

# ML – Gradient Descent

draft

# Gradient descent

- An algorithm used to find the minimum value of cost.
  - It works for any function, not only just the cost function of linear regression that we discussed so far.
- Used in almost all ML modes e.g., regression, NN, Deep learning.

# What we want to do with cost?

Have some function  $J(w, b)$  *for linear regression  
or any function*

Want  $\min_{w, b} J(w, b)$   $\min_{w_1, \dots, w_n, b} \underline{J(w_1, w_2, \dots, w_n, b)}$

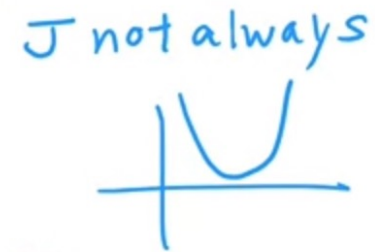
Outline:

Start with some  $w, b$  (set  $w=0, b=0$ )

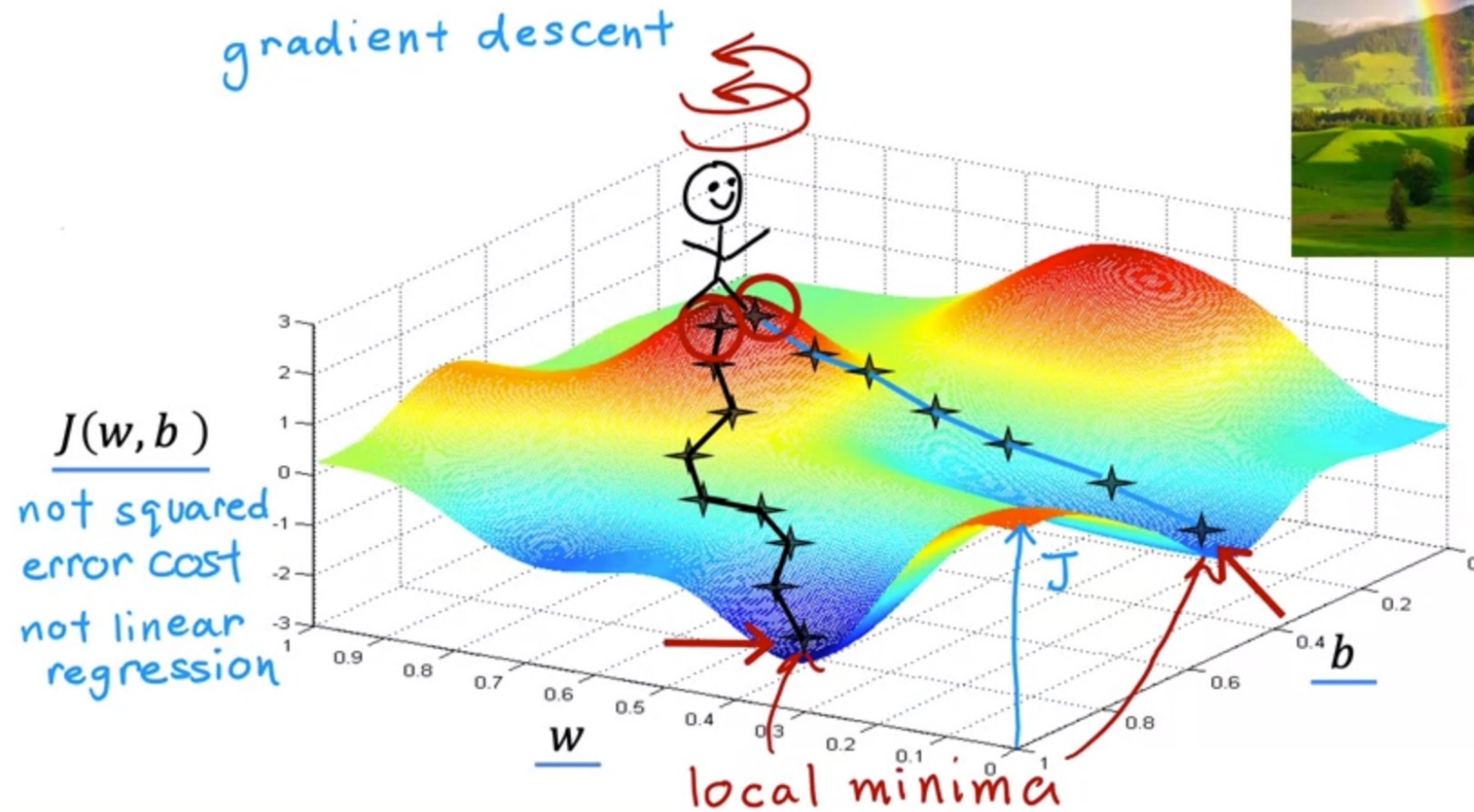
Keep changing  $w, b$  to reduce  $J(w, b)$

Until we settle at or near a minimum

*may have  $>1$  minimum*



# Gradient Descent



# Gradient descent algorithm

Repeat until convergence

$$\begin{cases} \underline{w} = w - \alpha \frac{d}{dw} J(w, b) \\ \underline{b} = b - \alpha \frac{d}{db} J(w, b) \end{cases}$$

Learning rate  
Derivative

Simultaneously  
update  $w$  and  $b$

Assignment

$$a = c$$

$$a = a + 1$$

Code

Truth assertion

$$a = c$$

$$a = a + 1$$

Math

$$a == c$$

# Gradient descent algorithm

Repeat until convergence

$$\left\{ \begin{array}{l} \underline{w} = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \underline{b} = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{array} \right.$$

Learning rate  
Derivative

Simultaneously  
update  $w$  and  $b$

Assignment

$$\begin{array}{l} a = c \\ \quad \quad \quad \swarrow \\ a = a + 1 \\ \text{Code} \end{array}$$

Truth assertion

$$\begin{array}{l} a = c \\ a = a + 1 \\ \quad \quad \quad \times \\ \text{Math} \\ a == c \end{array}$$

Correct: Simultaneous update

$$tmp\_w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$tmp\_b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$w = tmp\_w$$

$$b = tmp\_b$$

Incorrect

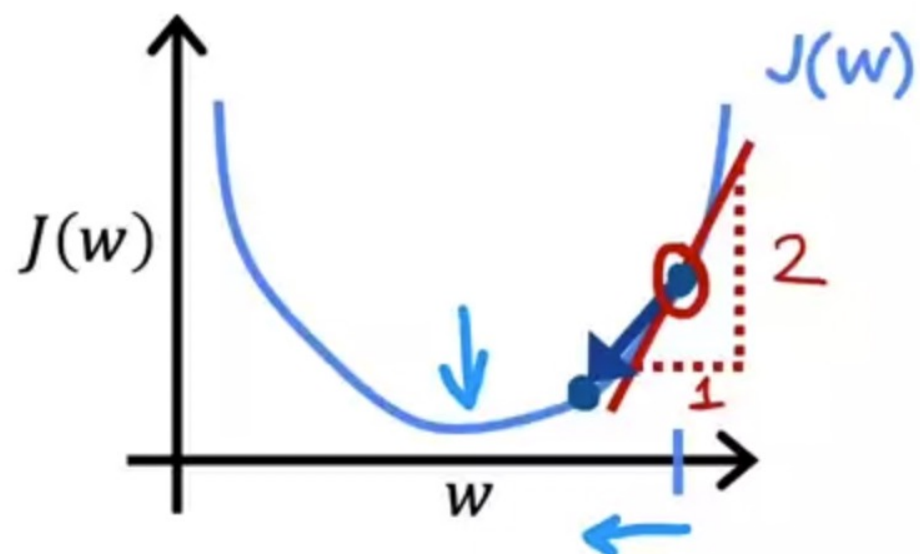
$$tmp\_w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$\underline{w} = tmp\_w$$

$$tmp\_b = b - \alpha \frac{\partial}{\partial b} J(\underline{w}, b)$$

$$\underline{b} = tmp\_b$$

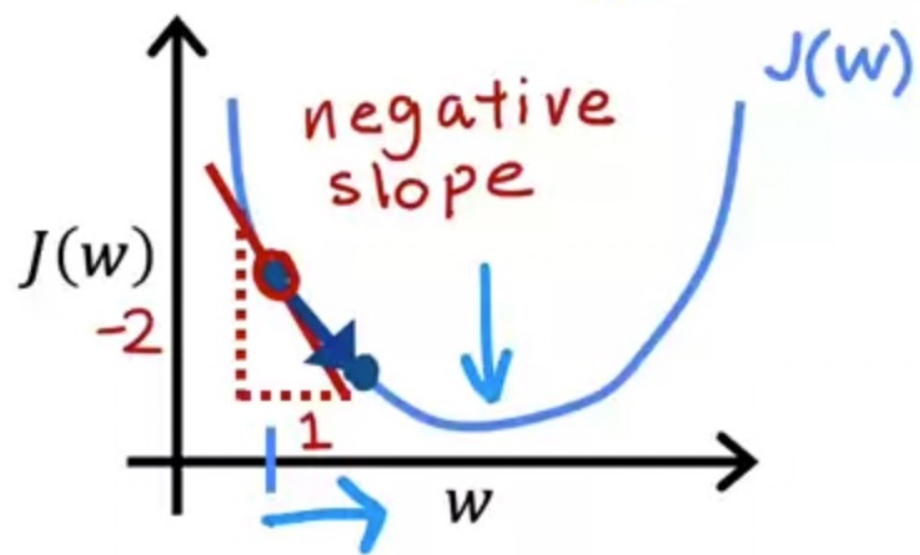




$$w = w - \alpha \underbrace{\frac{d}{dw} J(w)}_{>0}$$

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

$$\underline{\frac{d}{dw} J(w)} < 0$$



$$w = \underline{w} - \alpha \cdot (\text{negative number})$$

$\uparrow$                        $\uparrow$

$$w = w - \alpha \frac{d}{dw} J(w)$$

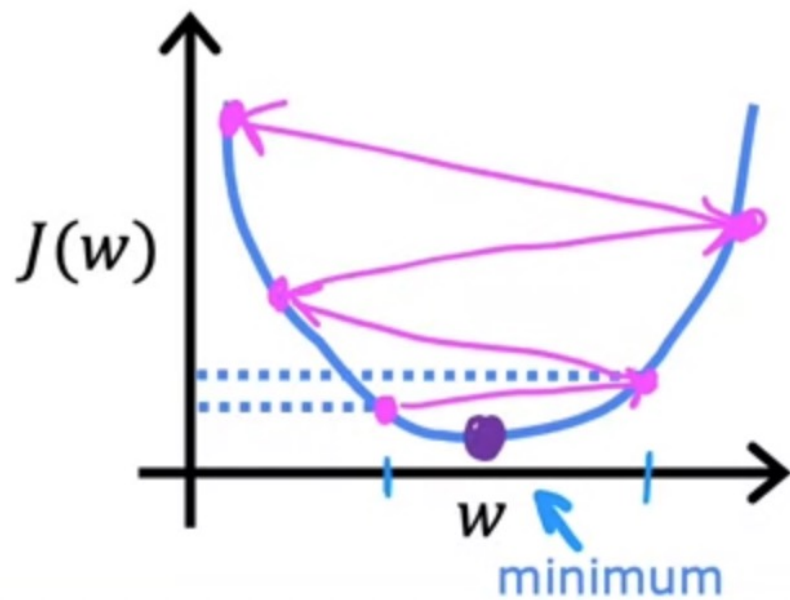
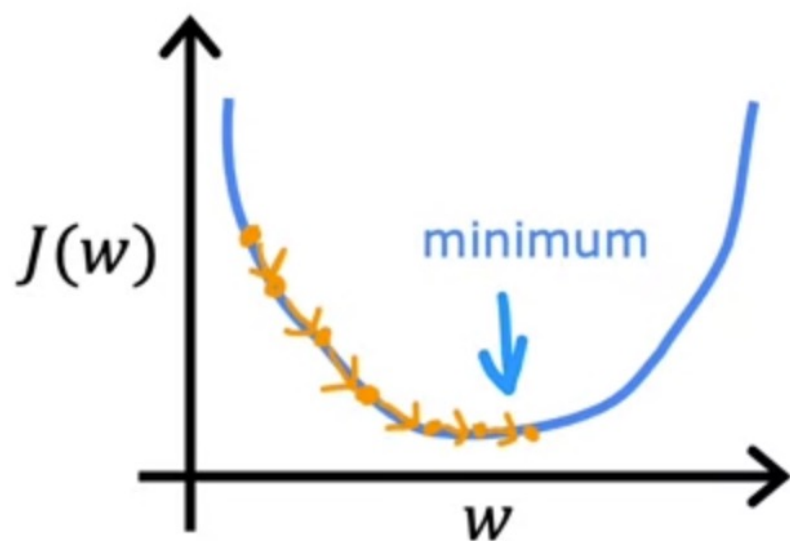
If  $\alpha$  is too small...

Gradient descent may be slow.

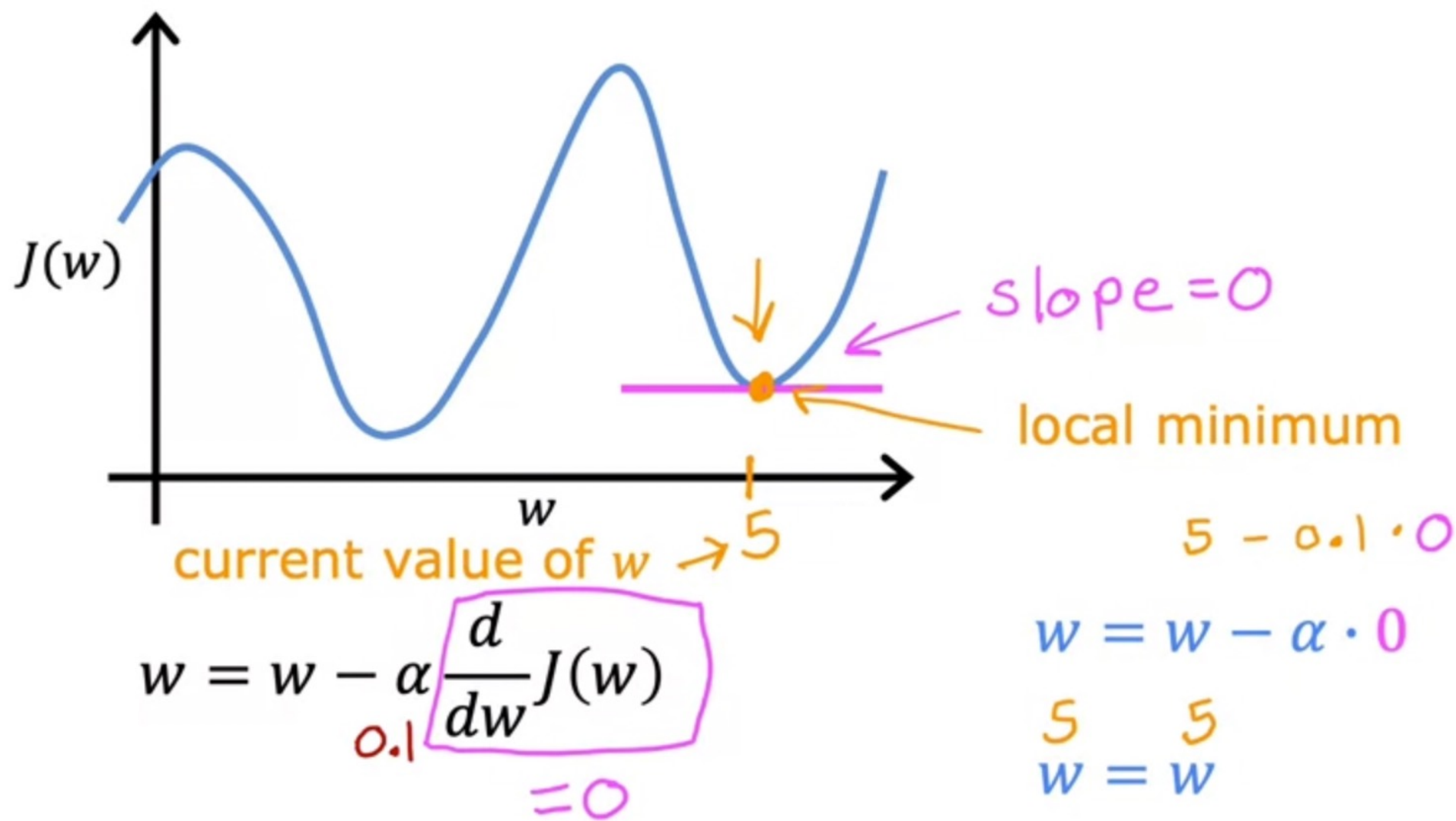
If  $\alpha$  is too large...

Gradient descent may:

- Overshoot, never reach minimum
- Fail to converge, diverge







Can reach local minimum with fixed learning rate

$\alpha$

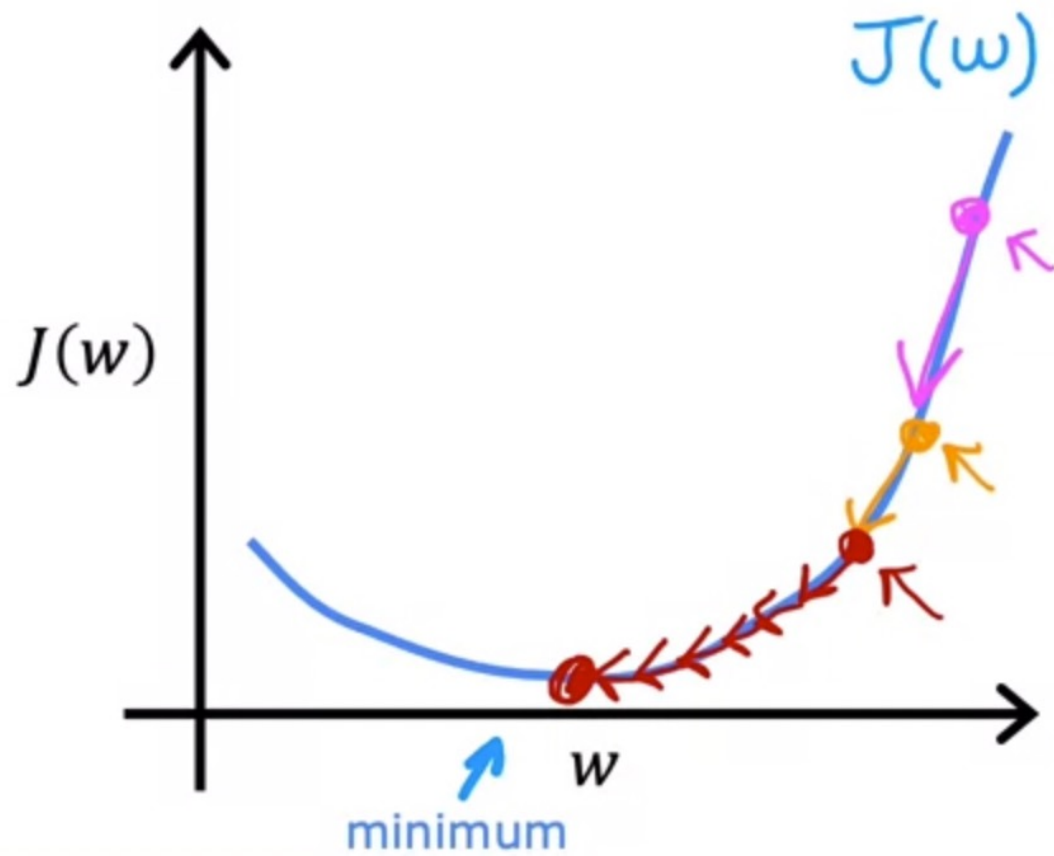
$$w = w - \underbrace{\alpha \underbrace{\frac{d}{dw} J(w)}}_{\text{large}}$$

smaller  
not as large

Near a local minimum,

- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate  $\alpha$



# Derived using calculus

Linear regression model

$$f_{w,b}(x) = wx + b$$

Cost function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Gradient descent algorithm

repeat until convergence {

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

}

# Derivative

$$\begin{aligned}\frac{\partial}{\partial w} J(w, b) &= \frac{d}{dw} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{d}{dw} \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)})^2 \\ &= \cancel{\frac{1}{2m}} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2} x^{(i)} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}}\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial b} J(w, b) &= \frac{d}{db} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{d}{db} \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)})^2 \\ &= \cancel{\frac{1}{2m}} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})} \\ &\quad \text{no } x^{(i)}\end{aligned}$$

# Gradient descent algorithm

$$\frac{\partial}{\partial w} J(w, b)$$

repeat until convergence {

$$w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

}


$$\frac{\partial}{\partial b} J(w, b)$$

Update  
w and b  
simultaneously

$$f_{w,b}(x^{(i)}) = wx^{(i)} + b$$



squared error cost

bowl shape   
convex function

