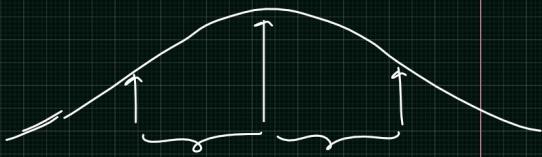
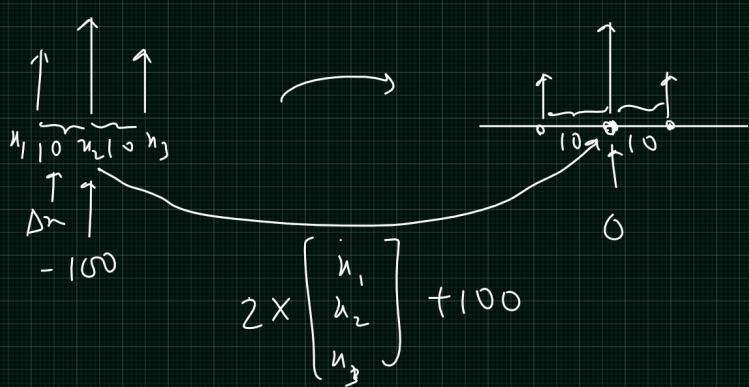
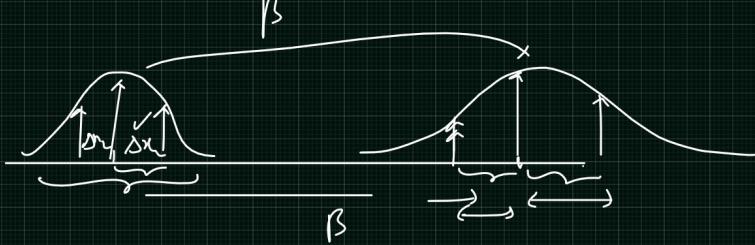
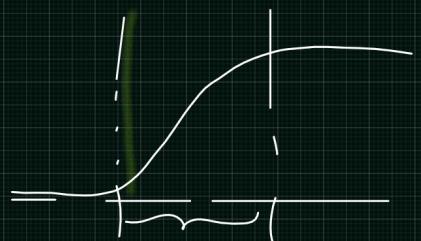
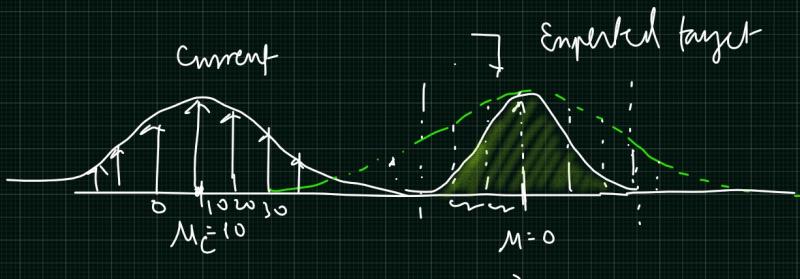


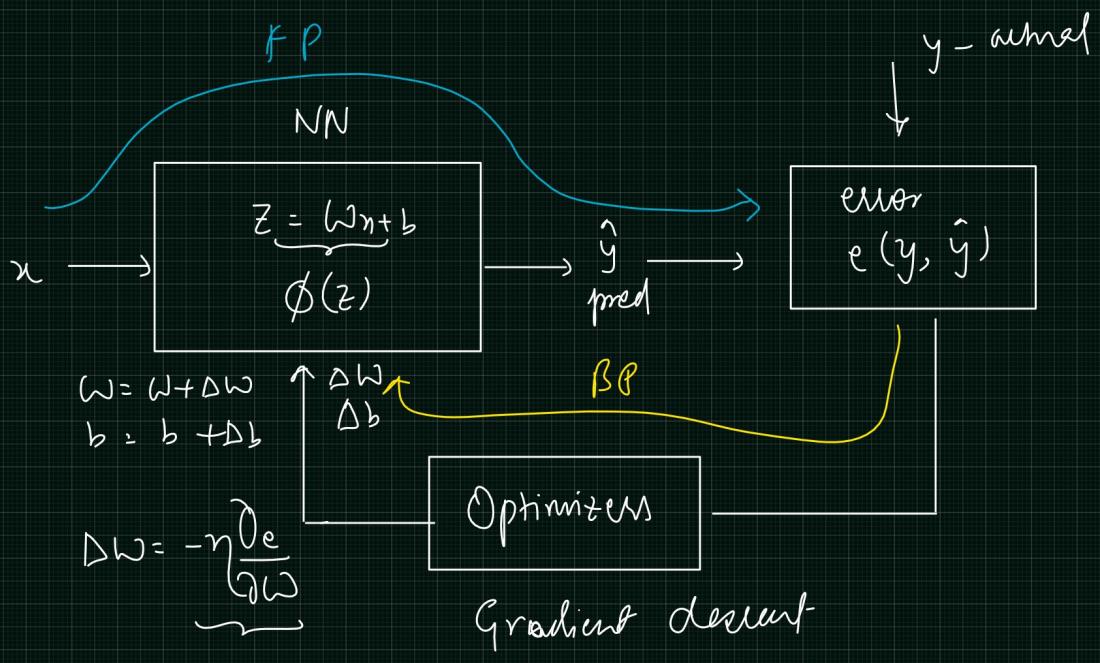
4)

$$\tilde{z}^{(i)} = \gamma \otimes \chi^{(i)} + \beta$$

Scaling      Shifting



# FAST OPTIMIZERS



Observation on Grad. descent

$$\omega_{\text{new}} = (\omega_{\text{old}} - \eta \underbrace{\frac{\partial e}{\partial \omega}}_{\omega = \omega_{\text{old}}})$$

$$\frac{\partial \mathcal{L}}{\partial \omega} (\omega = \omega_{\text{old}})$$

$$\theta = \theta - \eta \nabla_{\theta} J(\theta) \quad \begin{matrix} \text{or} \\ \text{Cost function} \end{matrix}$$

parameters  $(\omega, b)$

$$\begin{pmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ b_1 \\ b_2 \\ \vdots \\ \theta \end{pmatrix} = \begin{pmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ b_1 \\ b_2 \\ \vdots \\ \theta \end{pmatrix} - \eta \underbrace{\begin{pmatrix} \frac{\partial J}{\partial \omega_1} \\ \frac{\partial J}{\partial \omega_2} \\ \vdots \\ \frac{\partial J}{\partial b_1} \\ \frac{\partial J}{\partial b_2} \\ \vdots \\ \nabla_{\theta} J(\theta) \end{pmatrix}}_{\nabla_{\theta} J(\theta)}$$

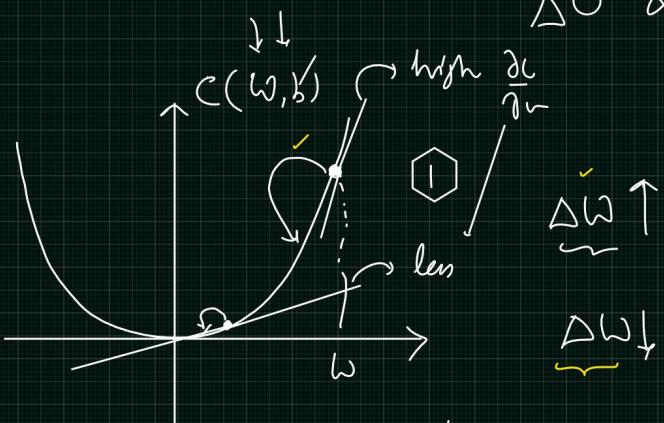
$$\omega = \omega - \eta \frac{\partial C}{\partial \omega}$$

$\Delta \omega = -\eta \frac{\partial C}{\partial \omega} = (-\eta) \frac{\partial C}{\partial \omega}$

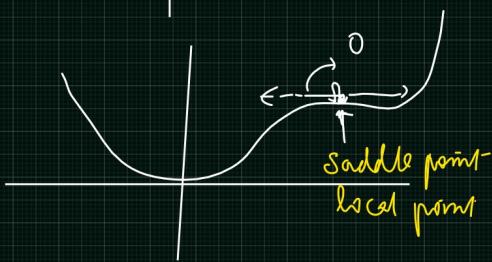
$\Delta \omega \propto \frac{\partial C}{\partial \omega} \rightarrow$  slope of the tangent

$$\Delta \theta \propto \nabla_{\theta} J(\theta)$$

cost fn.



if  $\frac{\partial C}{\partial \omega} \uparrow \rightarrow$  slope steep  
if  $\frac{\partial C}{\partial \omega} \downarrow \rightarrow$  near flat surface  
flur surface



②  $\Delta \omega$  only depends on current weight

$$C(y, \hat{y})$$

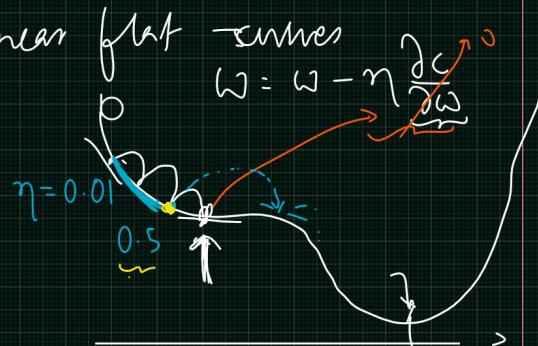
↑

$$C(y, \sigma(w_n + b))$$

$\Rightarrow$  it doesn't get influenced by previous wt-updates.

③ Grad. descent is very slow at near flat curves  
 $\Rightarrow$  slow training

because GD doesn't consider previous updates



$$100 \rightarrow 0.03$$

$\uparrow$

$$10\%$$

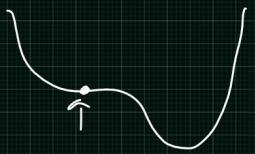
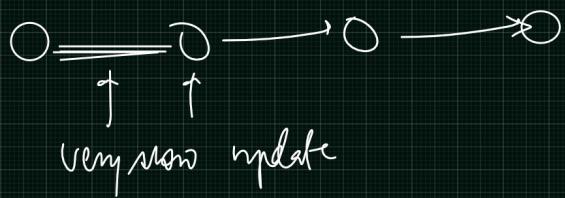
$$\Delta \omega \propto \{-\eta\} \frac{\partial C}{\partial \omega}$$

$\uparrow$

$$100 \rightarrow -\eta_1, \quad 200 \rightarrow \eta_2$$

④ Learning rate schedules { callbacks }

$$\textcircled{6} \quad \text{Chain rule} \rightarrow \frac{\partial e}{\partial w} = \frac{\partial e}{\partial c} \cdot \frac{\partial c}{\partial z} \cdot \frac{\partial z}{\partial w}$$

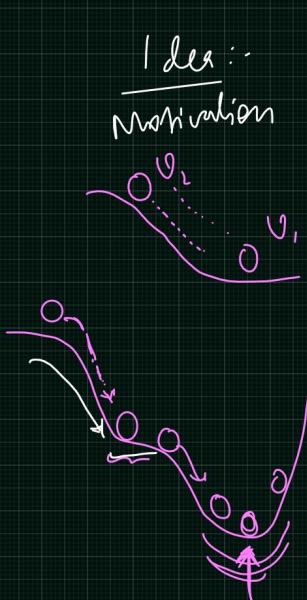


### Sol: Fast Optimizers

- i) Momentum Optimization ✓
- ii) Nesterov Accelerated Gradient (NAG)
- iii) AdaGrad ✓
- iv) RMSprop -
- v) Adam Optimizer -

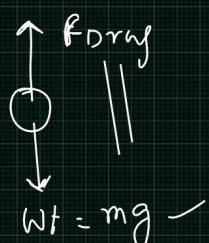
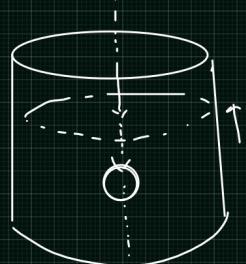
### Momentum Optimization :-

Given by Borts Polyak in 1964



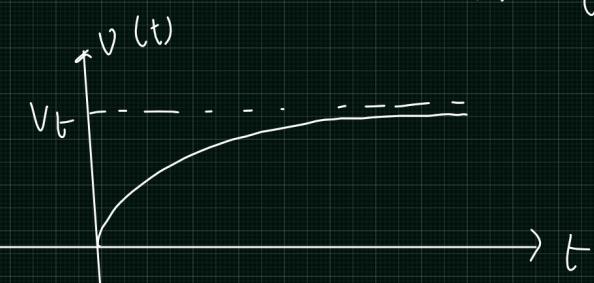
$v_{\text{terminal}}$  if friction is less

$$p = m v \\ F = \frac{dp}{dt} = m \frac{dv}{dt} = m a$$



$$F_{\text{drag}} = \omega r \Rightarrow a = 0$$

$v = \text{const}$   
 $v_{\text{terminal}}$  velocity

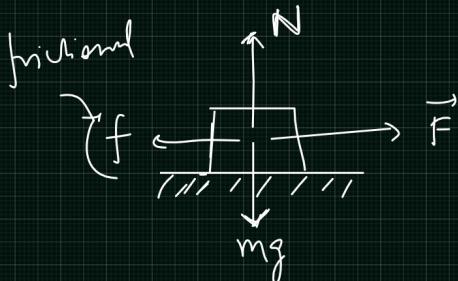


## Algorithm of momentum optimization

$m$  = momentum

$\beta$  = coeff of momentum

analogous to  
coefficients of friction



$$f = \mu N$$

$\downarrow$   
coeff of friction

$M = 0 \Rightarrow$  smooth surface

$$f = 0$$

Argo.

$$\begin{aligned} m &\leftarrow \beta m + \eta \nabla_{\theta} J(\theta) \\ \theta &\leftarrow \theta - m \end{aligned}$$

$$\begin{aligned} m &= \beta m + \eta \frac{\partial C}{\partial \omega} \\ \omega &= \omega - m \end{aligned}$$

$\omega = \omega_{current}$

$$\theta = \theta - \left[ \beta m + \eta \nabla_{\theta} J(\theta) \right]$$

assumption  
 $bias = 0$

if  $\beta > 0$

$$\theta = \theta - \beta m - \eta \nabla_{\theta} J(\theta)$$

$$\omega = \omega - \cancel{\beta m} - \eta \frac{\partial C}{\partial \omega}$$

$$\theta = \theta - \eta \nabla_{\theta} J(\theta)$$

GD

$$\omega = \omega - \eta \frac{\partial C}{\partial \omega}$$

GD

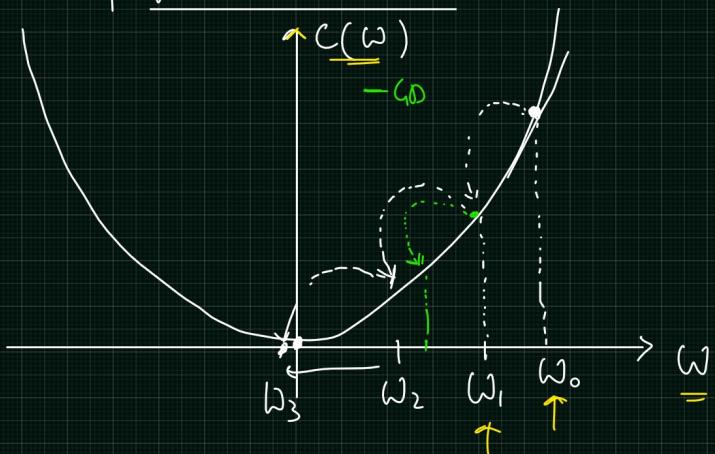
→ Through experimentation it has been observed that

—  $\beta = 0.9$  works well

but you can also make it weak

Graphical Observation :-

$$C \rightarrow C(\gamma, \dot{\gamma}) \rightarrow C(\gamma, \ddot{\omega}, b)$$



Equivalent GD

Step 1

initial,  
 $\underline{\omega = \omega_0}$ ,  $\underline{m_0 = 0}$ ,  $\beta = 0.9$

$$m_1 = \beta m_0 + \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_0}$$

$$m_1 = \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_0} \quad \text{--- (1)}$$

$$\underline{\omega_1 = \omega_0 - m_1}$$

$$\boxed{\omega_1 = \omega_0 - \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_0}} \quad \text{--- (2)}$$

Step 2 :

$$\underline{\omega = \omega_1}, \quad \underline{\beta = 0.9}, \quad \underline{m_1 = \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_0}}$$

$$m_2 = \beta m_1 + \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_1}$$

$$m_2 = \underbrace{\beta \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_0}}_{= \text{current part}} + \eta \left[ \frac{\partial C}{\partial \omega} \right]_{\omega=\omega_1}$$

$$= \eta \left[ \beta \left. \frac{\partial C}{\partial \omega} \right|_{\omega=\omega_0} + \left. \frac{\partial C}{\partial \omega} \right|_{\omega=\omega_1} \right]$$

0.9 part  
90% part

Current  
100% part

$$\rightarrow \omega_2 = \omega_1 - m_2$$

$$\left\{ \omega_2 = \omega_1 - \eta \left[ \beta \left. \frac{\partial C}{\partial \omega} \right|_{\omega=\omega_0} + \left. \frac{\partial C}{\partial \omega} \right|_{\omega=\omega_1} \right] \right\}$$

$\Delta \omega_{m_0}$

$m_0$  Eqn

$$\left\{ \omega_2 = \underline{\omega}_1 - \eta \left[ \frac{\partial c}{\partial \omega} \Big|_{\omega=\underline{\omega}_1} \right] \right\}_{GD}$$

$\Delta \omega_{GD}$

$$\underbrace{\Delta \omega_{MS}}_{\text{MS}} > \underbrace{\Delta \omega_{GD}}_{\text{GD}}$$

Why?:

$$m_2 = \eta \left[ \beta \frac{\partial c}{\partial \omega} \Big|_{\omega=\omega_0} + \frac{\partial c}{\partial \omega} \Big|_{\omega=\underline{\omega}_1} \right]$$

$$\omega = \omega_2 \quad , \quad \beta = 0.9$$

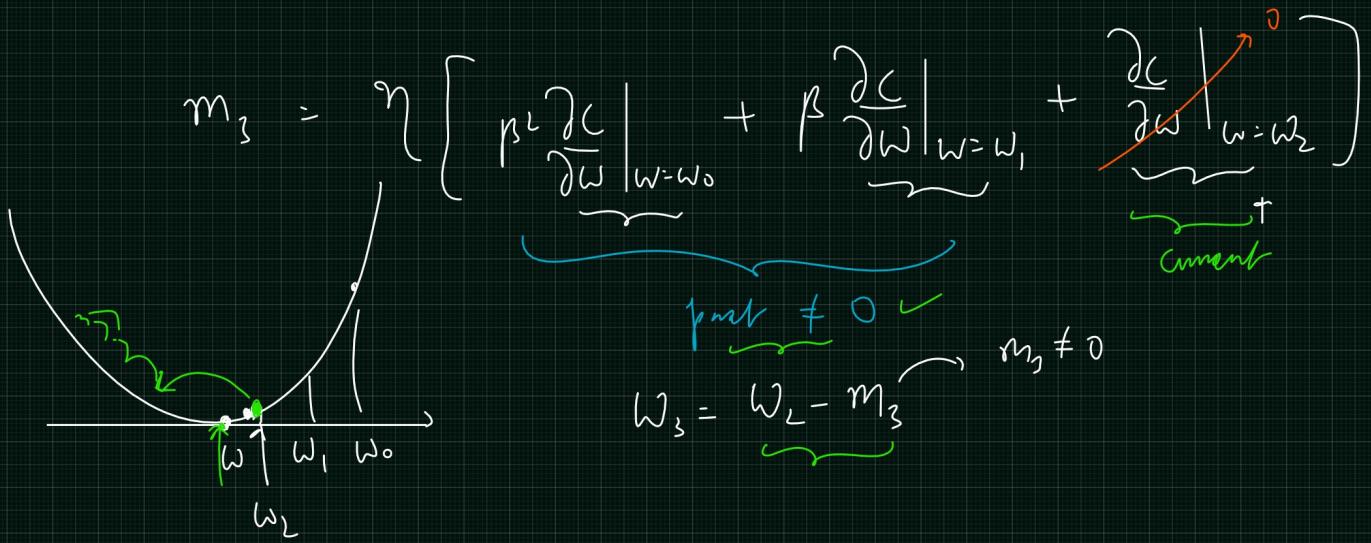
$$m_3 = \beta m_2 + \eta \frac{\partial c}{\partial \omega} \Big|_{\omega=\omega_2}$$

$$\begin{aligned} m_3 &= \beta \eta \left[ \beta \frac{\partial c}{\partial \omega} \Big|_{\omega=\omega_0} + \frac{\partial c}{\partial \omega} \Big|_{\omega=\underline{\omega}_1} \right] + \eta \frac{\partial c}{\partial \omega} \Big|_{\omega=\omega_2} \\ &= \eta \left[ \underbrace{\beta^2 \frac{\partial c}{\partial \omega} \Big|_{\omega=\omega_0}}_{0.9 \times 0.9} + \underbrace{\beta \frac{\partial c}{\partial \omega} \Big|_{\omega=\underline{\omega}_1}}_{0.9} + \underbrace{\frac{\partial c}{\partial \omega} \Big|_{\omega=\omega_2}}_{90\% \text{ of recent past}} \right] \end{aligned}$$

more unreliable

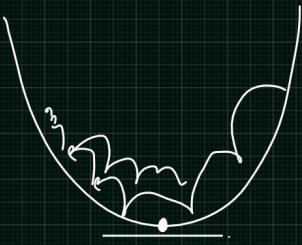
(no) Appear

$$\omega_3 = \omega_L - m_3$$



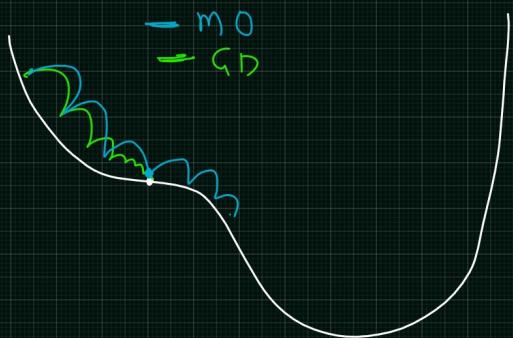
Drowbacks. ✓

- ① It oscillates when it reaches closer to minima (local / global) because of accumulation of past gradients



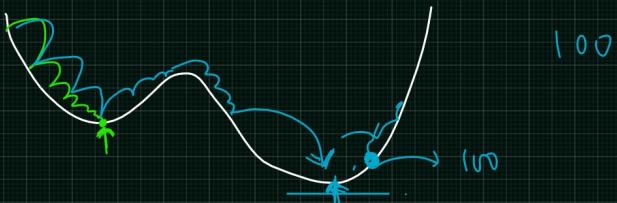
$$\begin{array}{c} \beta^2 \\ \uparrow \\ 0.81 \end{array} \quad \begin{array}{c} \beta^3 \\ \underbrace{0.81 \times 0.9} \\ \downarrow \end{array} \quad \begin{array}{c} \beta^4 \\ \dots \\ \dots \end{array}$$

$$0.9 \times 0.9 \rightarrow 0.81 \quad 0.9 > 0.81$$

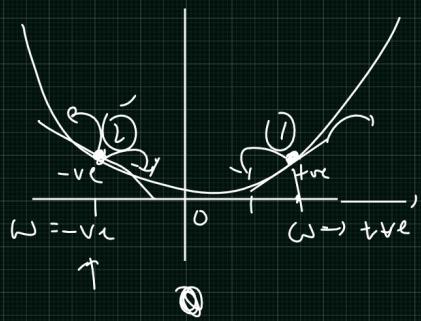
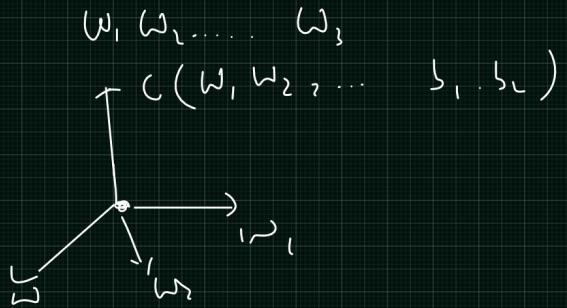
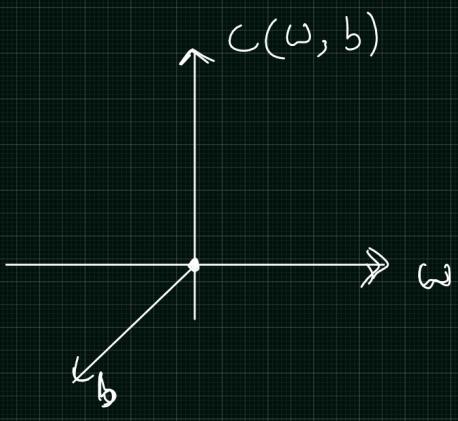


Advantages:

- ↳ Momentum can help in fast convergence
- ↳ Oscillation / Accumulation can help in overcoming local minima

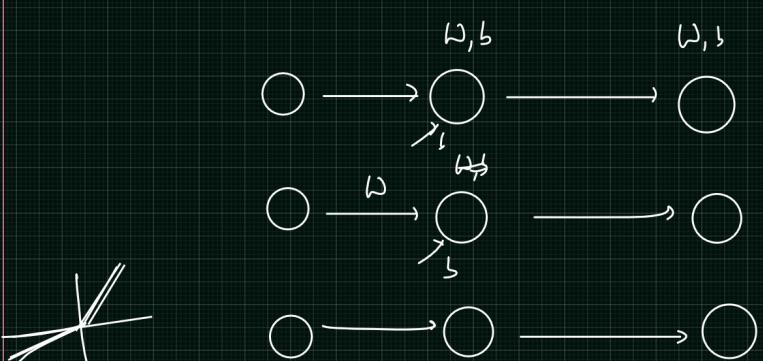


Keras  
 $\Rightarrow \text{optim} = \text{tf.keras.optimizers.SGD}(\beta=0.01, \text{momentum}=0.9)$

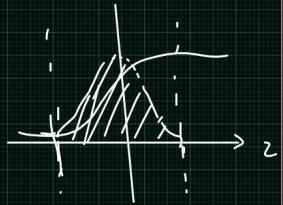


(1)      (2)

$$\begin{aligned}
 \omega &= \omega + \Delta\omega \\
 &\approx \omega - \eta \frac{\partial C}{\partial \omega} \Big|_{\omega=\omega_{\text{curr}}} \\
 &= \omega - (\underbrace{+v_e}_{+v_e} - \underbrace{-v_e}_{-v_e}) \\
 &= \underbrace{\omega - (+v_e)}_{-v_e} \underbrace{(-v_e)}_{+v_e} \\
 &= 10 - 15 = 8
 \end{aligned}$$



$$z = \underbrace{\omega_n z}_\sigma + b$$



PREW

d

time  
means  
memory  
flow

Traffic JPs  
flow rate  
speed, accuracy

