→ Here we will see how different learning rate effects the cost func graph.



$$w = w - \textcircled{$\alpha$} \frac{d}{dw} J(w)$$
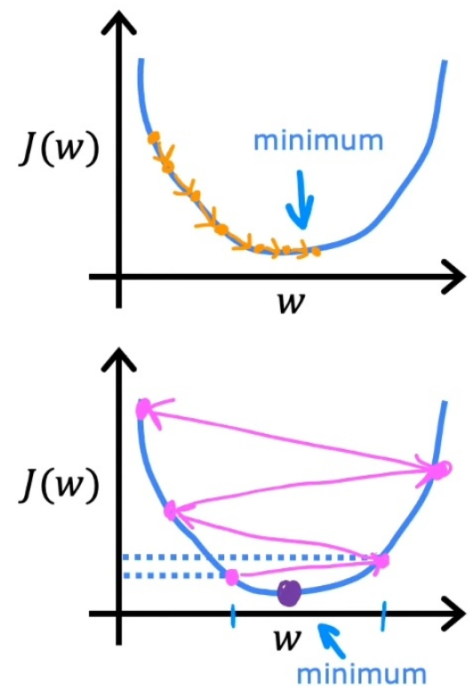
If $\alpha$ is too <u>small</u>...
Gradient descent may be <u>slow</u>.

If $\alpha$ is too <u>large</u>...

Gradient descent may:
 - <u>Overshoot</u>, never reach minimum
 - Fail to converge, diverge

→ If we choose $\alpha$ as a small value, if will take tiny steps toward local minima
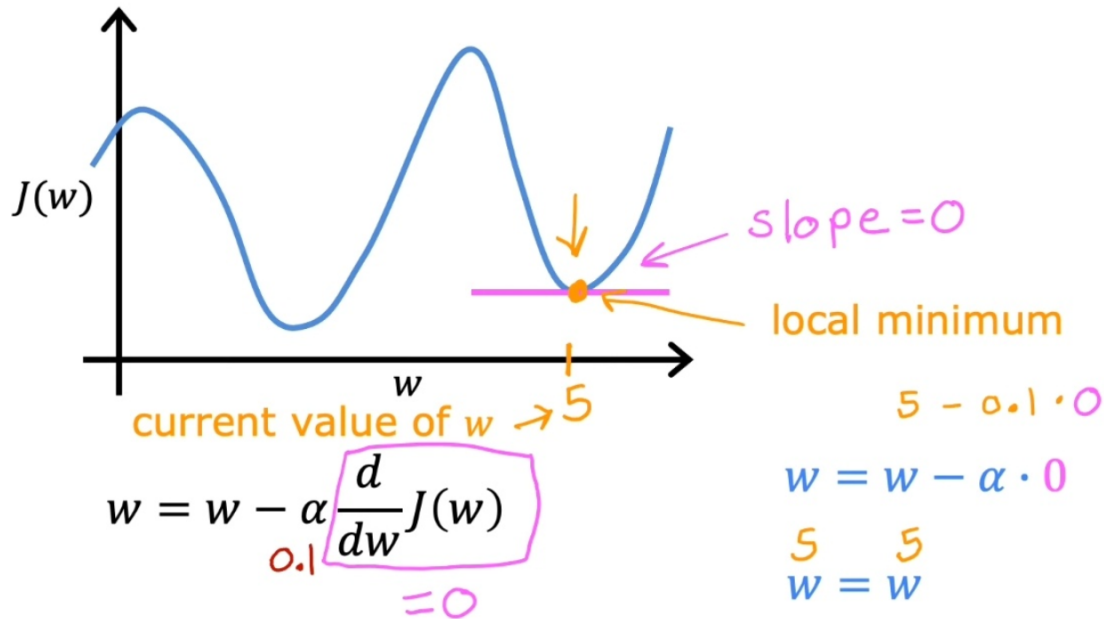
→ so it will be slow.

→ If $\alpha$ is too high, it may work opposite.

→ Cause it will take huge steps and it may increase the $J(w)$. since the

derivatives term is changing.

→ The value of w and b will increase.



→ The main drawback of this formula
is, it always stucks on local minima.

→ Cause, at the bottom of the curve
the Slope is 0. So, it becomes w=w-o

→ the value of w won't change.

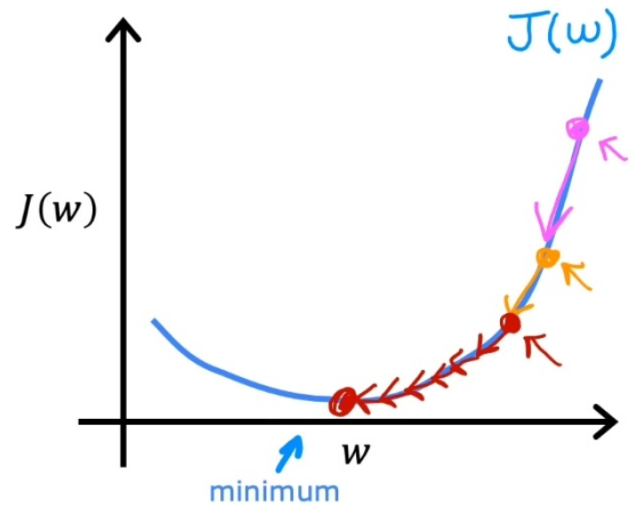# Can reach local minimum with fixed learning rate $\alpha$

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

- smaller
- not as large
- large

Near a local **minimum**,
- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate $\alpha$

→ We do not need to change the learning rate. It will work with fixed $\alpha$

→ At first step it will take large step. after that it will take small steps as reaching toward local minima.

→ cause the slope value is decreasing as we reaching toward local minima.