

Back propagation & Optimization

Plan:

Review

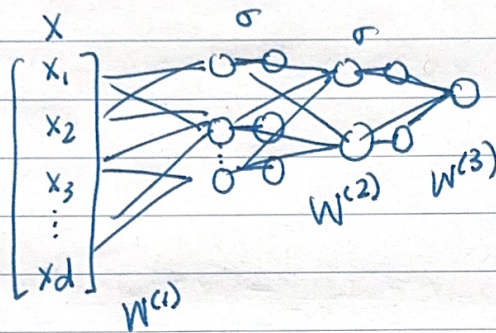
Backprop

Stochastic Gradient Descent

Logistics

Zoom + easbud

Review

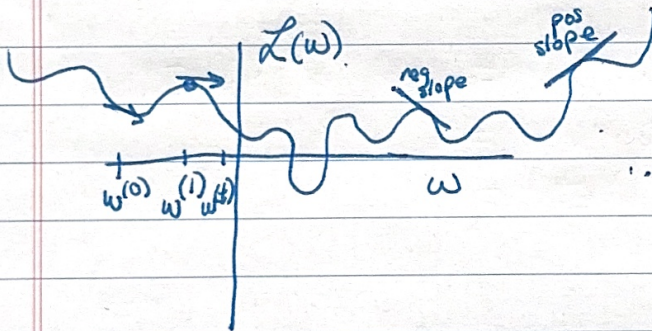


Activations

- ReLU
- sigmoid
- hyper tan
- ...

Layers

- FC
- Conv
- Residual
- Attention



$$w^{(t+1)} = w^{(t)} - \alpha \nabla_w L(w)$$

Backprop

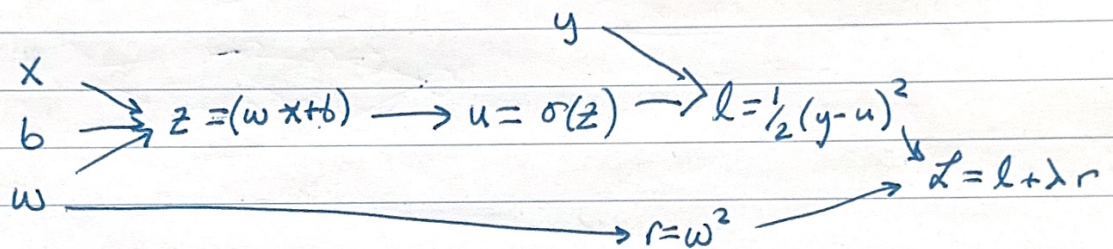
$$\mathcal{L}(w, b) = \frac{1}{2} (y - \sigma(wx+b))^2 + \lambda w^2 \quad \text{w, b} \in \mathbb{R}$$

← "regularization"

$$\frac{\partial \mathcal{L}}{\partial w} = (y - \sigma(wx+b)) \cdot \sigma'(wx+b) \cdot x + 2\lambda w$$

$$\frac{\partial \mathcal{L}}{\partial b} = (y - \sigma(wx+b)) \cdot \sigma'(wx+b)$$

⊖ complicated ⊖ Redundant



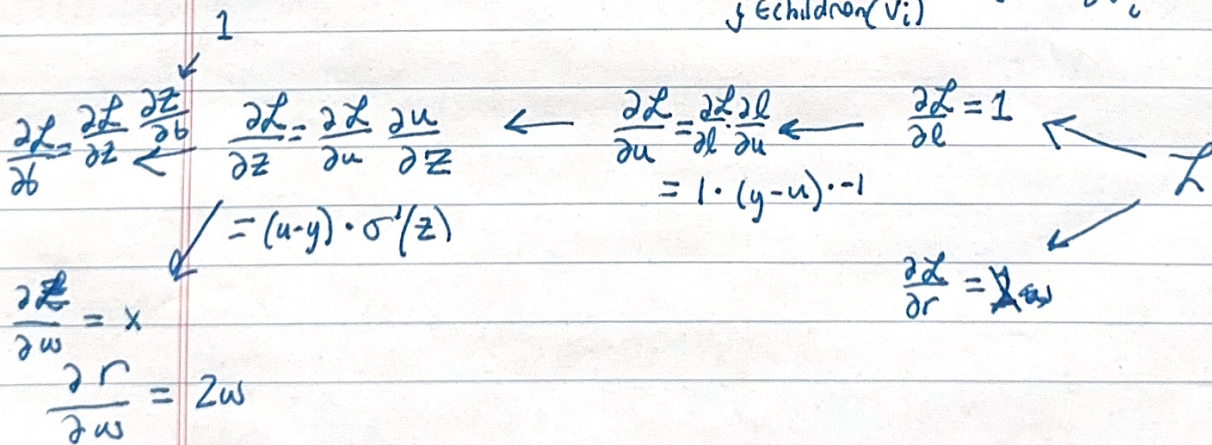
Forward

for $i \in \{1, \dots, N\}$:

compute v_i as function of parents (v_i)

Backward $i \in \{N, \dots, 1\}$:

compute $\frac{\partial \mathcal{L}}{\partial v_i} = \sum_{j \in \text{children}(v_i)} \frac{\partial \mathcal{L}}{\partial v_j} \cdot \frac{\partial v_j}{\partial v_i}$



Autograd

$$\frac{\partial \mathcal{L}}{\partial v_i} = \sum_{j \in \text{children}(i)} \frac{\partial \mathcal{L}}{\partial v_j} \cdot \frac{\partial v_j}{\partial v_i}$$

$$= \left[\dots \frac{\partial \mathcal{L}}{\partial v_j} \dots \right] \begin{bmatrix} \frac{\partial v_j}{\partial v_i} \\ \vdots \\ \frac{\partial v_j}{\partial v_i} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial v_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial v_n} \end{bmatrix} = \begin{bmatrix} \dots \frac{\partial \mathcal{L}}{\partial v_1} \dots \\ \vdots \\ \dots \frac{\partial \mathcal{L}}{\partial v_n} \dots \end{bmatrix} \begin{bmatrix} \frac{\partial v_1}{\partial v_1} \\ \vdots \\ \frac{\partial v_n}{\partial v_1} \end{bmatrix}$$

matrix multiplication

↳ forward!

↳ backward!

GPU* really good at this

- lots of cores → parallel
- lots of memory

* valuable so load in batches

$$\text{GD} : \frac{1}{|\mathcal{S}|} \sum_{i=1}^n \ell_i(x^{(i)}, y^{(i)})$$

vs.

$$\text{Stochastic GD} : \mathcal{S} \subseteq [n] \\ \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \ell_i(x^{(i)}, y^{(i)})$$

SGD Variants

$$w^{(t+1)} = w^{(t)} - \alpha \nabla_w \mathcal{L}(w^{(t)})$$

Momentum!

↳ get out of local optima

↳ "average" out noise

$$v_t^{(t+1)} = \beta v_t^{(t)} + (1-\beta) \nabla_w \mathcal{L}(w^{(t)})$$

$$w^{(t+1)} = w_t^{(t)} - \alpha v_t^{(t)}$$

Adaptive!

↳ adjust step with progress

$$s^{(t+1)} = \beta s^{(t)} + (1-\beta) \|\nabla_w \mathcal{L}(w^{(t+1)})\|_2^2$$

$$w^{(t+1)} = w^{(t)} - \frac{\alpha}{\sqrt{s^{(t+1)} + \epsilon}} \nabla_w \mathcal{L}(w^{(t)})$$

Adam!!

↳ both!!

$$w^{(t+1)} = w^{(t)} - \frac{\alpha}{\sqrt{s^{(t+1)} + \epsilon}} v_t^{(t)}$$

Initialization

$$f(x) \approx \sigma(W^{(L)} \sigma(W^{(L-1)} \dots \sigma(W^{(1)} x) \dots)$$

↳ Random so no symmetry

↳ Var $\neq 1$ so values don't explode or vanish

$$z^{(l+1)} = W^{(l)} z^{(l)}$$

$$W^{(l)} \in \mathbb{R}^{n_{in} \times n_{out}}$$

$$z_i^{(l+1)} = \sum_{j=1}^{n_{in}} W_{ij}^{(l)} z_j^{(l)}$$

$$\text{Var}(W_{ij}^{(l)}) = \sigma^2$$

$$\text{Var}(z_i^{(l+1)}) \stackrel{\text{indep}}{=} \sum_{j=1}^{n_{in}} \sigma^2 \text{Var}(z_j^{(l)}) = n_{in} \sigma^2 \text{Var}(z_j^{(l)})$$

$$\sigma^2 n_{in} \approx 1$$

$$\text{Next layer: } \sigma^2 n_{out} \approx 1 \quad \text{so } \sigma^2 = \frac{1}{2n_{in} + 2n_{out}}$$

Generalization

Test vs training data

- bottle necks
- early stopping
- weight decay
- dropout
- transfer learning