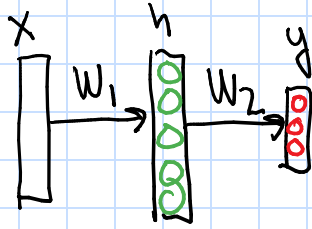


Question from Piazza

Exam #2 : 04/09
-04/12

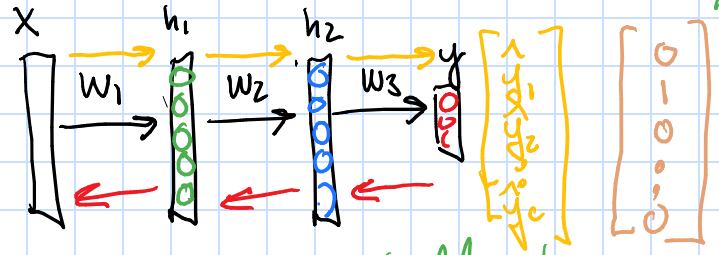


MLP w/ one hidden layer

$$h = f(W_1 x + b_1)$$

$$\hat{y} = g(W_2 h + b_2)$$

where $x \in \mathbb{R}^D$, $W_1 \in \mathbb{R}^{K \times D}$
 $b_1 \in \mathbb{R}^{K \times 1}$, $h \in \mathbb{R}^{K \times 1}$
 $W_2 \in \mathbb{R}^{C \times K}$, $y \in \{0, 1\}^{C \times 1}$



MLP w/ two hidden layers

$$h_1 = f_1(W_1 x + b_1)$$

$$h_2 = f_2(W_2 h_1 + b_2)$$

$$y = g(W_3 h_2 + b_3)$$

$$\delta_k(n) = e_k(n) Q'_k(v_k(n))$$

$$W_{ji} = \mathcal{N}(0, \sigma)$$

Local Gradients

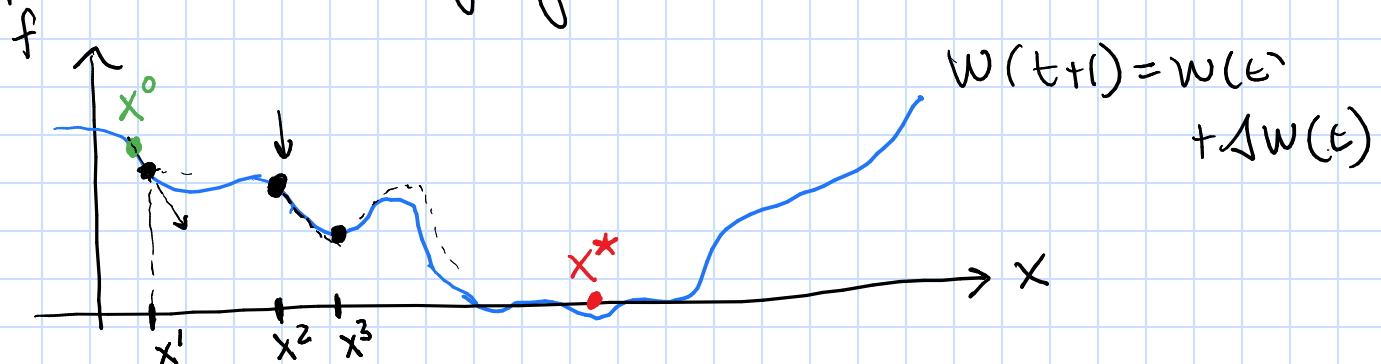
$$\delta_j(n) = e_j(n) Q'_j(v_j(n))$$

- or -

$$\delta_j(n) = Q'_j(\underline{v_j(n)}) \sum_k \delta_k(n) W_{kj}(n)$$

Momentum

- Find the weights of a neural net is a non-convex optimization problem. We are going to find a local minimum.

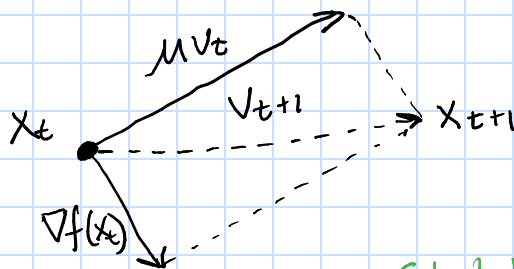


Classical Momentum

$$V_{t+1} = \mu V_t - \epsilon \nabla f(x_t)$$

$$x_{t+1} = x_t + V_{t+1}$$

convergence $\rightarrow \mathcal{O}(1/T)$

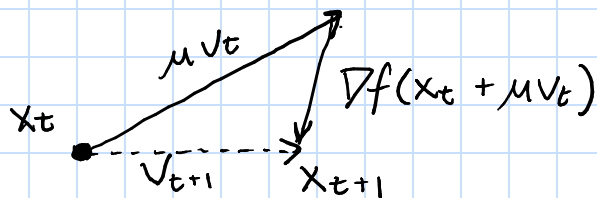


Nesterov's Accelerated Gradient

$$V_{t+1} = \mu V_t - \epsilon \nabla f(x_t + \mu V_t)$$

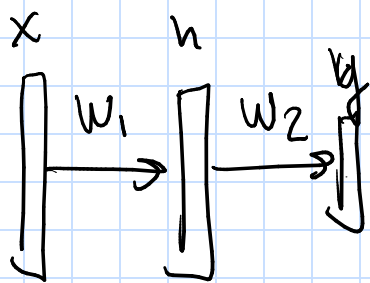
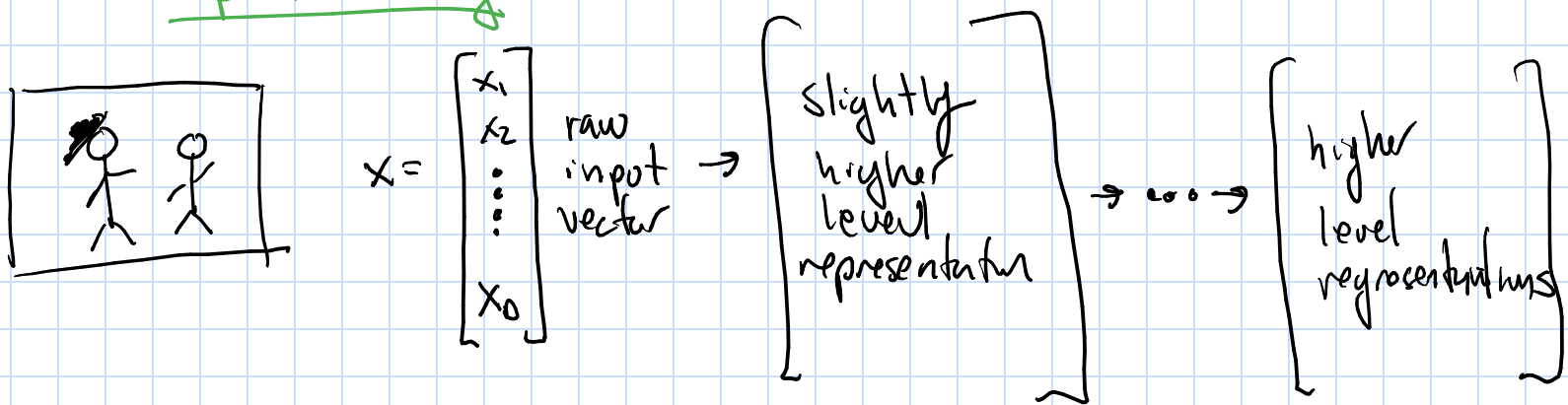
$$x_{t+1} = x_t + V_{t+1}$$

convergence $\rightarrow \mathcal{O}(1/T^2)$

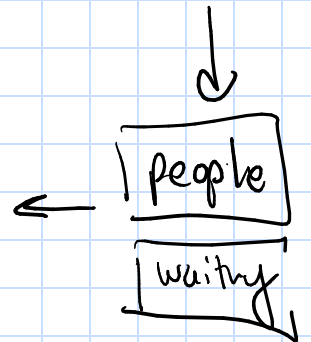


Stutskever et al. "On the importance of initialization and momentum", ICML, 2013.

Deep Learning



(even higher level of representation)



$$S_j(n) = \alpha_j(u_j(n)) \sum_k S_k(n) W_{kj}(n)$$

↑
previous layers

Traditional ML : $f: X \rightarrow y$

New Idea : $X \rightarrow H \rightarrow y$

↑
Model has hidden/latent layers.