

}

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

theta 1

J (Theta 1)

5.0
4.5
4.0
3.5
3.0
2.5
2.0
1.5
1.0
0.5
0

-1    -0.5    0    0.5    1.0    1.5    2.0    2.5    3.0    theta 1

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

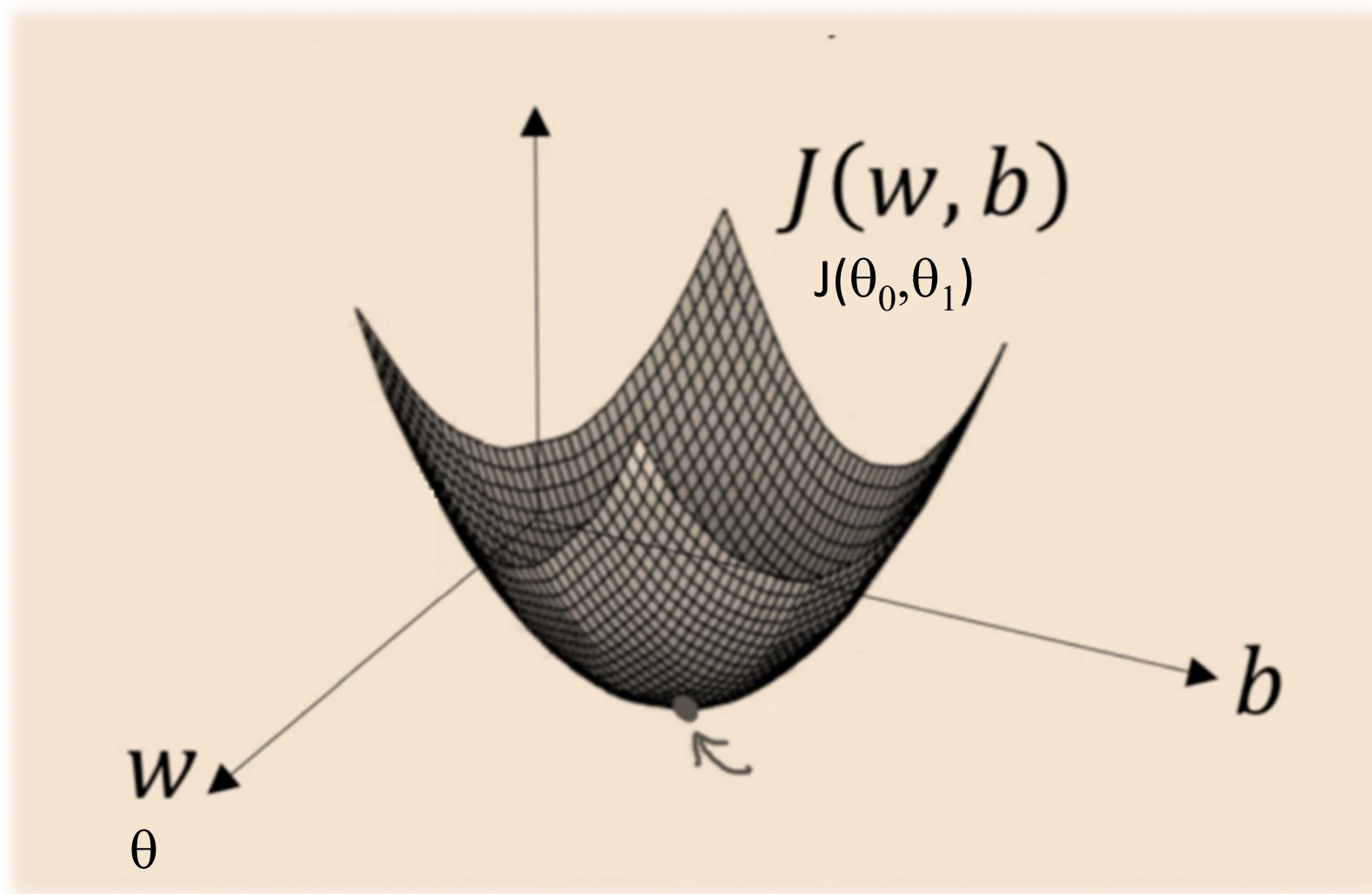Repeat this step until achieve convergence
{
theta sub j = theta sub j - alpha partial derivative J(theta's)/j
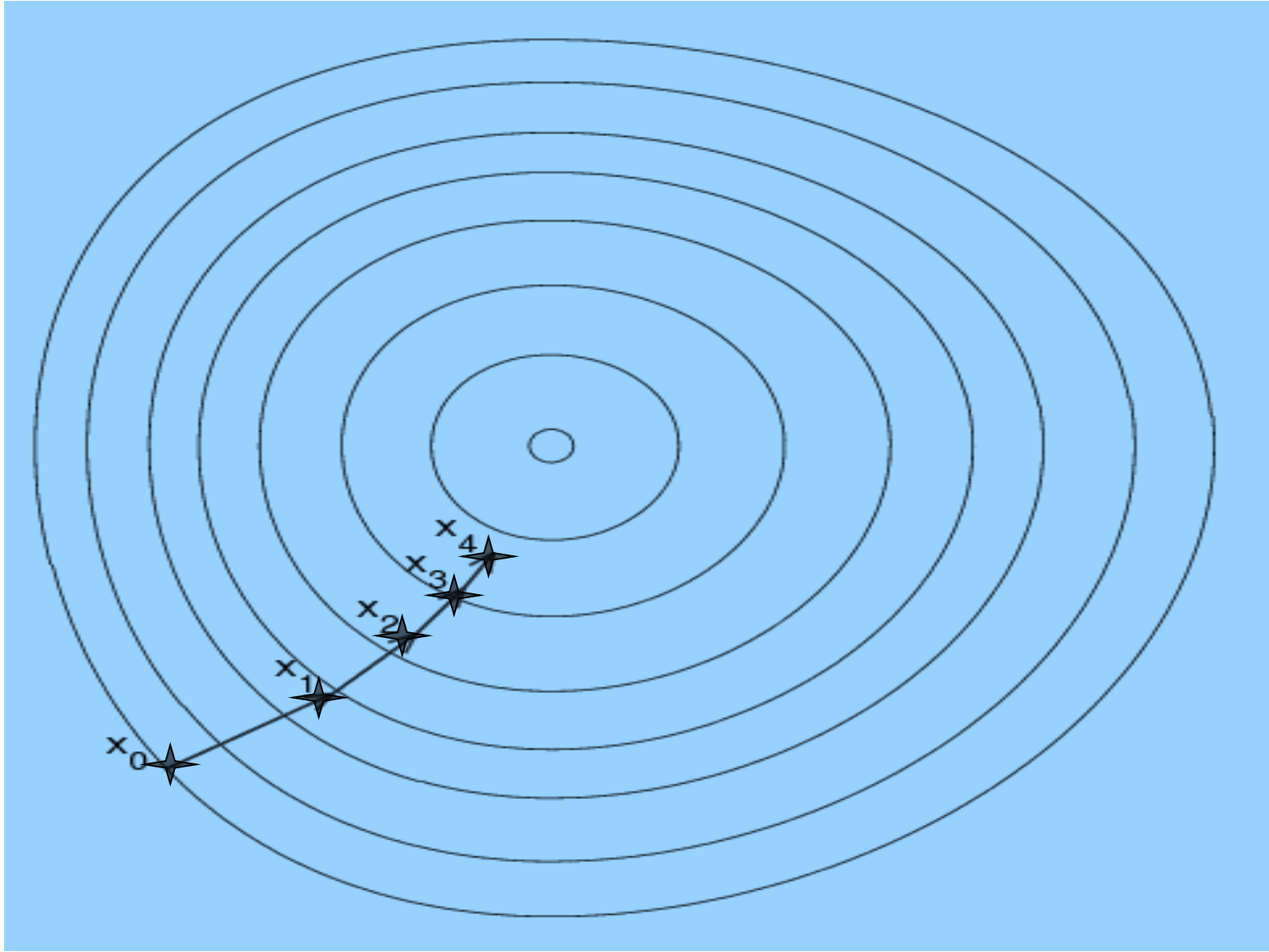
(for j = 0 & j = 1)

}

Gradient Descent algorithm

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

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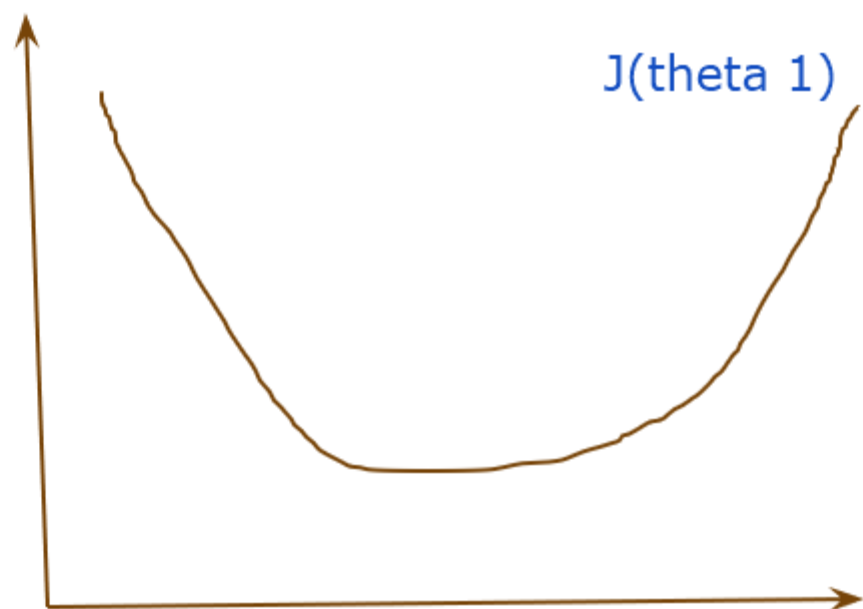
$J(\theta_0, \theta_1)$

$b$

$w$

$\theta$

# Gradient Descent algorithm
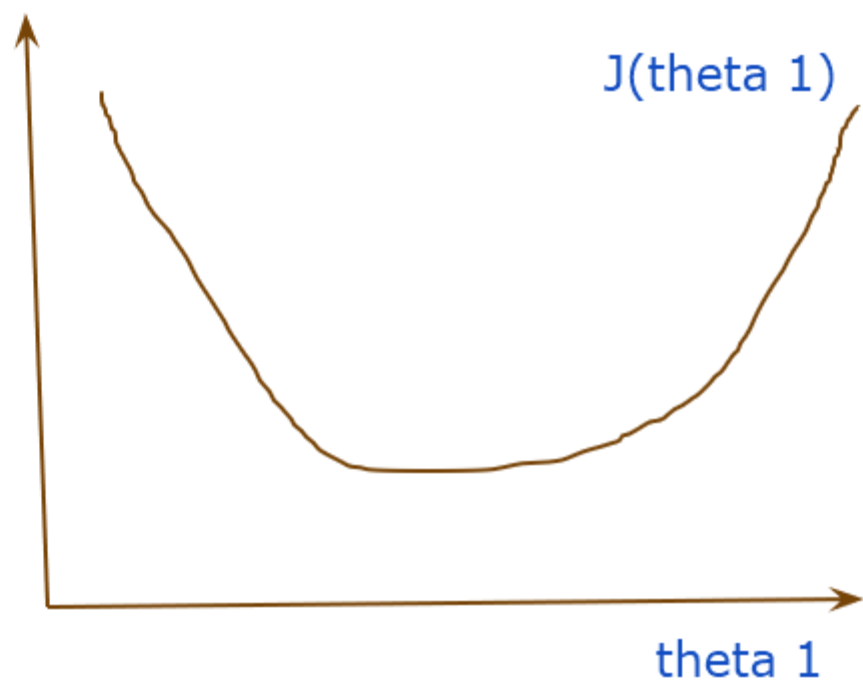
# Gradient Descent algorithm

Gradient Descent algorithm

J(theta 1)

# Gradient Descent algorithm



J(theta 1)

theta 1

Gradient Descent algorithm

$J(\theta_1)$

$J(\theta_1)$

theta 1

Gradient Descent algorithm

$J(\theta_1)$

$J(\theta_1)$

theta 1

# Gradient Descent algorithm



J (theta 1)

J(theta 1)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

Gradient Descent algorithm

J (theta 1)

J(theta 1)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

Gradient Descent algorithm

J(theta 1)    alpha: learning rate

J (theta 1)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

Gradient Descent algorithm

J (theta 1)

J(theta 1)

alpha: learning rate

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values

Gradient Descent algorithm

J(theta 1)

J (theta 1)

alpha: learning rate

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values

# Gradient Descent algorithm



J(theta 1)          alpha: learning rate

J (theta 1)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values

Gradient Descent algorithm

J (theta 1)

J(theta 1)     alpha: learning rate

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values

Gradient Descent algorithm

J(theta 1)

alpha: learning rate

J (theta 1)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values

Gradient Descent algorithm

J (theta 1)

J(theta 1)    alpha: learning rate

theta 1

theta 1:= theta 1 - alpha (positive slope)

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
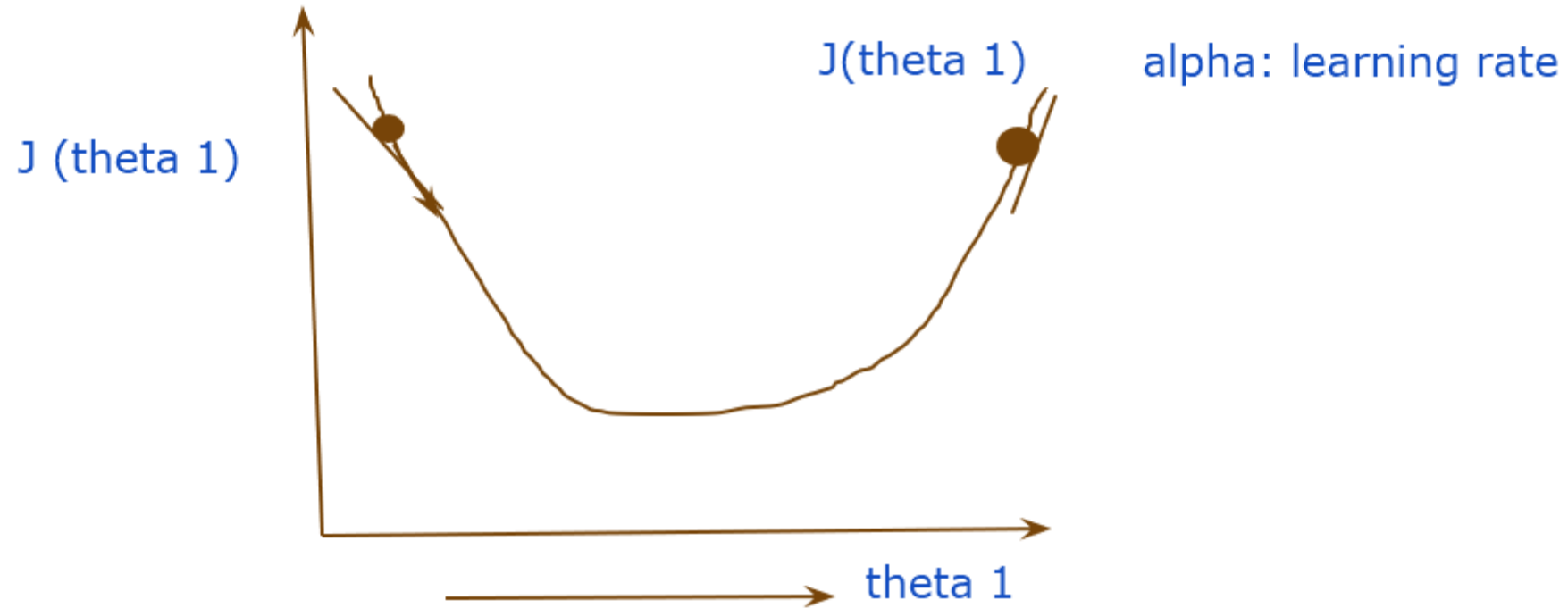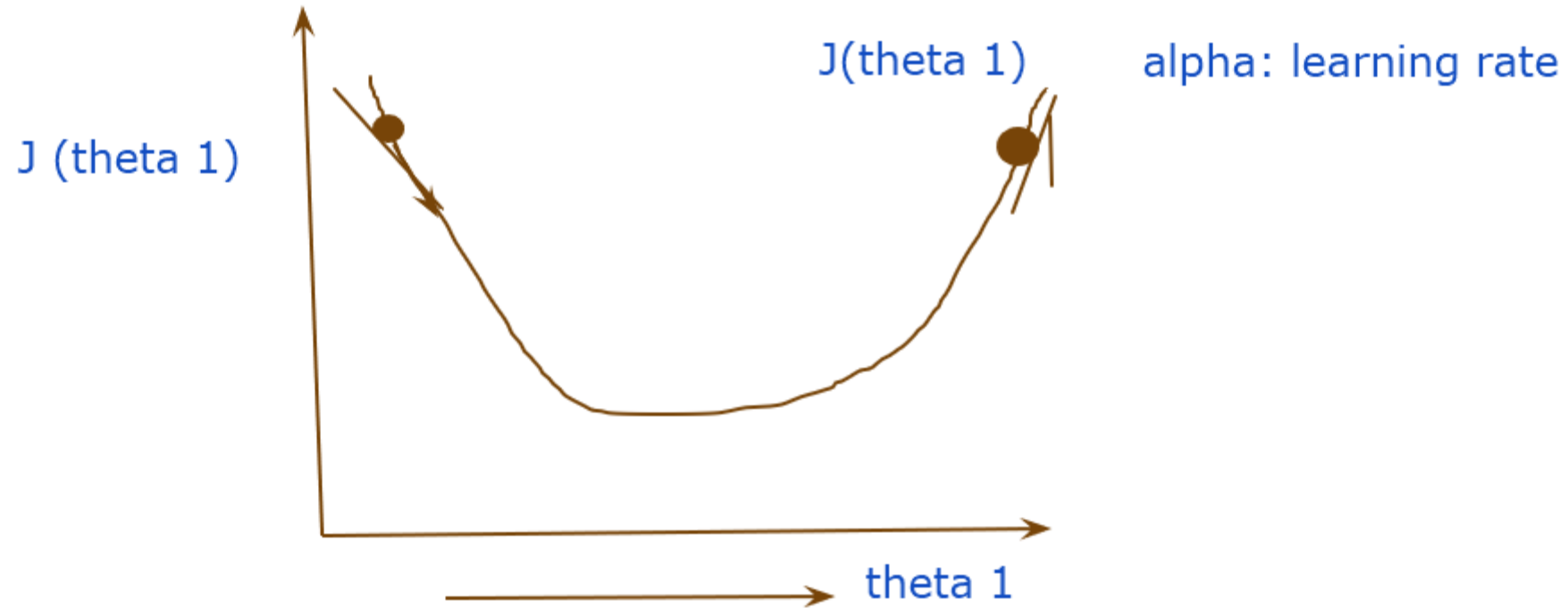
# Gradient Descent algorithm



J(theta 1)

alpha: learning rate

J (theta 1)

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
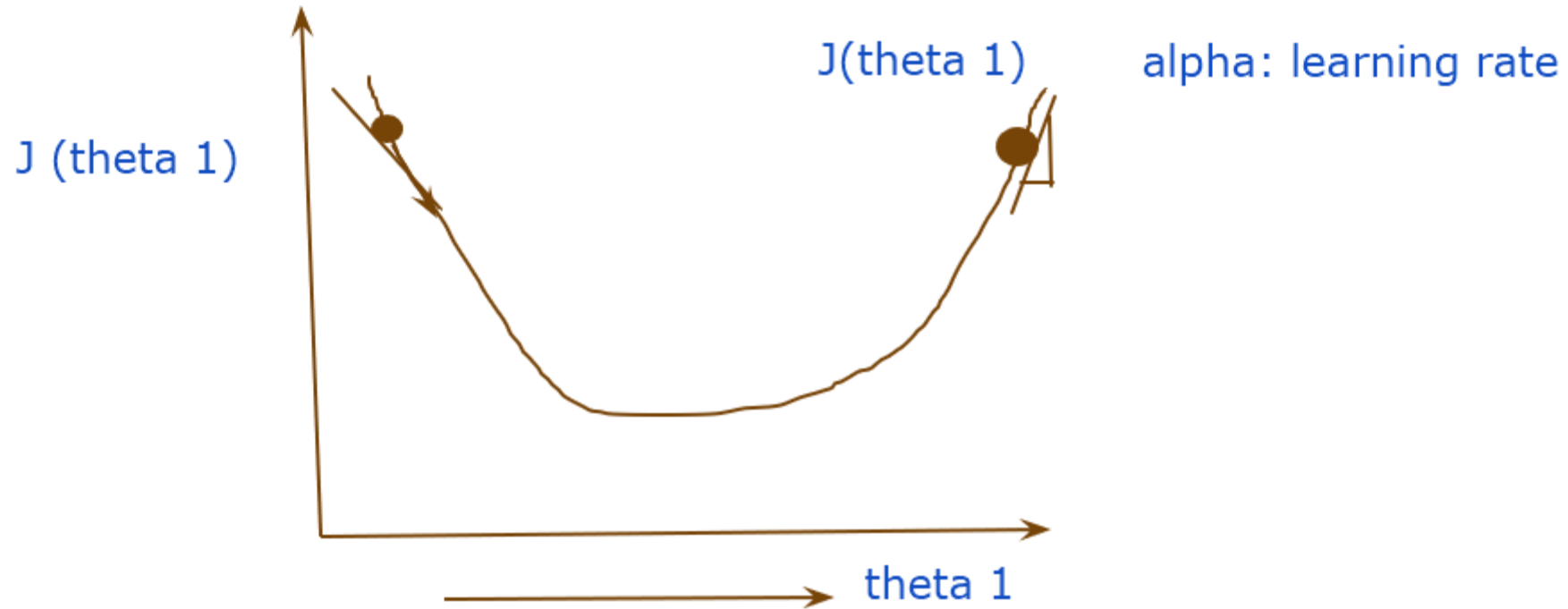
Gradient Descent algorithm



J (theta 1)

J(theta 1)        alpha: learning rate

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values

Gradient Descent algorithm

J (theta 1)

J(theta 1)     alpha: learning rate

alpha: very small

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
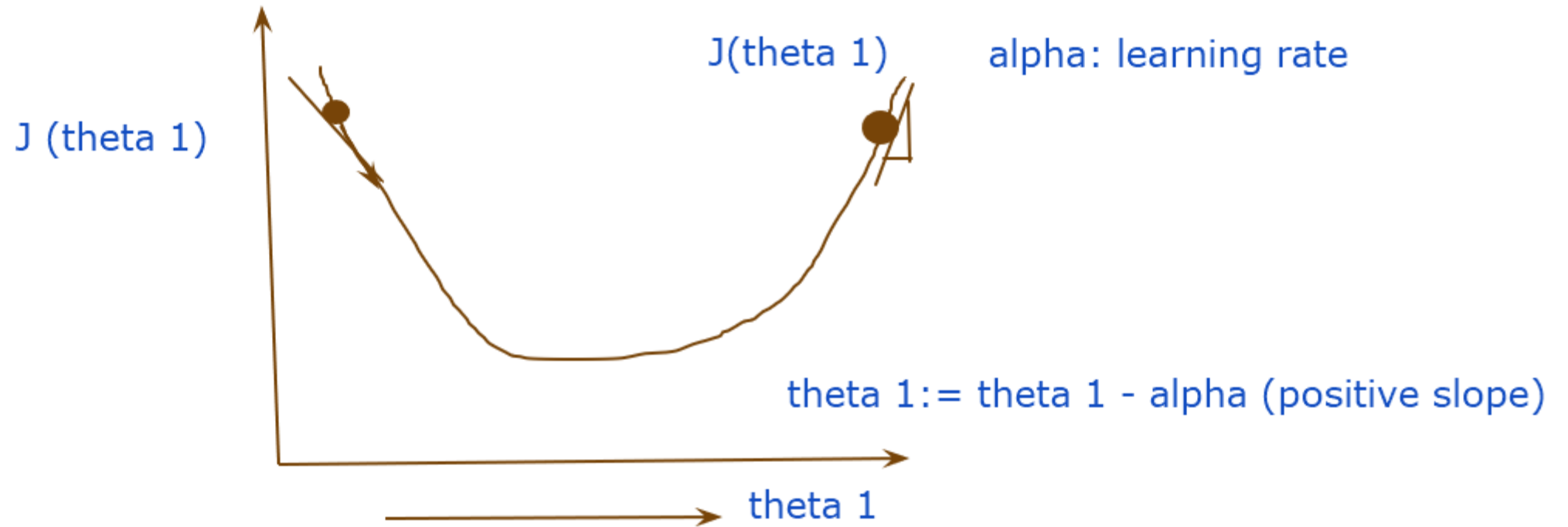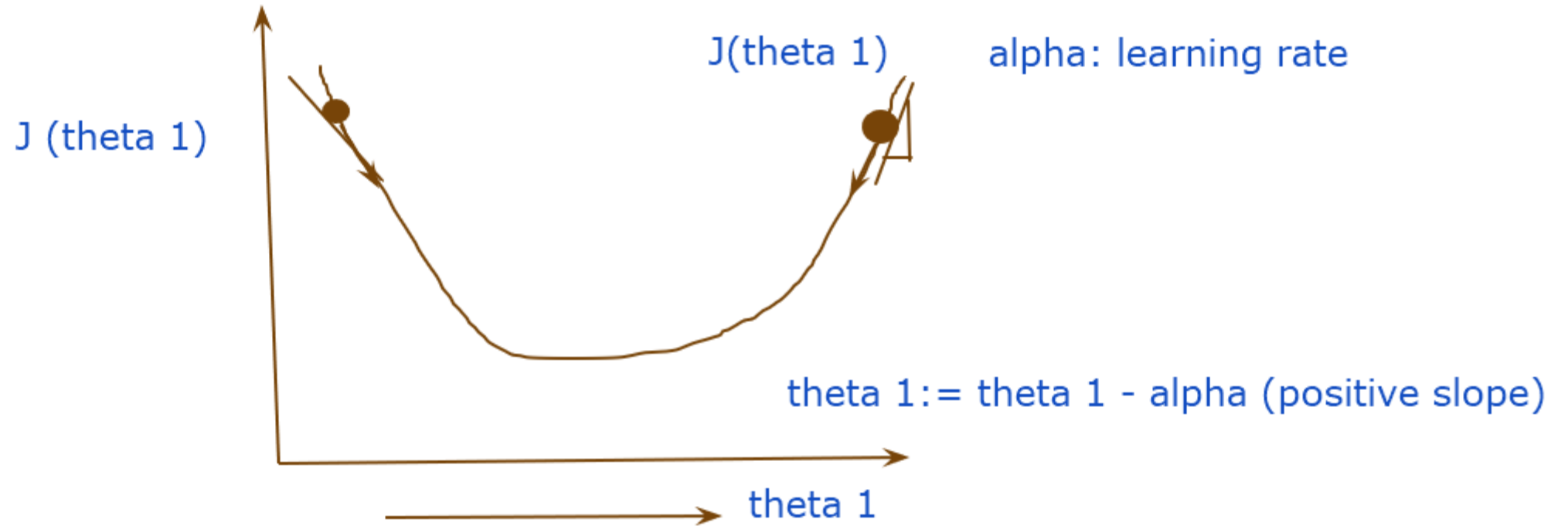
Gradient Descent algorithm

J (theta 1)

J(theta 1)

alpha: learning rate

alpha: very small

it will be very slow

theta 1 := theta 1 - alpha (positive slope)

theta 1

theta 1 := theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1 := theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
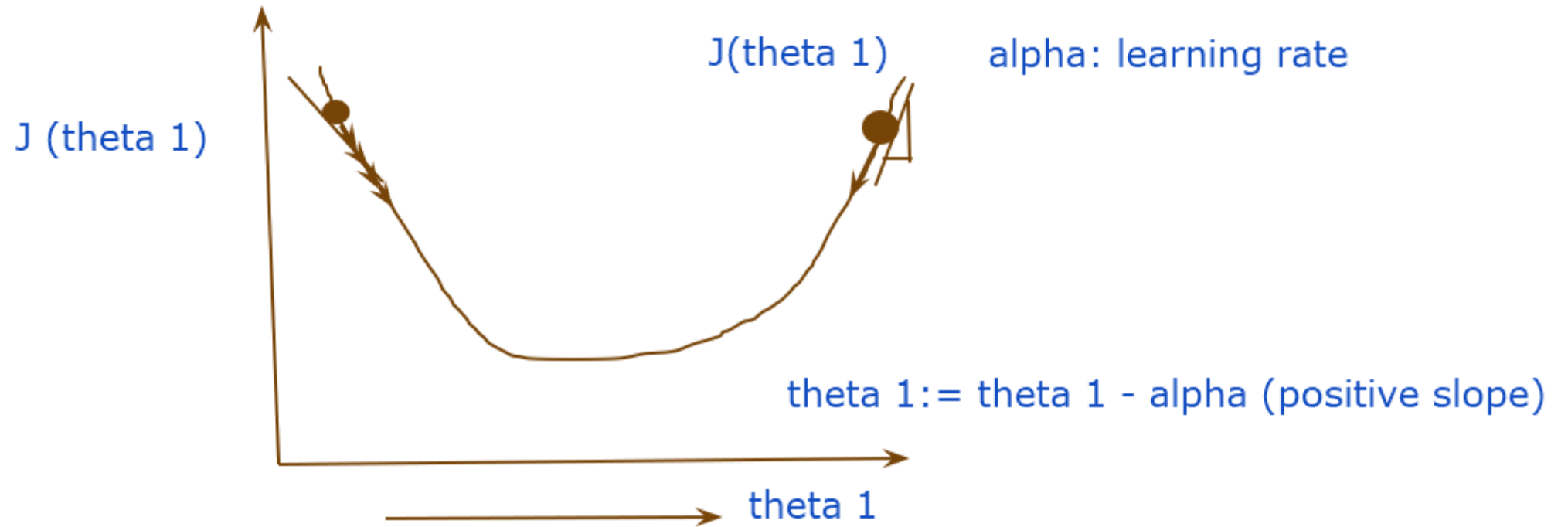
Gradient Descent algorithm

J(theta 1)    alpha: learning rate

J (theta 1)    alpha: very small

it will be very slow
alpha: very high

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
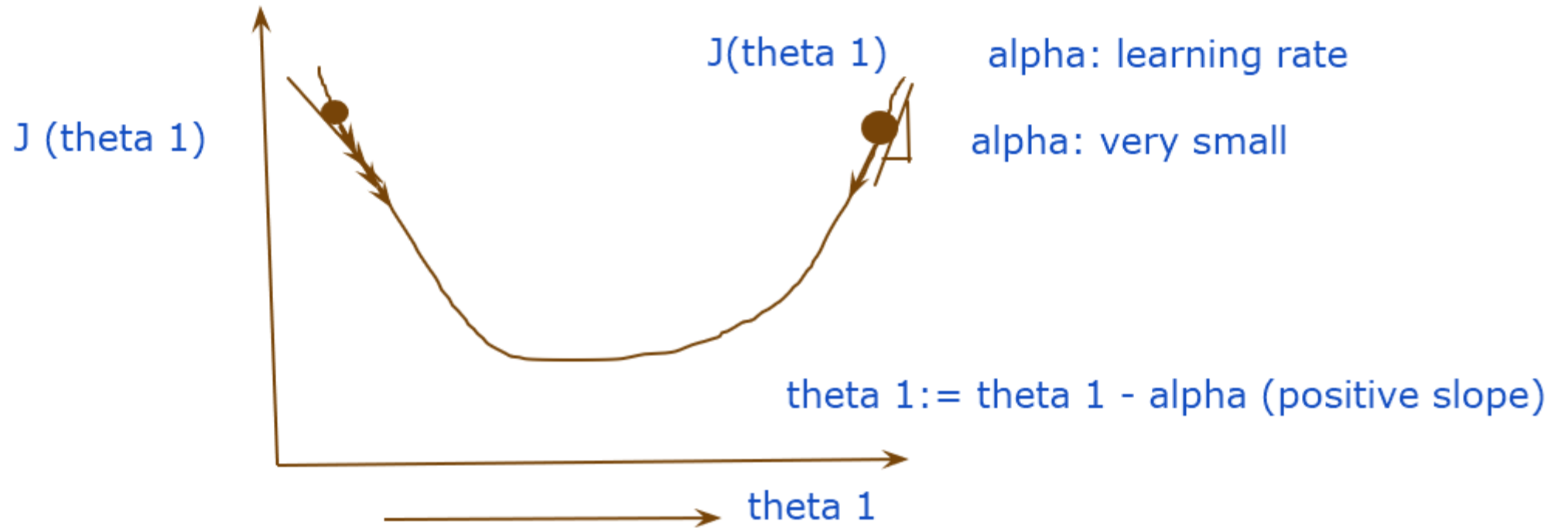
Gradient Descent algorithm

J (theta 1)

J(theta 1)

alpha: learning rate

alpha: very small

it will be very slow
alpha: very high

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
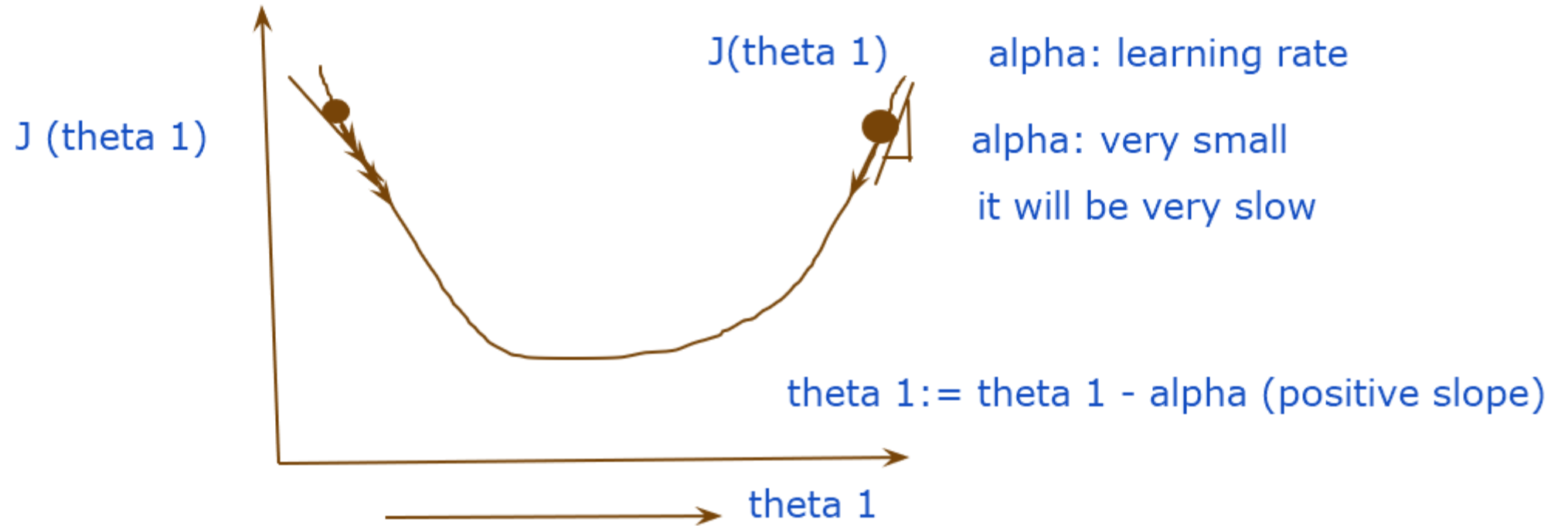
Gradient Descent algorithm

J (theta 1)

J(theta 1)     alpha: learning rate

alpha: very small

it will be very slow
alpha: very high

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input
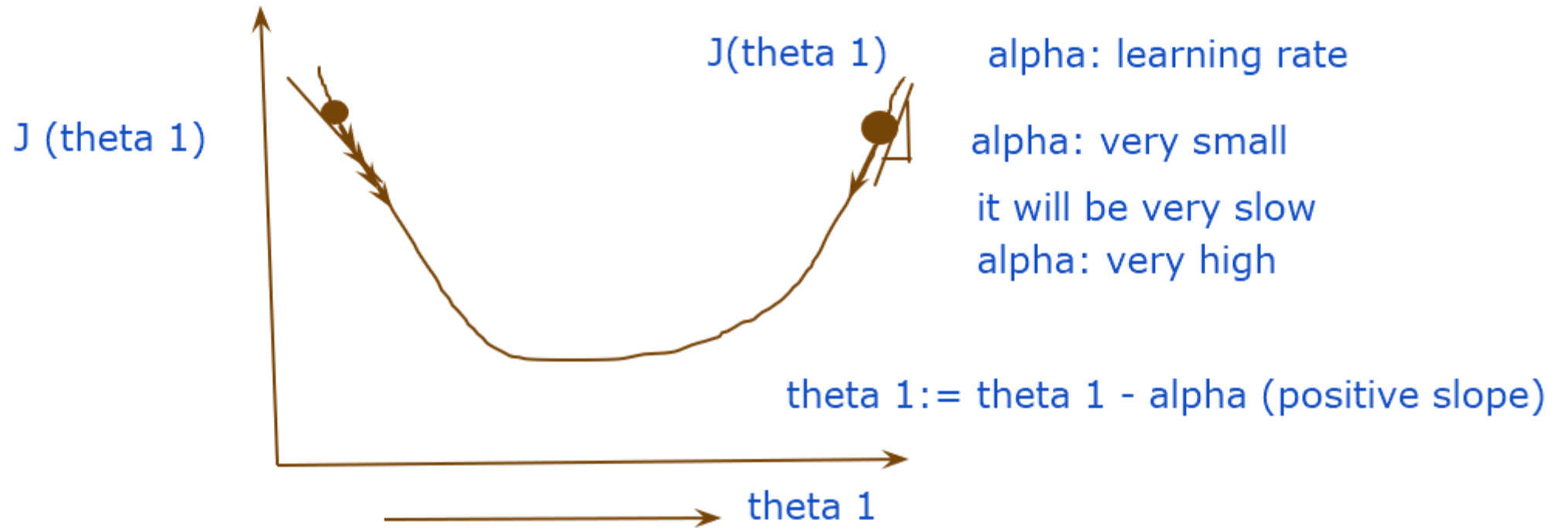with small values

# Gradient Descent algorithm



J(theta 1)

J (theta 1)

alpha: learning rate

alpha: very small

it will be very slow
alpha: very high

theta 1

theta 1:= theta 1 - alpha (positive slope)

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
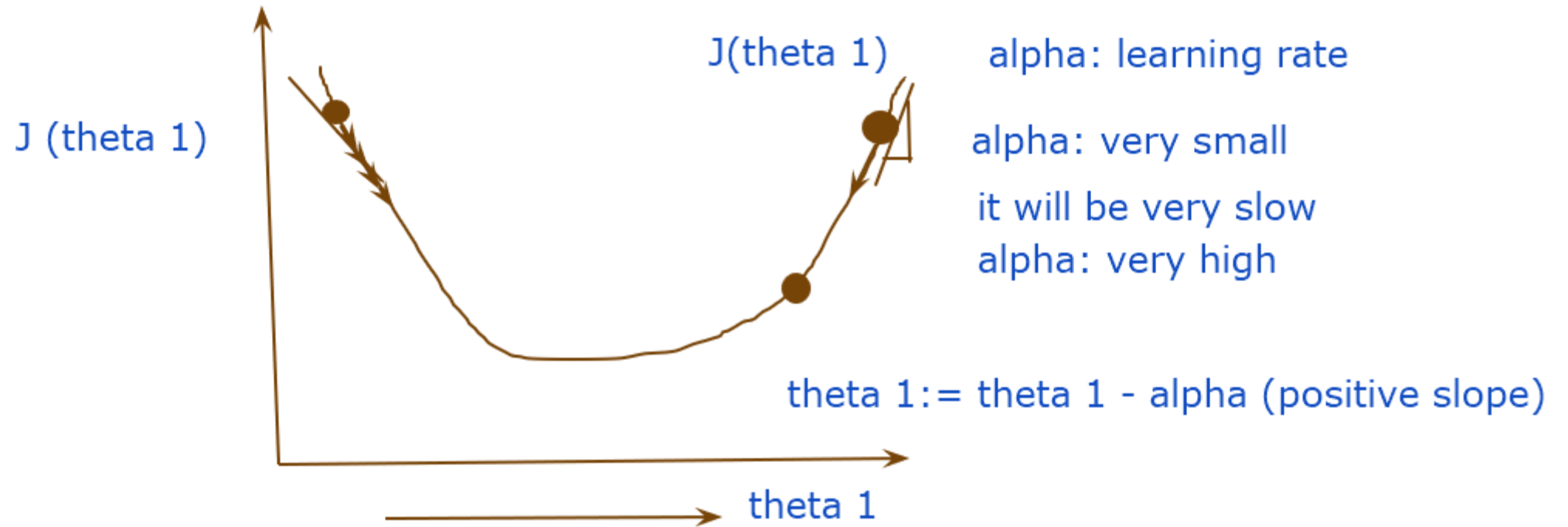
Gradient Descent algorithm

J(theta 1)        alpha: learning rate

J (theta 1)        alpha: very small

it will be very slow
alpha: very high

it may fail to converge or even diverge

theta 1:= theta 1 - alpha (positive slope)

theta 1

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
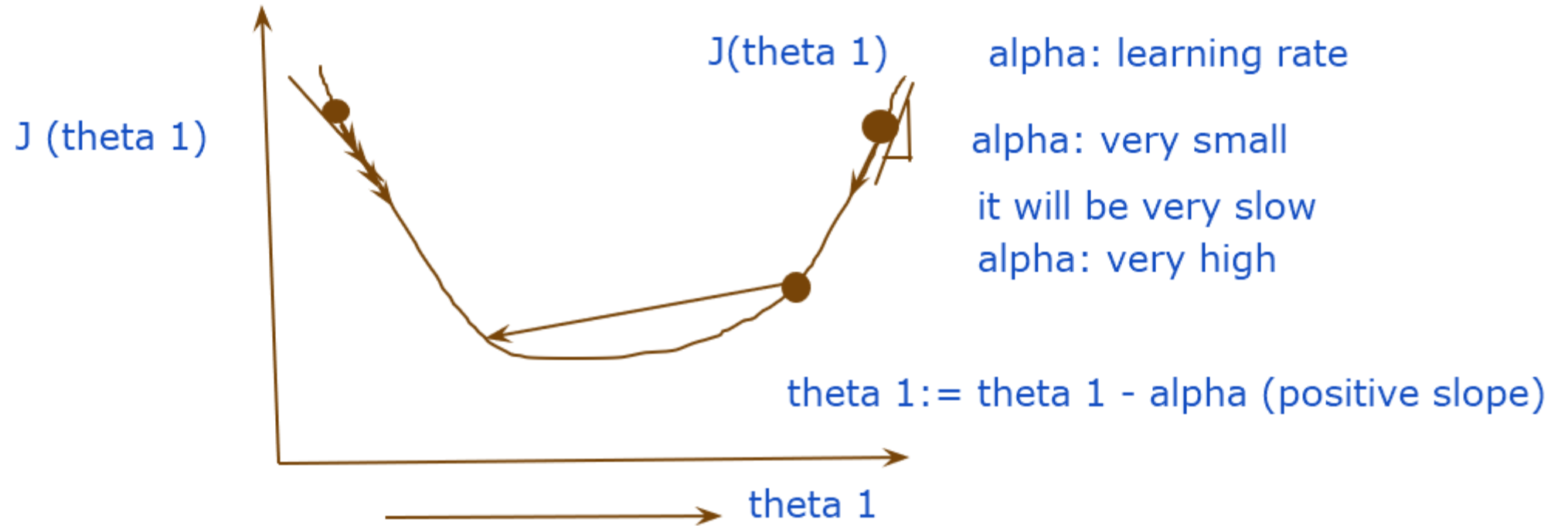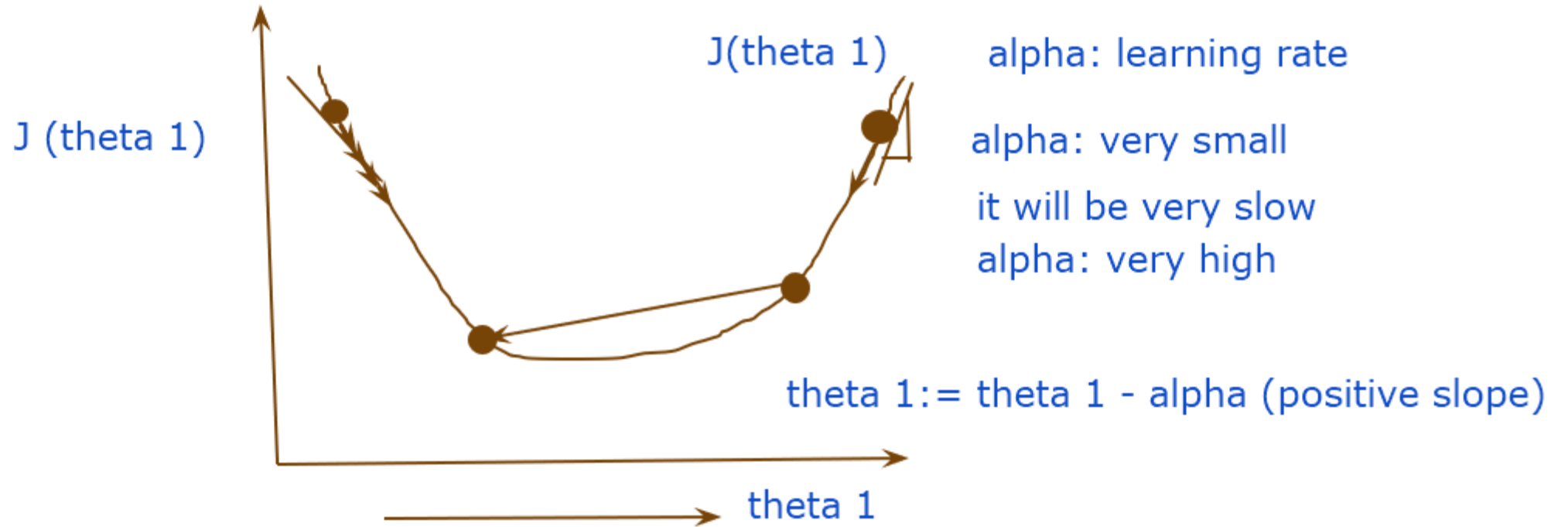
Gradient Descent algorithm

J (theta 1)

J(theta 1)

theta 1

alpha: learning rate

alpha: very small

it will be very slow
alpha: very high

it may fail to converge or even diverge

theta 1:= theta 1 - alpha (positive slope)

theta 1:= theta 1 - alpha(partial derivative J(theta 1)/theta 1 )

theta 1:= theta 1 - alpha (negative slope)

gradient: measures how much the output of a function changes if we change the input with small values
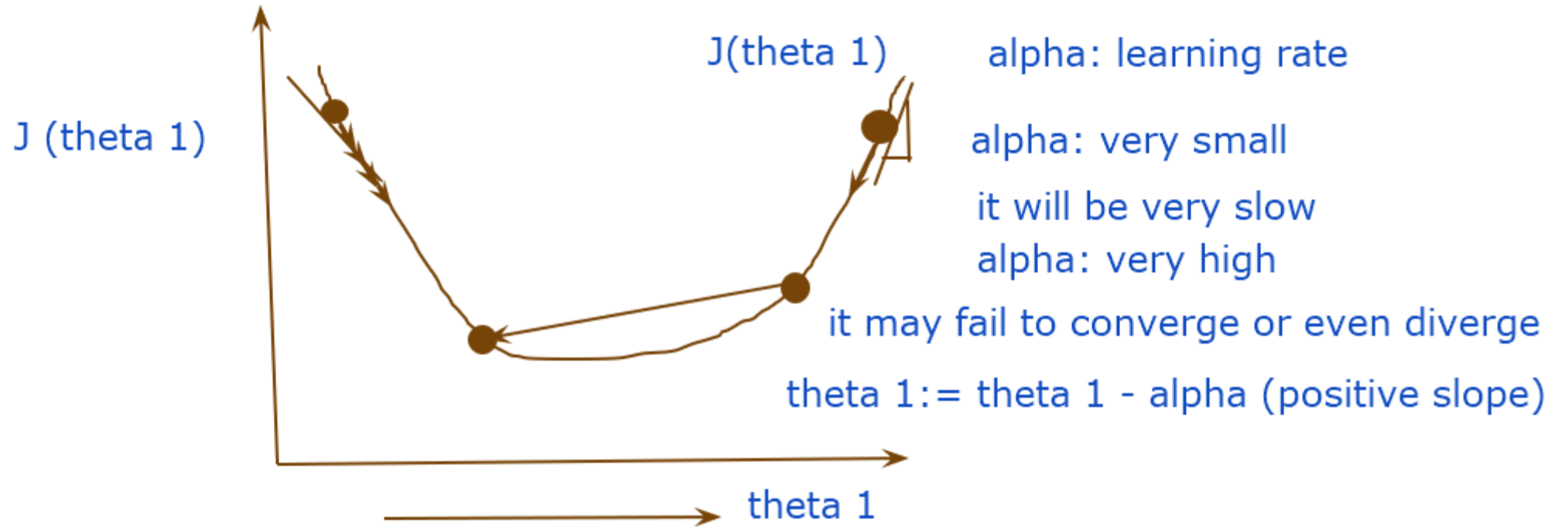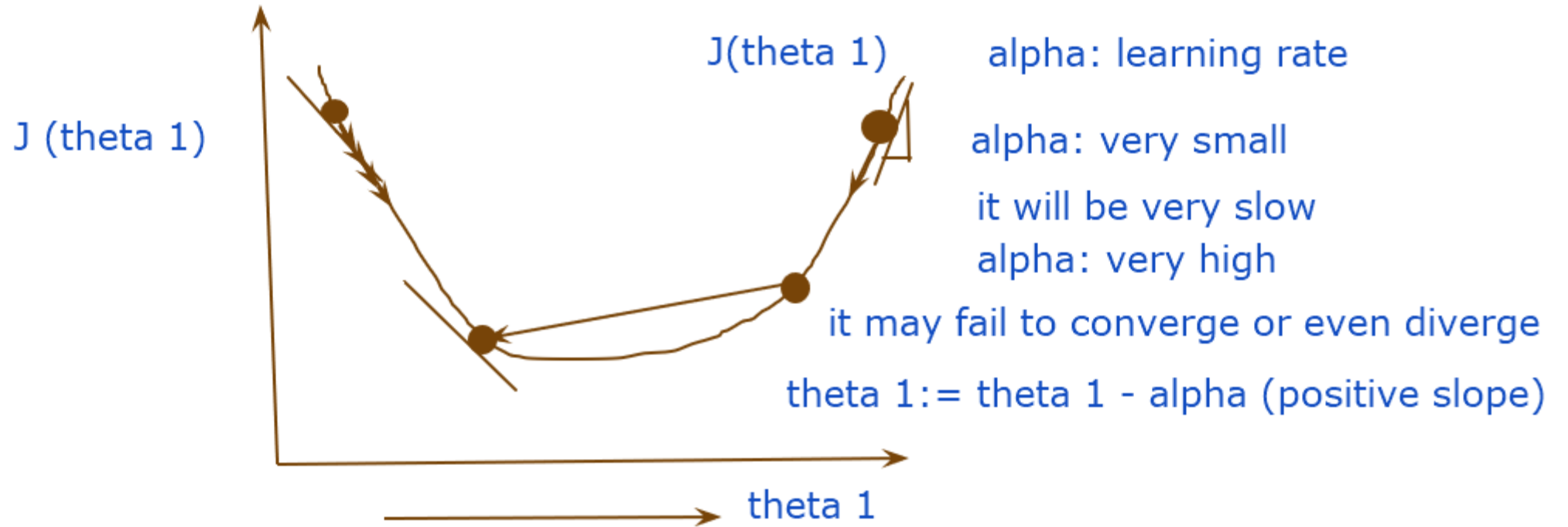
# Gradient Descent Algorithm

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

Want $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ This is true for all
For i = 0, 1, 2, 3, ...., n $\theta_i$

## Steps:

- Start with some $\theta_0, \theta_1$

- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

  until we hopefully end up at a minimum

# Gradient descent algorithm

repeat until convergence {

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

}

**Learning Rate**

---

Correct Update:

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

update Simultaneously: $\theta_0$ and $\theta_1$

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

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

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

"Batch" gradient Descent algorithm:
"Stochastic" gradient descent algorithm

"Batch" gradient Descent algorithm:
"Stochastic" gradient descent algorithm

Each step of gradient descent uses all the training example: batch gradient descent algorithm
Each step of gradient decent uses single training example
Mini-batch: each step of gradient descent uses subset of training example

"Batch" gradient Descent algorithm:
"Stochastic" gradient descent algorithm

Each step of gradient descent uses all the training example: batch gradient descent algorithm
Each step of gradient decent uses single training example
Mini-batch: each step of gradient descent uses subset of training example

# "Batch" Gradient Descent

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

# "Stochastic" Gradient Descent

"Stochastic": Each step of gradient descent use only single training example.

# "Mini-batch" Gradient Descent

It lies in between of these two extremes, and can use a mini-batch (small portion) of training data examples in each step

# Batch Gradient Descent Algorithm

# Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]$$

Batch Gradient Descent Algorithm

$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Batch Gradient Descent Algorithm

$J(\theta_1) := \frac{1}{2m} \text{summation } (i = 1 \text{ to } m) [h(x)i - (y)i]^2$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Batch Gradient Descent Algorithm

$J(\text{theta } 1){:} = 1/2m$ summation $(i = 1 \text{ to } m) [h(x)i - (y)i]^2$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]^2$$

$\theta_1 : \theta_1 - (\alpha)$ (partial derivative of cost function with respect to the parameters)

Training data points

Batch Gradient Descent Algorithm

$$J(\theta 1) := \frac{1}{2m} \text{summation } (i = 1 \text{ to } m) [h(x)i - (y)i]^2$$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]^2$$

$\theta_1: \theta_1 - (alpha)$ (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \text{summation } (i = 1 \text{ to } m) [h(x)i - (y)i]^2$$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

update the parameters for each training example one by one

Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]^2$$

$\theta_1 : \theta_1 - (alpha)$ (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

update the parameters for each training example one by one

Mini-batch gradient descent algorithm

Batch Gradient Descent Algorithm

$$J(\theta 1) := \frac{1}{2m} \text{summation } (i = 1 \text{ to } m) [h(x)i - (y)i]^2$$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

update the parameters for each training example one by one

Mini-batch gradient descent algorithm : splits the training dataset into small batches & perform an update

Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]^2$$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

update the parameters for each training example one by one

Mini-batch gradient descent algorithm : splits the training dataset into small batches & perform an update
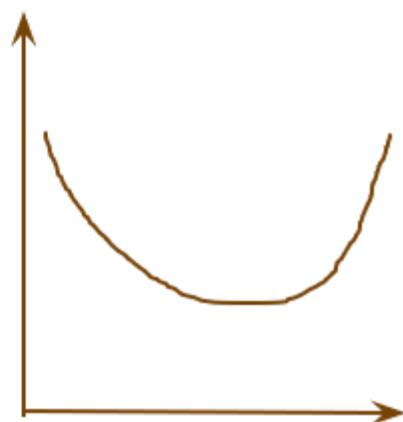
Batch Gradient Descent Algorithm

$$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]^2$$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm
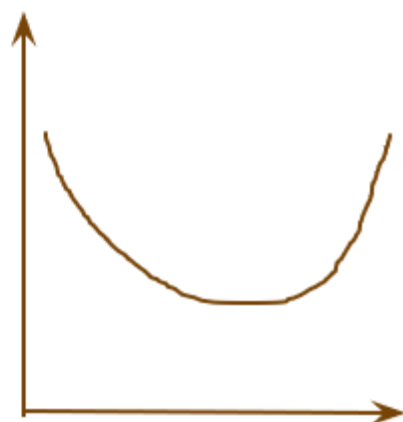
update the parameters for each training example one by one

Mini-batch gradient descent algorithm : splits the training dataset into small batches & perform an update

Batch Gradient Descent Algorithm

$J(\theta_1) := \frac{1}{2m} \sum_{i=1}^{m} [h(x)_i - (y)_i]^2$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

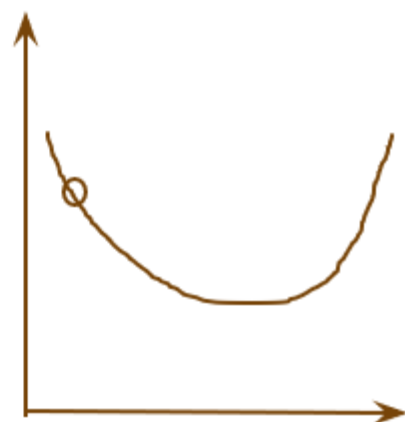update the parameters for each training example one by one

Mini-batch gradient descent algorithm : splits the training dataset into small batches & perform an update

Batch Gradient Descent Algorithm

$J(\theta 1) := \frac{1}{2m}$ summation $(i = 1$ to $m) [h(x)i - (y)i]^2$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

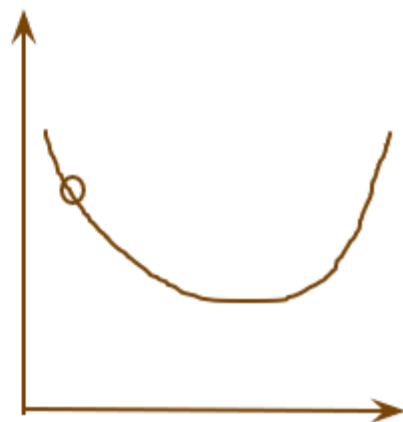update the parameters for each training example one by one

Mini-batch gradient descent algorithm : splits the training dataset into small
batches & perform an update

## Batch Gradient Descent Algorithm

$J(\theta_1) := \frac{1}{2m} \text{summation } (i = 1 \text{ to } m) [h(x)i - (y)i]^2$

theta 1: theta 1 - (aplha) (partial derivative of cost function with respect to the parameters)

Training data points

Vanilla gradient descent algorithm/gradient descent algorithm

Stochastic Gradient descent algorithm

update the parameters for each training example one by one

Mini-batch gradient descent algorithm : splits the training dataset into small
batches & perform an update