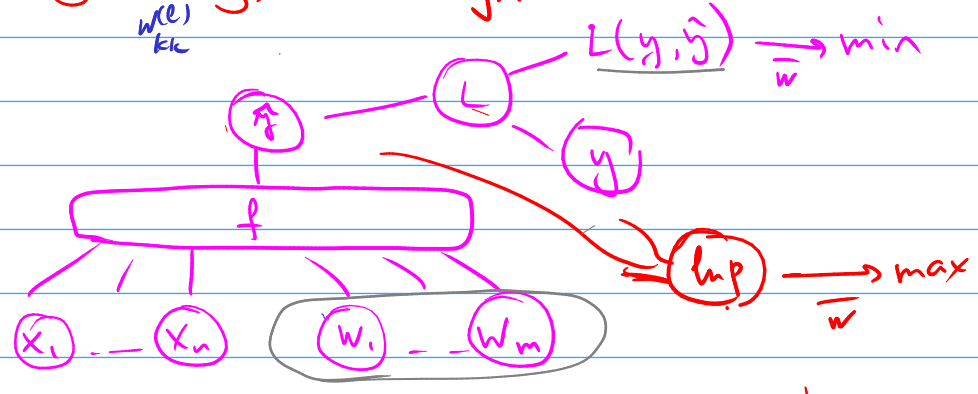


Computational graph

$$\nabla_{\bar{w}} L$$

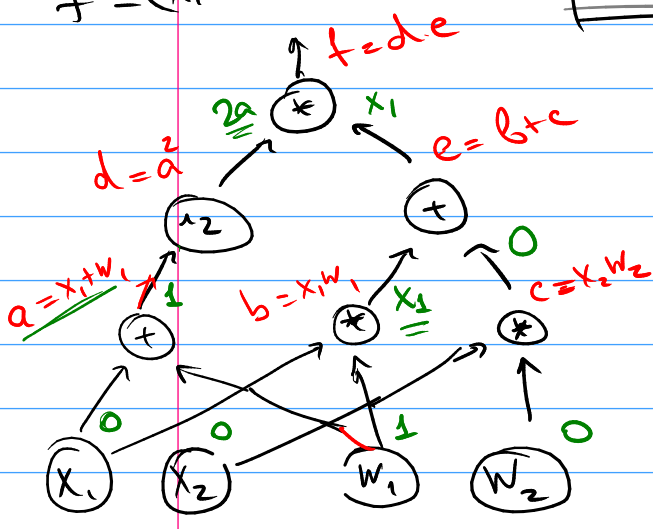


$$\bar{w} := \bar{w} - \eta \nabla_{\bar{w}} L$$

$$H \equiv \nabla_{\bar{w}} (\nabla_{\bar{w}} L)$$

quasi-Newton $H \approx UV^T$ L-BFGS

$$f = (x_1 + w_1)^2 (x_1 w_1 + x_2 w_2)$$



$$\nabla_{\bar{w}} f = \begin{pmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \end{pmatrix}$$

$$\frac{\partial f(g_1(x), \dots, g_k(x))}{\partial x} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x} + \dots + \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial x}$$

forward propagation

$$\frac{\partial a}{\partial w_1} = \frac{\partial a}{\partial x_1} \left(\frac{\partial x_1}{\partial w_1} \right) + \frac{\partial a}{\partial w_1} \left(\frac{\partial w_1}{\partial w_1} \right) = 1$$

$$\frac{\partial b}{\partial w_1} = x_1$$

$$\frac{\partial c}{\partial w_1} = 0$$

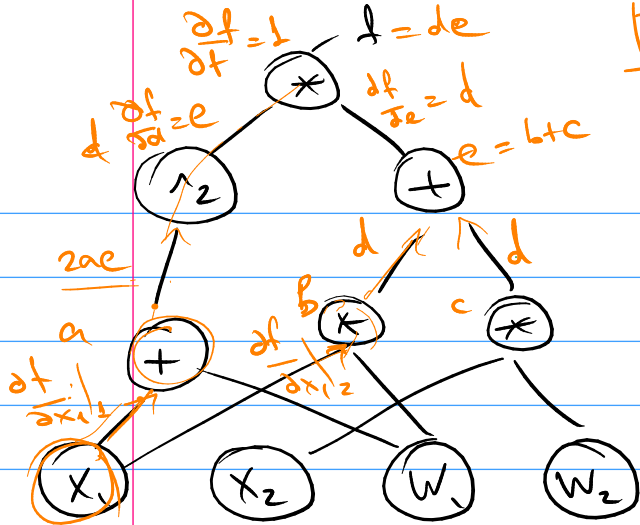
$$\frac{\partial d}{\partial w_1} = \frac{\partial d}{\partial a} \frac{\partial a}{\partial w_1} = 2a$$

$$\frac{\partial e}{\partial w_1} = x_1$$

$$\frac{\partial f}{\partial w_1} = d \cdot \frac{\partial e}{\partial w_1} + e \cdot \frac{\partial d}{\partial w_1} = 2ae + dx_1 =$$

$$= 2(x_1 + w_1)(x_1 w_1 + x_2 w_2) + (x_1 + w_1)^2 x_1$$

Backpropagation



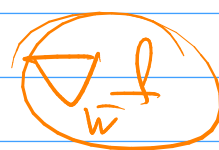
$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial e} \frac{\partial e}{\partial b} = d \cdot 1 = d$$

$$\frac{\partial f}{\partial c} = d$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial a} = 2a \cdot e$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial f}{\partial x_1} \Big|_a + \frac{\partial f}{\partial x_1} \Big|_b = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x_1} + \frac{\partial f}{\partial b} \frac{\partial b}{\partial x_1} = 2ae + dx_1$$

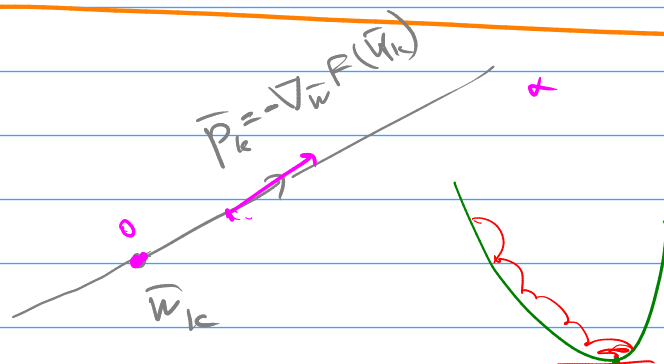
$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial w_1} + \frac{\partial f}{\partial b} \frac{\partial b}{\partial w_1} = 2ae + dx_1$$



$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial c} \frac{\partial c}{\partial w_2} = d \cdot x_2 = (x_1 + w_1)^2 \cdot x_2$$

$$F(\bar{w}) \rightarrow \min$$

$$\bar{w}_{k+1} = \bar{w}_k - \alpha \nabla_{\bar{w}} F(\bar{w}_k)$$



$$\min_{\alpha} F(\bar{w}_k + \alpha \bar{P}_k) = \psi_k(\alpha)$$

Armijo's rule

$$\psi_k(\alpha) \leq \psi_k(0) + c_1 \alpha \cdot \psi_k'(0)$$

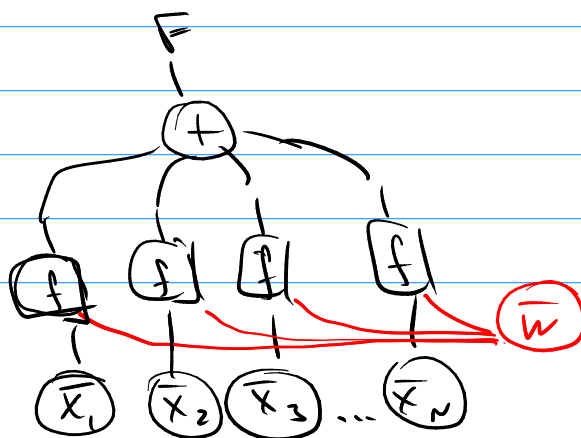
$$0 < c_1 < 1/2$$

$$|\psi_k'(\alpha)| \leq c_2 |\psi_k'(0)|$$

$$F(\bar{w}, D) = \sum_{d \in D} f(\bar{w}, d)$$

$$\ln p(D|\bar{w}) = \sum_{d \in D} \ln p(d|\bar{w})$$

$$\nabla_{\bar{w}} F = \sum_d \nabla_{\bar{w}} f(\bar{w}, d)$$



Stochastic gradient descent

mini-batch $\bar{w}_{k+1} = \bar{w}_k - \alpha \nabla_{\bar{w}} F^{(k)}(\bar{w})$

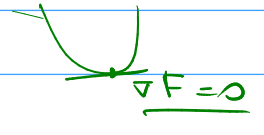
Stochastic optimization

$$F^{(k)}(\bar{w}) = \sum_{d \in \text{Batch}_k} f(\bar{w}, d)$$

$$F(\bar{w}) = \mathbb{E}_{q(y)} [f(\bar{w}, y)] \xrightarrow{\bar{w}} \min$$

$$F(\bar{w}) = \mathbb{E}_{\text{Unif}(D)} [f(\bar{w}, d)] \xrightarrow{\bar{w}} \min$$

$$= \frac{1}{N} \sum f(\bar{w}, d)$$



$$\hat{F}(\bar{w}) = \frac{1}{m} \sum_{i=1}^m f(\bar{w}, d_i), \text{ ye } d_i \sim \text{Unif}(D)$$

SGD

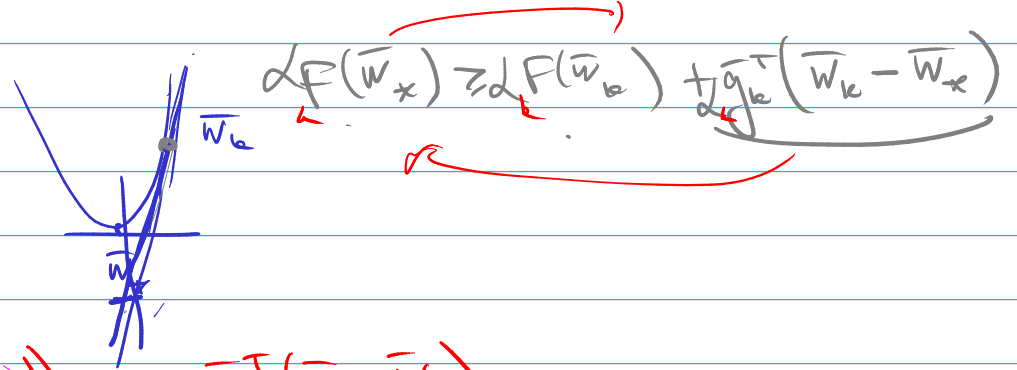
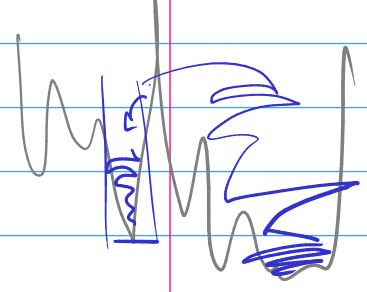
$$F(\bar{w}) = \mathbb{E}_q [f(\bar{w}, y)]$$

$$\bar{w}_{k+1} = \bar{w}_k - \alpha_k \hat{g}_k, \quad \mathbb{E}[\hat{g}_k] = \bar{g}_k = \nabla F(\bar{w}_k)$$

$$\mathbb{E}_q [\|\bar{w}_{k+1} - \bar{w}_*\|^2] = \mathbb{E} [\|\bar{w}_k - \alpha_k \hat{g}_k - \bar{w}_*\|^2] =$$

$$= \|\bar{w}_k - \bar{w}_*\|^2 - 2\alpha_k \hat{g}_k^T (\bar{w}_k - \bar{w}_*) + \alpha_k^2 \|\hat{g}_k\|^2$$

$$\xrightarrow{\mathbb{E}} \|\bar{w}_k - \bar{w}_*\|^2 - 2\alpha_k \bar{g}_k^T (\bar{w}_k - \bar{w}_*) + \alpha_k^2 \mathbb{E} [\|\hat{g}_k\|^2]$$



$$\alpha_k (F(\bar{w}_k) - F(\bar{w}_*)) \leq \alpha_k \bar{g}_k^T (\bar{w}_k - \bar{w}_*) =$$

$$= \frac{1}{2} \|\bar{w}_k - \bar{w}_*\|^2 + \frac{1}{2} \alpha_k^2 \mathbb{E} [\|\hat{g}_k\|^2] - \frac{1}{2} \mathbb{E} [\|\bar{w}_{k+1} - \bar{w}_*\|^2]$$

$$\underbrace{\tau_0}_{\tau_0} \rightarrow \tau_1 + \tau_2 - \tau_2 + \dots + \tau_k - \tau_{k+1}$$

$$\sum_{i=0}^k d_i (E[F(\bar{w}_i)] - F(\bar{w}_*)) \leq \frac{1}{2} \|\bar{w}_0 - \bar{w}_*\|^2 + \frac{1}{2} \sum_{i=0}^k d_i^2 E[\|\hat{g}_i\|^2] - \frac{1}{2} E[\|\bar{w}_{k+1} - \bar{w}_*\|^2]$$

$$E\left[\frac{\sum_{i=0}^k d_i F(\bar{w}_i)}{\sum d_i}\right] \geq E\left[F\left(\frac{\sum d_i F(\bar{w}_i)}{\sum d_i}\right)\right]$$

$$E\left[F\left(\frac{\sum d_i F(\bar{w}_i)}{\sum d_i}\right) - F(\bar{w}_*)\right] \leq \frac{\frac{1}{2} \|\bar{w}_0 - \bar{w}_*\|^2 + \frac{1}{2} \sum_{i=0}^k d_i^2 E[\|\hat{g}_i\|^2]}{\sum_{i=0}^k d_i}$$

$$\begin{cases} \|\bar{w}_0 - \bar{w}_*\| \leq R \\ E[\|\hat{g}_i\|^2] \leq G^2 \end{cases}$$

$$E[F(\bar{w}_0) - F(\bar{w}_*)] \leq \frac{R^2 + G^2 \sum_{i=0}^k d_i^2}{2 \sum_{i=0}^k d_i} \xrightarrow{k \rightarrow \infty} 0$$

$\sum_{i=0}^k d_i^2 = h^{2 \cdot (k+1)}$
 $\sum_{i=0}^k d_i = h(k+1)$

$d_i = h$

$$E[F(\bar{w}_0) - F(\bar{w}_*)] \leq \frac{R^2}{2h(k+1)} + \frac{1}{2} G^2 h \xrightarrow{k \rightarrow \infty} \frac{1}{2} G^2 h$$

$\sum_{i=0}^{\infty} d_i = \infty$
 $\sum_{i=0}^{\infty} d_i^2 < \infty$

$$d_i = \frac{1}{i}$$

