

1.(a). $W_{t+1} - W^* = (1-\gamma\sigma^2)(W_t - W^*)$
 $\therefore W^* = \frac{\gamma}{\sigma^2}$

$$\therefore W_t = W^* - W^*(1-\gamma\sigma^2)^t$$

\because stable

$$\therefore |1-\gamma\sigma^2| < 1$$

$$\therefore -1 < 1-\gamma\sigma^2 < 1$$

$$\therefore 0 < \gamma < \frac{1}{\sigma^2}, \text{ converge to } W^*.$$

(b). $W_t = (1 - (1-\gamma\sigma^2)^t) W^*$

$$W_t - W^* = (1-\gamma\sigma^2)^t \cdot W^*$$

$$\therefore W^*(1-\gamma\sigma^2)^t \leq \epsilon \cdot W^*$$

$$\therefore (1-\gamma\sigma^2)^t \leq \epsilon$$

$$\therefore t \geq \frac{\log \epsilon}{\log(1-\gamma\sigma^2)}$$

$$\therefore t_{\min} = \left\lceil \frac{\log \epsilon}{\log(1-\gamma\sigma^2)} \right\rceil$$

(c). $\begin{cases} \delta_e w[1] = y[1] \\ \delta_s w[2] = y[2] \end{cases}$

$$\therefore |1-\gamma\delta_e^2| < 1, |1-\gamma\delta_s^2| < 1$$

$$\therefore \gamma < \frac{1}{\delta_e^2}, \gamma < \frac{1}{\delta_s^2}$$

$$\because \delta_e \gg \delta_s$$

$$\therefore \gamma < \frac{1}{\delta_e^2}$$

(d). δ_e faster.

$$\begin{aligned} \text{a)} \quad A. \quad & \beta_1 \theta_{t-1} + (1-\beta_1) \hat{\theta}_t \\ B. \quad & \beta_2 v_{t-1} + (1-\beta_2) \hat{\theta}_t \end{aligned}$$

$$(b) \quad \nabla f_t'''(\theta) = \nabla f_t(\theta) + \lambda \theta$$

$$\begin{aligned} \theta_t = \theta_{t-1} - \eta \nabla f_{t-1}^{\text{reg}}(\theta_{t-1}) &= \theta_{t-1} - \eta (\nabla f_{t-1}(\theta_{t-1}) + \lambda \theta_{t-1}) \\ &= ((-\eta \lambda) \theta_{t-1} - \eta) \nabla f_{t-1}(\theta_{t-1}) \end{aligned}$$

$$\theta_t = ((-\gamma) \theta_{t-1} - \gamma) \nabla f_{t-1}(\theta_{t-1})$$

identical if $\gamma = -\eta \lambda$.

$$\text{i.e. } \gamma = \eta \lambda$$

\therefore SGD with weight decay using $\gamma = \eta \lambda$ on $f_t(\theta)$ is equal to
vanilla SGD on L_2 -regularized loss $f_t^{\text{reg}}(\theta)$.