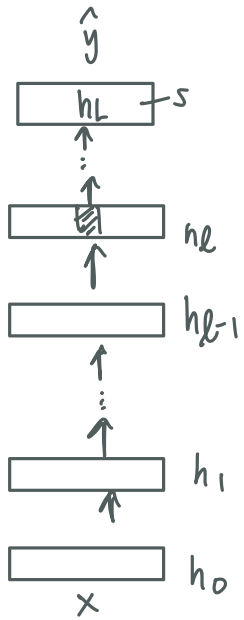


classification



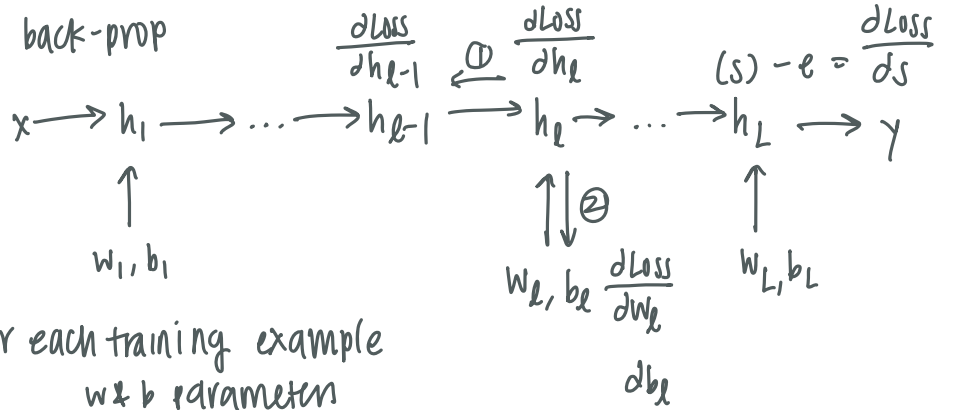
$$h_l = \sigma(s_l)$$

$$s_l = w_l h_{l-1} + b_l$$

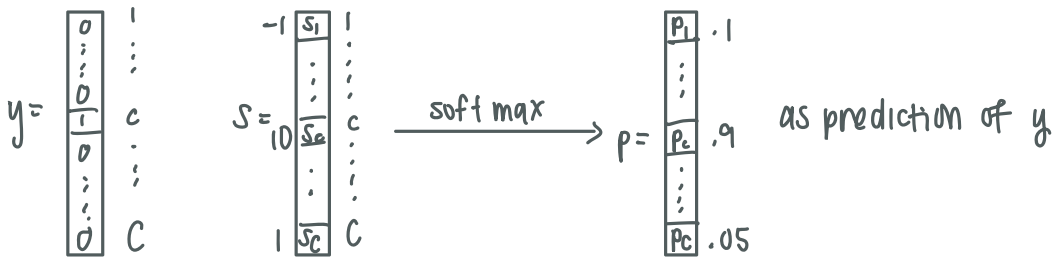
Top layer



back-prop



for each training example  
 $w$  &  $b$  parameters  
 $h$  hidden variables



$$p_c = \frac{e^{s_c}}{\sum_{c'=1}^C e^{s_{c'}}$$

summation over all categories

$$= \frac{e^{s_c}}{Z}$$

← normalizing constant

$$\sum_{c=1}^C p_c = 1$$

$$\text{Likelihood} = \prod_{c=1}^C p_c^{y_c}$$

$$= \text{prob}(\text{observed category})$$

$$\text{log-lik} = \sum_{c=1}^C y_c \log p_c$$

$$= \sum_{c=1}^C y_c (s_c - \log Z)$$

$$= \sum_{c=1}^C y_c s_c - \log Z \quad (\text{b/c } \sum_{c=1}^C y_c = 1)$$

$$\text{Loss} = -(\text{log-lik}) = -\left(\sum_{c=1}^C y_c s_c - \log Z\right)$$

$$\frac{d\text{Loss}}{ds} = \begin{pmatrix} \frac{d\text{Loss}}{ds_k} \\ \vdots \\ \frac{d\text{Loss}}{ds_k} \\ \vdots \\ \frac{d\text{Loss}}{ds_k} \\ \vdots \\ \frac{d\text{Loss}}{ds_k} \end{pmatrix} = \begin{pmatrix} y_k - \frac{d}{ds_k} \log Z \\ \vdots \\ y_k - \frac{d}{ds_k} \log Z \\ \vdots \\ y_k - \frac{d}{ds_k} \log Z \\ \vdots \\ y_k - \frac{d}{ds_k} \log Z \end{pmatrix} = -(\vec{y} - \vec{p}) = -\vec{e}$$

(top layer vector  $h_L$ )

multinomial  
logistic  
regression

$$\frac{d}{ds_k} \log Z = \frac{\frac{d}{ds_k} Z}{Z} = \frac{\frac{d}{ds_k} \sum_{c=1}^C e^{s_c}}{Z} = \frac{e^{s_k}}{Z} = p_k$$

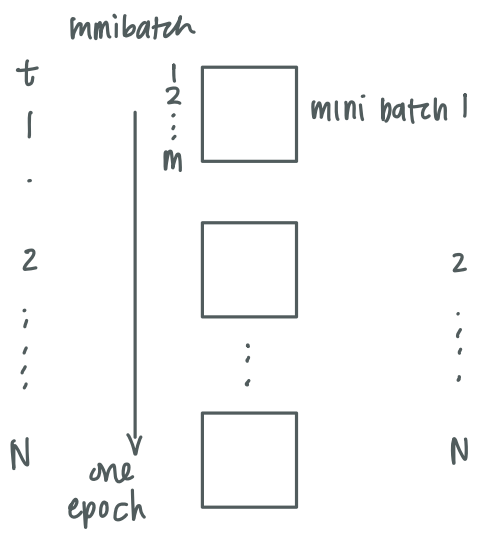
binary logistic

$$y = \begin{matrix} \boxed{1} & + \\ \boxed{0} & - \end{matrix}$$

$$s = \begin{matrix} \text{vector} \\ \boxed{s^+} & + \\ \boxed{0} & - \end{matrix} \quad \text{scalar}$$

$$p = P = \begin{matrix} \text{sigmoid} \\ \boxed{e^s} & + \\ \boxed{e^0 + e^s} & + \\ \boxed{e^0} & - \\ \boxed{e^0 + e^s} & - \end{matrix}$$

# Stochastic gradient descent (SGD)



step size, learning rate  $\eta_t \propto \frac{1}{t}$

$$\theta_{t+1} = \theta_t - \eta_t \frac{1}{m} \sum_{i=1}^m \frac{d\text{Loss}_i}{d\theta_t}$$

(mini batch)

$$\theta = (w_\ell, b_\ell, \ell = 1, \dots, L)$$

## Gradient

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_k \\ \vdots \\ x_d \end{pmatrix} \quad \Delta X = \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_k \\ \vdots \\ \Delta x_d \end{pmatrix}$$

$f(x) = f(x_1, \dots, x_k, \dots, x_d)$   
(surface, hill, valley)

## Taylor expansion

$$f(x + \Delta x) = f(x_1 + \Delta x_1, \dots, x_k + \Delta x_k, \dots, x_d + \Delta x_d)$$

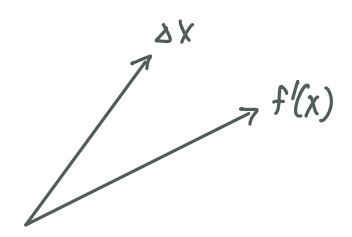
$$\doteq f(x_1, \dots, x_k, \dots, x_d) + \frac{df}{dx_1} \Delta x_1 + \dots + \frac{df}{dx_k} \Delta x_k + \dots + \frac{df}{dx_d} \Delta x_d$$

1st order Taylor expansion approximation

$$\doteq f(x) + \left\langle \begin{pmatrix} \frac{df}{dx_1} \\ \vdots \\ \frac{df}{dx_k} \\ \vdots \\ \frac{df}{dx_d} \end{pmatrix}, \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_k \\ \vdots \\ \Delta x_d \end{pmatrix} \right\rangle$$

$$\doteq f(x) + \langle f'(x), \Delta x \rangle$$

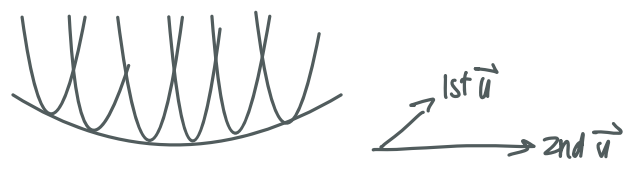
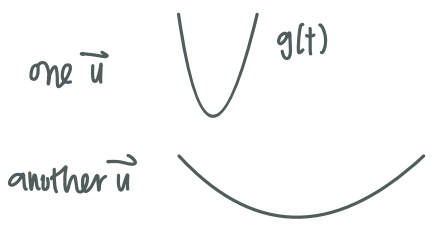
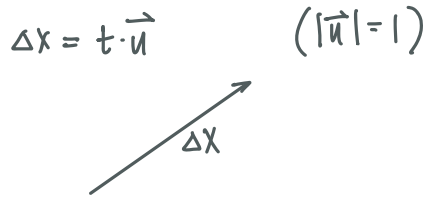
$$\doteq f(x) + |f'(x)| |\Delta x| \cos \theta$$



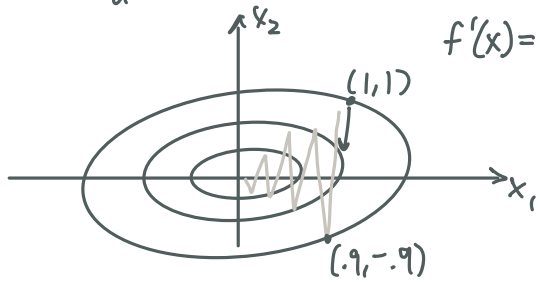
max  $\theta = 0$   
 $\Delta x \propto f'(x)$   
steepest direction  
gradient

$$g(t) = f(x + t \cdot \vec{u}) \doteq g(0) + g'(0)t + \frac{1}{2} g''(0)t^2$$

$$= f(x) + \langle f'(x), \Delta x \rangle + \frac{1}{2} \Delta x^2 f''(x) \Delta x$$



$$f(x) = \frac{x_1^2}{a^2} + x_2^2$$



$$f'(x) = 2 \begin{pmatrix} \frac{x_1}{a^2} \\ x_2 \end{pmatrix}$$

$$= 2 \begin{pmatrix} \frac{1}{100} \\ 1 \end{pmatrix} \text{ for } (x_1, x_2) = (1, 1)$$

$$\rightarrow 2 \begin{pmatrix} .9 \\ \frac{.9}{100} \\ -.9 \end{pmatrix}$$

Adam optimizer

$$g_t = \frac{1}{m} \sum_{i=1}^n \frac{dLoss_i}{d\theta_t}$$

SGD:  $\theta_{t+1} = \theta_t - \eta_t g_t$

Adam:  $v_t = \gamma v_{t-1} + (1-\gamma) g_t$

momentum velocity

$$\hat{g}_t = \beta \hat{g}_{t-1} + (1-\beta) g_t^2$$

magnitude (elementwise  $g_t^2$ )

$\beta$  between 0 & 1

$$\tilde{v}_t = v_t / (1-\gamma) \quad \tilde{g}_t = \hat{g}_t / (1-\beta) \quad \text{adaptive gradient}$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{\tilde{v}_t}{\sqrt{\tilde{g}_t + \epsilon}} \quad (\text{element wise})$$

$v_t$ : SGD  
momentum

$\hat{g}_t$ : SGD  
adaptive gradient

$$\hat{g}_0 = 0$$

$$\hat{g}_1 = (1-\beta) g_1^2$$

$$\hat{g}_2 = \beta(1-\beta) g_1^2 + (1-\beta) g_2^2$$

$$= (1-\beta)(\beta g_1^2 + g_2^2)$$

$$\hat{g}_3 = \beta(1-\beta)(\beta g_1^2 + g_2^2) + (1-\beta) g_3^2$$

$$= (1-\beta)(\beta^2 g_1^2 + \beta g_2^2 + g_3^2)$$

...

$$\hat{g}_t = (1-\beta)(\beta^{t-1} g_1^2 + \beta^{t-2} g_2^2 + \dots + \beta g_{t-1}^2 + g_t^2)$$

sparse features

accumulating recent  $g_t^2$   
recursive formula allows accumulation of magnitude.

