

Backpropagation & SGD

parameters

$$h(x) = w^T \phi(x)$$

$$\phi(x) = \sigma(A \phi(x))$$

$$\phi'(x) = \sigma'(A \phi(x))$$

$$\phi''(x) = \sigma''(A \phi(x))$$

function

$$\mathcal{L}(w) = \sum_{i=1}^n \ell(h(x_i), y_i) = \sum_{i=1}^n (h(x_i) - y_i)^2$$

$$\frac{\partial \mathcal{L}}{\partial w} = \sum_{i=1}^n (w^T \phi(x_i) - y_i) \phi(x_i)$$

Backpropagation

$$w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$

$$A \leftarrow A - \alpha \frac{\partial \mathcal{L}}{\partial A}$$

$$\frac{\partial \mathcal{L}}{\partial A} = \left(\frac{\partial \mathcal{L}}{\partial a} \right) \left(\frac{\partial a}{\partial A} \right) = \phi'(x)$$

$$\frac{\partial \mathcal{L}}{\partial A'} = \left(\frac{\partial \mathcal{L}}{\partial a} \right) \left(\frac{\partial a}{\partial A'} \right) = \phi'(x)$$

Convex Optimization (SVMs, Ridge Regression, ...)

$$\phi'(x) = \sigma'(a)$$

$$a = A \phi(x)$$

$$\frac{\partial a}{\partial a'} = A \frac{\partial \sigma(a')}{\partial a'} = \text{ReLU}'(\cdot) = \max\{0, x\}$$

Non-convex Optimization (NNs)

Normal GD

$$\frac{\partial \mathcal{L}}{\partial A} = \sum_{i=1}^n \frac{\partial \ell_i}{\partial A}$$

SGD

$$\frac{\partial \mathcal{L}}{\partial A} \approx \frac{\partial \ell_i}{\partial A}$$

Important Changes

- Initialization: random!
- SGD

- Small but steady progress
- Wide and deep minima

