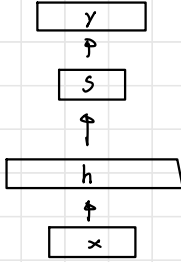


Lecture 16



- Multilayer Perceptron:
 - supervised learning

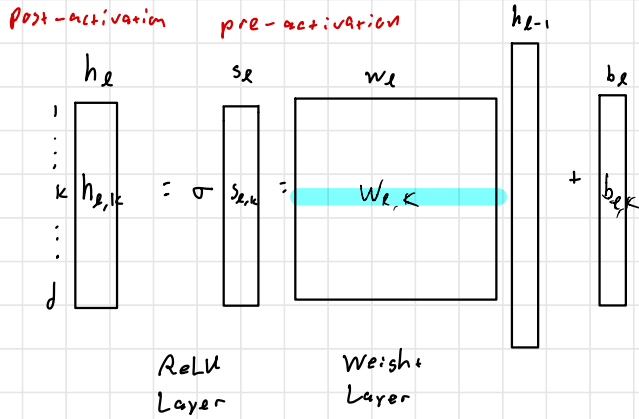
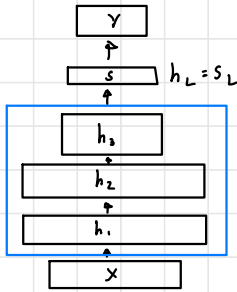


$$\log P_{\theta}(y|x):$$

$$\text{Regression: } [y|s] \sim \mathcal{N}(s, \sigma^2 I)$$

$$\text{Classification: } [y|s] \sim \text{Multinomial}(p = \text{softmax}(s))$$

• MLP:



• Error Back-Prop:

$$\frac{\partial \log_{\theta} P(y|x)}{\partial \theta} \quad \ell$$

$$\text{Loss} = -\log P_{\theta}(y|x)$$

$$\frac{\partial \text{Loss}}{\partial s} = -\text{error} = \begin{cases} -(y-s), & \text{regression} \\ -(y-p), & \text{classification} \end{cases}$$

$$\frac{\partial \text{Loss}}{\partial h_L}$$

↓ ①

$$\frac{\partial \text{Loss}}{\partial s_L} \xrightarrow{\textcircled{3}} \frac{\partial \text{Loss}}{\partial w_L}, \frac{\partial \text{Loss}}{\partial b_L}$$

↓ ②

$$\frac{\partial \text{Loss}}{\partial h_{L-1}}$$

↓
⋮

$$\textcircled{1} \quad \frac{\partial \text{Loss}}{\partial s_{\ell, k}} = \frac{\partial \text{Loss}}{\partial h_{\ell, k}} \frac{\partial h_{\ell, k}}{\partial s_{\ell, k}} = \frac{\partial \text{Loss}}{\partial h_{\ell, k}} \sigma'(s_{\ell, k})$$

$$\begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ d \end{matrix} \begin{pmatrix} \vdots \\ \vdots \\ \frac{\partial \text{Loss}}{\partial s_{\ell, k}} \\ \vdots \\ \vdots \end{pmatrix} = \begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ d \end{matrix} \begin{pmatrix} \vdots \\ \vdots \\ \frac{\partial \text{Loss}}{\partial h_{\ell, k}} \cdot \sigma'(s_{\ell, k}) \\ \vdots \\ \vdots \end{pmatrix}, \quad \frac{\partial \text{Loss}}{\partial s_{\ell}} = \frac{\partial \text{Loss}}{\partial h_{\ell}} \odot \sigma'_{\ell}$$

$$\textcircled{2} \quad \frac{\partial \text{Loss}}{\partial h_{\ell-1}^T} = \frac{\partial \text{Loss}}{\partial s_{\ell}^T} \cdot \frac{\partial s_{\ell}}{\partial h_{\ell-1}^T} = \frac{\partial \text{Loss}}{\partial s_{\ell}^T} w_{\ell}$$

w_{ℓ}

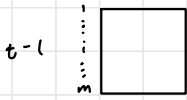
$$\textcircled{3} \quad \frac{\partial \text{Loss}}{\partial w_{\ell, k}} = \frac{\partial \text{Loss}}{\partial s_{\ell, k}} \cdot \frac{\partial s_{\ell, k}}{\partial w_{\ell, k}} = \frac{\partial \text{Loss}}{\partial s_{\ell, k}} h_{\ell-1}^T$$

$$\begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ d \end{matrix} \begin{pmatrix} \vdots \\ \vdots \\ \frac{\partial \text{Loss}}{\partial w_{\ell, k}} \\ \vdots \\ \vdots \end{pmatrix} = \begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ d \end{matrix} \begin{pmatrix} \vdots \\ \vdots \\ \frac{\partial \text{Loss}}{\partial s_{\ell, k}} \\ \vdots \\ \vdots \end{pmatrix} h_{\ell-1}^T \quad \text{so} \quad \frac{\partial \text{Loss}}{\partial w_{\ell}} = \frac{\partial \text{Loss}}{\partial s_{\ell}} h_{\ell-1}^T$$

- Stochastic Gradient Descent:

- mini-batch at iteration t :

$$\{x_i, y_i, i=1, \dots, m\}$$



$$g_t = -\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta} \log P_{\theta_t}(y_i | x_i)$$



$$\nabla_{\theta} \log P_{\theta_t}(y_i | x_i)$$

⋮



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



• Adam

- momentum + adaptive gradient

elementwise
operations

$$v_t = \gamma v_{t-1} + (1-\gamma)g_t \longrightarrow \text{cancel oscillations \& randomness}$$

$$G_t = \beta G_{t-1} + (1-\beta)g_t^2 \longrightarrow \text{crude estimate of curvature}$$

- sparse features

$$\tilde{v}_t = \frac{v_t}{1-\gamma}, \quad \tilde{G}_t = \frac{G_t}{1-\beta}$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{\tilde{v}_t}{\sqrt{\tilde{G}_t + \epsilon}}$$

: Example: $G_0 = 0$

$$G_1 = (1-\beta)g_1^2$$

$$G_2 = \beta(1-\beta)g_1^2 + (1-\beta)g_2^2$$
$$= (1-\beta)(\beta g_1^2 + g_2^2)$$

$$G_3 = \beta(1-\beta)(\beta g_1^2 + g_2^2) + (1-\beta)g_3^2$$
$$= (1-\beta)(\beta^2 g_1^2 + \beta g_2^2 + g_3^2)$$

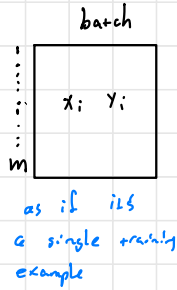
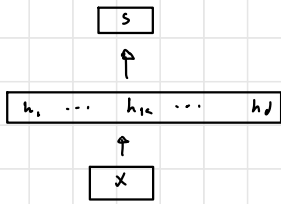
⋮

$$G_t = (1-\beta)(\beta^{t-1}g_1^2 + \beta^{t-2}g_2^2 + \dots + g_t^2)$$

- windowed - average (exponential moving average)
- accumulation of recent past

• Batch normalization:

- avoid covariance shift



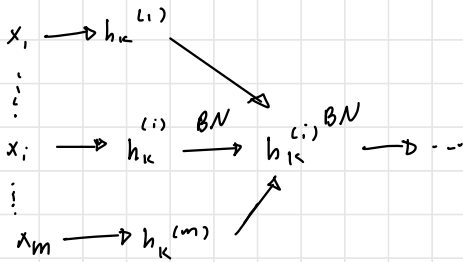
as if it's
a single training
example

$$\mu = \frac{1}{m} \sum_{i=1}^m h_{1k}(x_i)$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (h_{1k}(x_i) - \mu)^2$$

$$\tilde{h}_{1k}(x_i) = \frac{h_{1k}(x_i) - \mu}{\sigma}$$

$$h_{1k}^{BN}(x_i) = \beta \tilde{h}_{1k}(x_i) + \gamma$$



• Layer Norm:

$$\mu = \frac{1}{d} \sum_{k=1}^d h_k$$

$$\sigma^2 = \frac{1}{d} \sum_{k=1}^d (h_k - \mu)^2$$

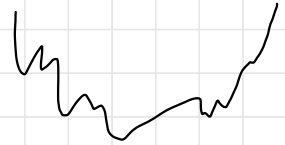
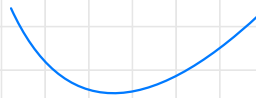
$$\tilde{h}_k = \frac{h_k - \mu}{\sigma}$$

error-correction
fault-tolerance

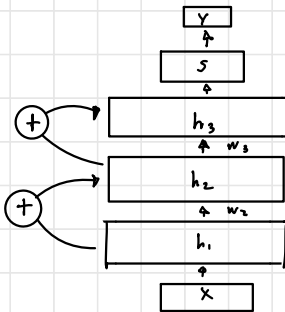
$$h_k^{LN} = \beta \tilde{h}_k + \gamma$$

Loss^{LN}(θ)

Loss(θ)



• Residual / Skip:



$$h_e = F_e(h_{e-1}) \rightarrow \text{Reconstruct existing features \& Improve}$$

(weight + BN + ReLU)

ReLU



Identity



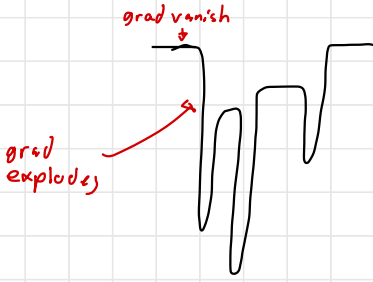
Residual parameterization:

$$h_e = h_{e-1} + f_x(h_{e-1}) \rightarrow \text{Only need improvement \& keep existing good features}$$

(weight + BN + ReLU) revision

• Original Loss

Residual



• Residual:

$$\frac{\partial h_l}{\partial h_{l-1}^T} = \mathbf{I} + \frac{\partial f_l}{\partial h_{l-1}^T}$$

• residual: $h_3 = h_2 + \sigma(w_3 h_2) = h_1 + \sigma(w_1 h_1) + \sigma(w_3 h_2)$

↑

$h_2 = h_1 + \sigma(w_2 h_1)$

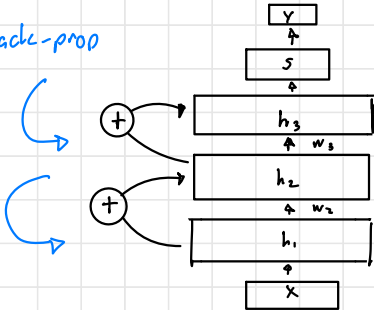
• original

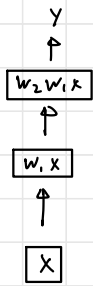
$$h_3 = \sigma(w_3 h_2)$$

↑

$$h_2 = \sigma(w_2 h_1)$$

back-prop





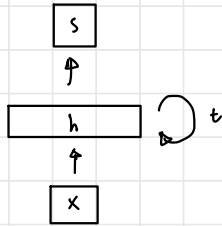
original $y = w_2 w_1 x$

ResNet :

$$\begin{aligned}
 y &= (I + w_2)(I + w_1)x \\
 &= (I + w_2)x + (I + w_2)w_1x \\
 &= x + w_2x + w_1x + w_2w_1x
 \end{aligned}$$

• Recurrent / iterative algorithm

$$h_t = h_{t-1} + f_b(h_{t-1})$$



• Slogan 2: Learned Computation (algorithm)

