

- Iterative optimization ; iteratively set the gradient to zero

* gradient descent (general algorithm)

$$\rightarrow \text{direction? } \Delta w = -\nabla E_{in}(w(t))$$

how much? η = learning rate

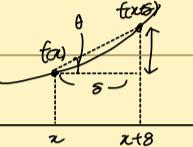
\rightarrow caution : local optimum \rightarrow start point! \Rightarrow

but-3D surface cross-entropy \Rightarrow A local optimum \rightarrow convexified.

$$W(t+1) = W(t) + \eta \cdot \nabla E_{in}(W(t)) \Rightarrow \text{pick } \eta. \quad E_{in}(W(t+1)) \text{ of } 253, \text{ vs. } E_{in}(W(t)) \text{ of } 253$$

$$\Rightarrow \text{Taylor 2nd. } f(x+s) - f(x) \approx \nabla f(x)^T s$$

$$\|s\|^2 \geq x^T y \geq -\|x\|^2, \quad b = -\frac{x}{\|x\|} \text{ grad } x \text{ at } 180^\circ$$



$$\Delta E_{in} = E_{in}(W(t+1)) - E_{in}(W(t)) \Rightarrow x \cdot y \neq 180^\circ$$

$$\approx \eta \|\nabla E_{in}(W(t))\| \quad \text{best } \eta = -\frac{\nabla E_{in}(W(t))}{\|\nabla E_{in}(W(t))\|} \text{ of } \Delta E_{in} = 34^\circ$$

$\Rightarrow \eta$: too small \rightarrow far from local optimum
too large \rightarrow bouncing around

$$\eta_t = \eta \|\nabla E_{in}(W(t))\|, \text{ norm of gradient}$$

$\Rightarrow W(t+1) = W(t) - \eta \nabla E_{in}(W(t))$ approximate OCL by stochastic gradient descent

optimization

Iter upper-bound \rightarrow $\|\nabla E_{in}(W(t))\|$ threshold \rightarrow E_{in} lower bound \rightarrow

\rightarrow batch · stochastic · mini-batch gradient descent \rightarrow

Information $h(x) := -\log p(x)$, $p(x)$ 확률밀도함수

Entropy $H[X] := -\sum p_i \log p_i$, p_i 확률, 확률밀도의 확률. $\begin{cases} -(0.5 \log 0.5 + 0.5 \log 0.5) = 0.5 \\ -(1 \log 1 + 0 \log 0) = 1 \end{cases}$

KL-Divergence $KL(p||q) := (-\sum p_i \log q_i) - (-\sum p_i \log p_i) = -\sum p_i \log \frac{q_i}{p_i}$, 확률밀도의 차이.

Cross-Entropy $CE(p, q) := -\sum p_i \log q_i$, p 확률 \rightarrow true, q 확률 \rightarrow how

$$p \log \frac{1}{q} - (1-p) \log \frac{1}{1-q}$$

$p, 1-p$ observed
 $q, 1-q$ fitted

IV. NN

function signal (forward propagation) · error signal (backward propagation)

soft max function : $\sigma = p^k \mapsto p^k, \sigma(h)_j = \frac{e^{h_j}}{\sum_k e^{h_k}}$, multinomial prob. distribution

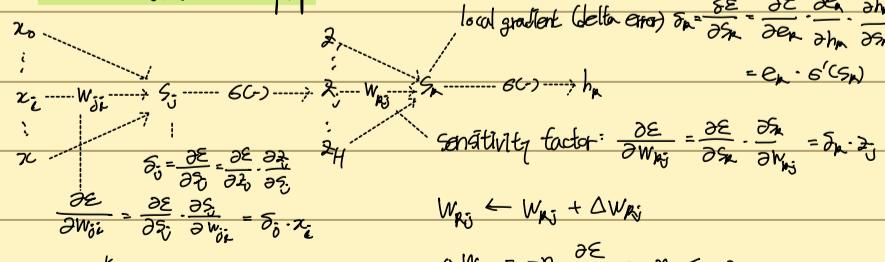
hidden unit $z_j(x) = \sigma(\sum_i w_{ij} x_i)$

output unit $h_k(x) = \sigma(\sum_j w_{kj} z_j) = \sigma(\sum_j w_{kj} \sigma(\sum_i w_{ij} x_i))$

$$E_R = h_R - y_R, \quad E_n = \frac{1}{2} \sum E_R^2, \quad E_D = \frac{1}{N} \sum E_n \rightarrow E_n = E_n(W) = \frac{1}{2} J(h_R - y_R)^2$$

$$\text{find } W^* = \arg \min_W \frac{1}{N} \sum E_n(W), \quad E(W) = E_n(W) - \eta \nabla E(W)$$

$\nabla E(W)$ aka back prop



\Rightarrow Widrow-Hoff learning rule aka LMS rule

$$\Delta w_{kj} = -\eta \cdot \epsilon_k \cdot z_j$$

single gate $\bar{t} \rightarrow f \rightarrow \text{oat}$ $\frac{\partial E}{\partial t} = \frac{\partial E}{\partial \text{oat}} \cdot \frac{\partial \text{oat}}{\partial t} \Rightarrow \frac{\partial E}{\partial \text{oat}} = F'(t)$

multiplication $\bar{t}_1, \bar{t}_2 \rightarrow f \rightarrow \text{oat}$ $\frac{\partial E}{\partial t_1} = \frac{\partial E}{\partial \text{oat}} \cdot \frac{\partial \text{oat}}{\partial t_1}, \quad \frac{\partial E}{\partial t_2} = \frac{\partial E}{\partial \text{oat}} \cdot \frac{\partial \text{oat}}{\partial t_2}$

V. CH4

VI. RNN