

$$J(w, b) = \frac{1}{2m} \sum (\hat{y}^i - y^i)^2$$

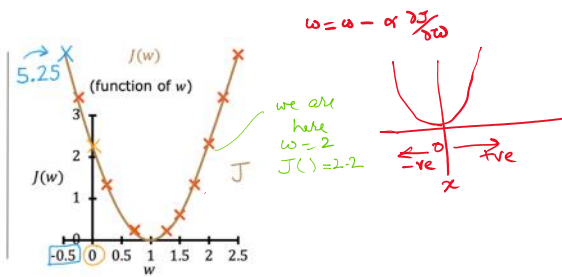
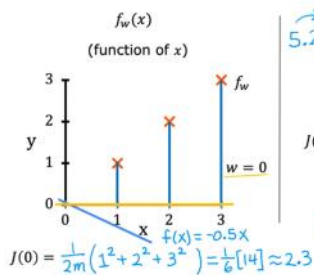
$$J(w, b) = \frac{1}{2m} \sum (wx^i + b - y^i)^2$$

$$J(w) = \frac{1}{2m} \sum (wx^i + b - y^i)^2$$

$w = w - \alpha \frac{\partial J}{\partial w}$

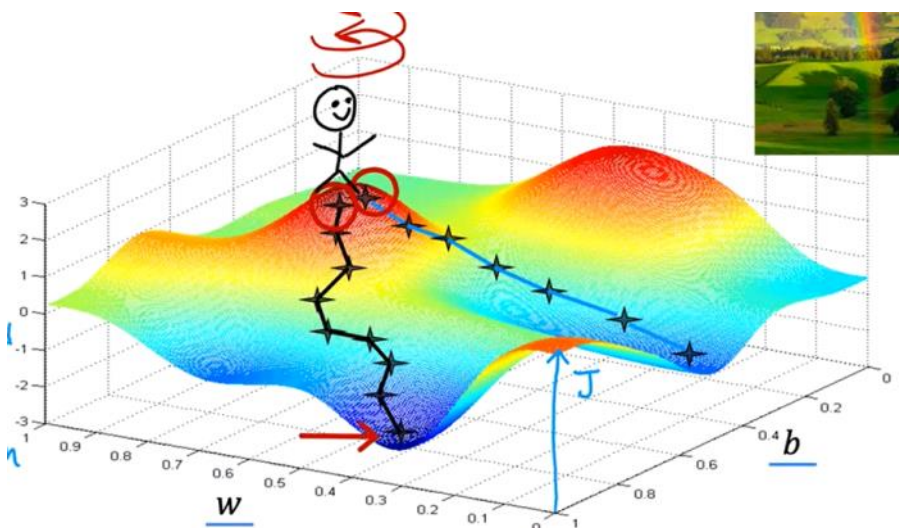
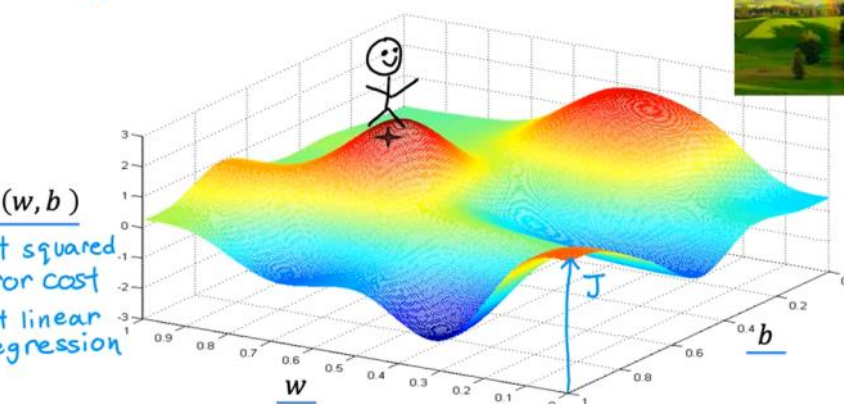
$J(w)$

$w \rightarrow$



gradient descent

$J(w, b)$
not squared
error cost
not linear
regression



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

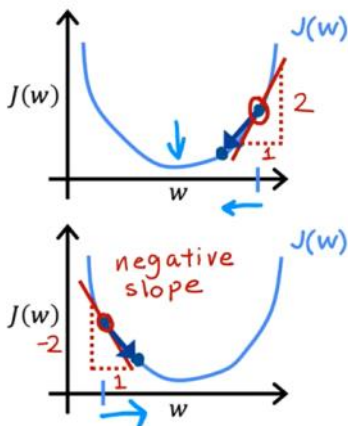
Simultaneously update w and b

Correct: Simultaneous update

$$\begin{aligned} tmp_w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ tmp_b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w &= tmp_w \\ b &= tmp_b \end{aligned}$$

Incorrect

$$\begin{aligned} tmp_w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \underline{w} &= tmp_w \\ tmp_b &= b - \alpha \frac{\partial}{\partial b} J(\underline{w}, b) \\ b &= tmp_b \end{aligned}$$



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

> 0

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

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

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

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

no $x^{(i)}$

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

If α is too small...

Gradient descent may be slow.

If α is too large...

Gradient descent may:

- Overshoot, never reach minimum

