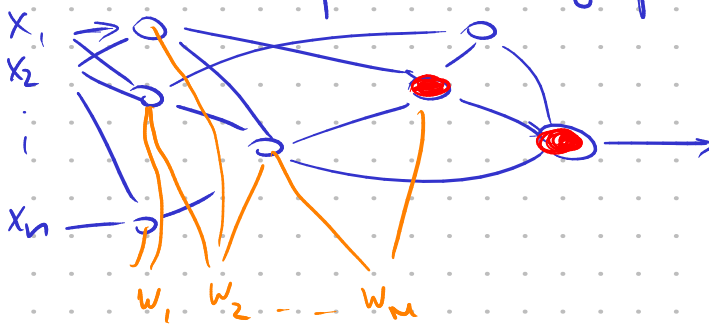


Computational graph

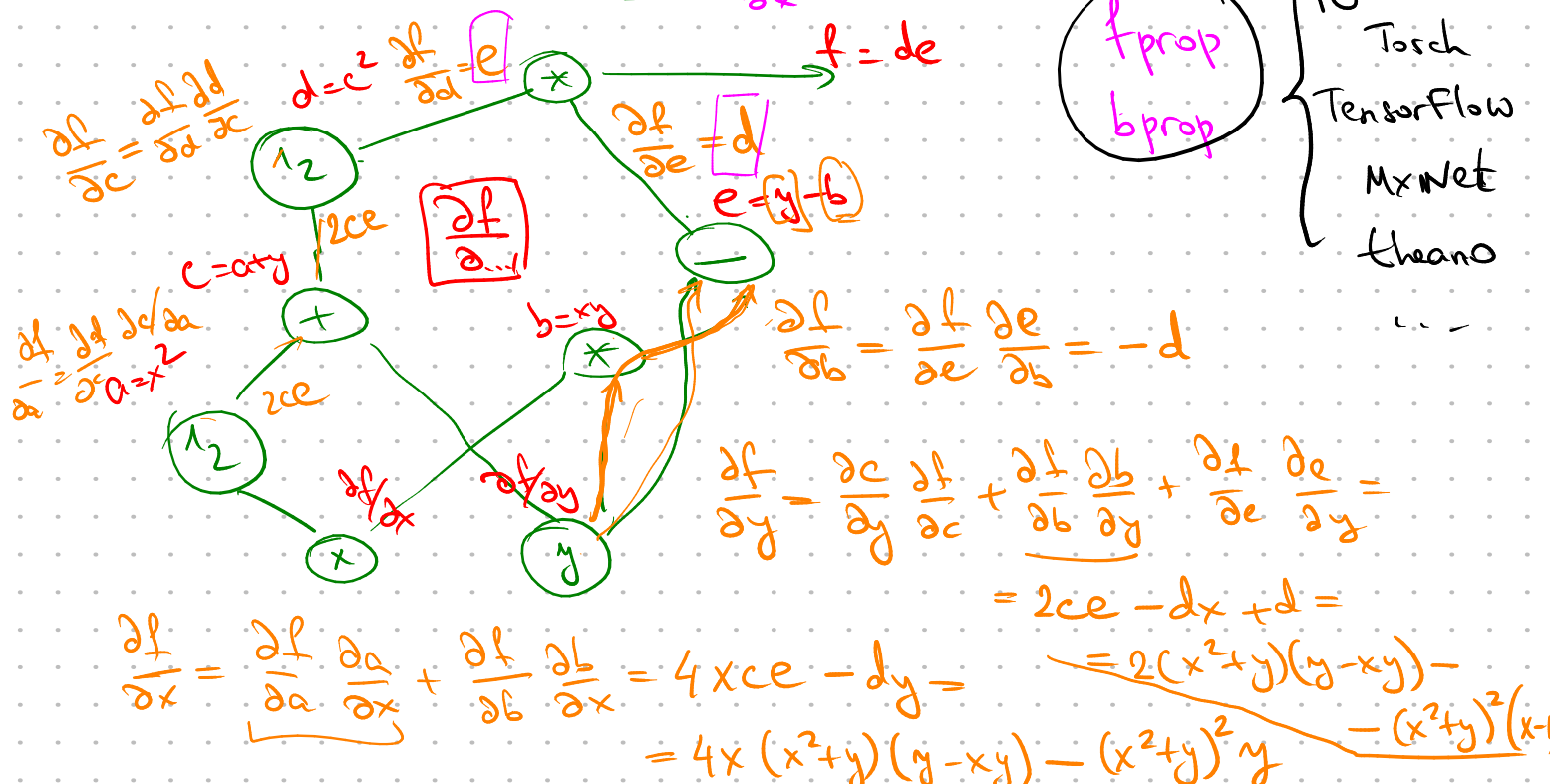
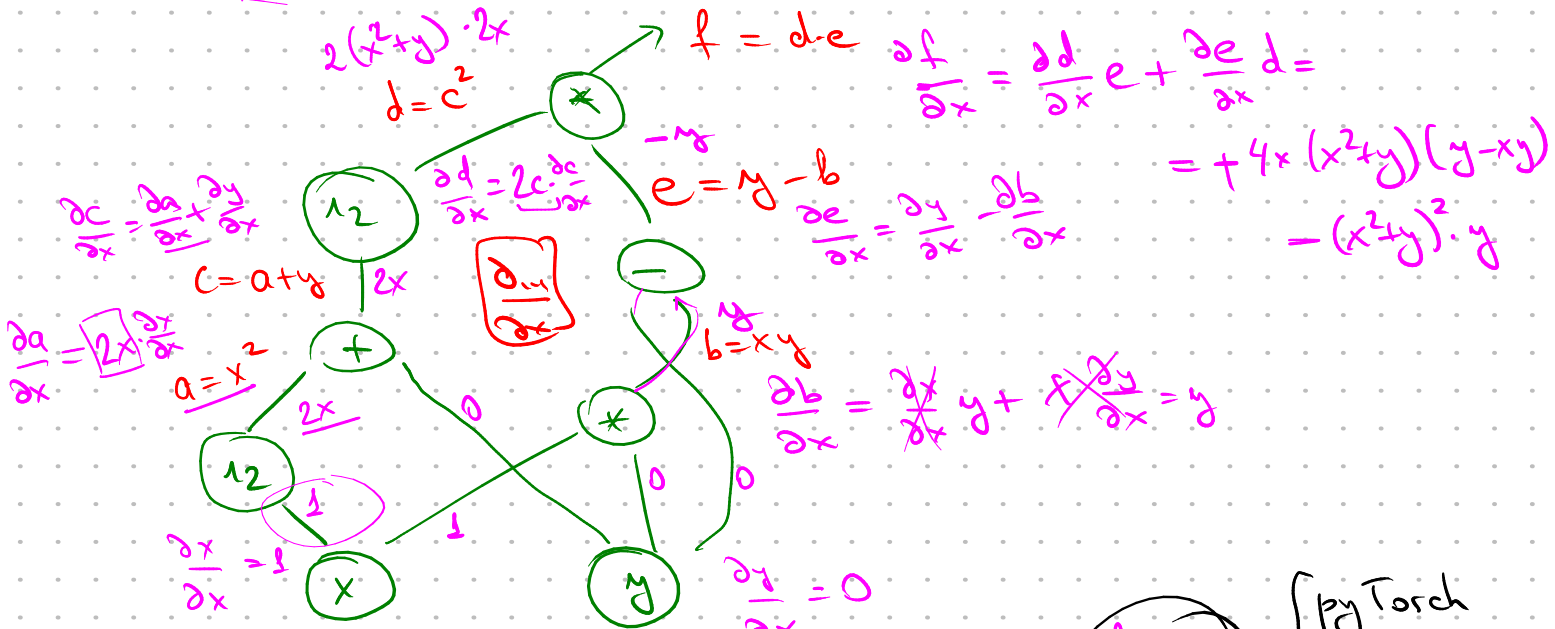


$L \rightarrow \min$
 $\log p(\mathcal{D}|\theta) \rightarrow \max$

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

$$f = (x^2 + y)^2 \cdot (y - xy)$$

$\frac{\partial f}{\partial x}$ $\frac{\partial f}{\partial y}$



$\left. \begin{matrix} \text{fprop} \\ \text{bprop} \end{matrix} \right\} \begin{matrix} \text{pyTorch} \\ \text{Torch} \\ \text{TensorFlow} \\ \text{MxNet} \\ \text{theano} \end{matrix}$

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

$$L = \frac{1}{N} \sum_{n=1}^N L(y_n, F(\bar{x}_n))$$

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

quasi-Newton methods

L-BFGS

$q(y) = \text{Unif}(D)$

Stochastic gradient descent

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

$$G(\bar{x}) = \nabla_{\bar{x}} F(\bar{x})$$

mini-batch

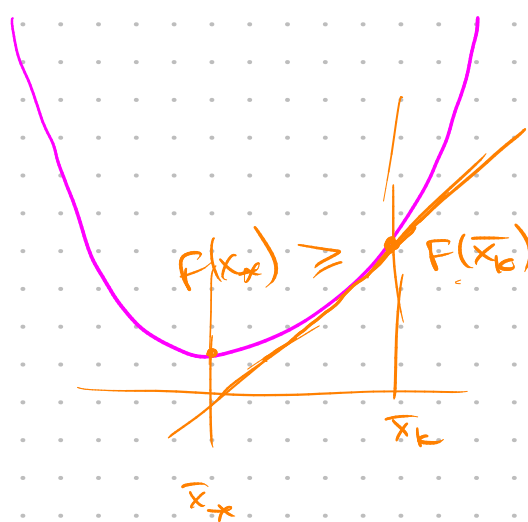
$$\hat{F}(\bar{x}) = \frac{1}{m} \sum_{i=1}^m f(\bar{x}, \bar{y}_i), \quad \hat{g}(\bar{x}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\bar{x}} f(\bar{x}, \bar{y}_i)$$

$$\bar{x}_{k+1} = \bar{x}_k - \alpha_k \hat{g}_k$$

$$g_k = \mathbb{E} \hat{g}_k = \nabla F(\bar{x}_k)$$

$$\begin{aligned} \|\bar{x}_{k+1} - \bar{x}_*\|^2 &= \|\bar{x}_k - \alpha_k \hat{g}_k - \bar{x}_*\|^2 = \\ &= \|\bar{x}_k - \bar{x}_*\|^2 - 2\alpha_k \hat{g}_k^T (\bar{x}_k - \bar{x}_*) + \alpha_k^2 \|\hat{g}_k\|^2 \end{aligned}$$

$$\mathbb{E} \|\bar{x}_{k+1} - \bar{x}_*\|^2 = \mathbb{E} \|\bar{x}_k - \bar{x}_*\|^2 - 2\alpha_k \mathbb{E} \hat{g}_k^T (\bar{x}_k - \bar{x}_*) + \alpha_k^2 \mathbb{E} \|\hat{g}_k\|^2$$



$$\alpha_k F(\bar{x}_*) \geq F(\bar{x}_k) + \bar{g}_k^T (\bar{x}_* - \bar{x}_k)$$

$$\begin{aligned} \alpha_k \bar{g}_k^T (\bar{x}_k - \bar{x}_*) &\geq \\ &\geq \alpha_k (F(\bar{x}_k) - F(\bar{x}_*)) \end{aligned}$$

$$\mathbb{E} \|\bar{x}_{k+1} - \bar{x}_*\|^2 - \mathbb{E} \|\bar{x}_k - \bar{x}_*\|^2 \leq$$

$$\leq -\alpha_k (F(\bar{x}_k) - F(\bar{x}_*)) + \alpha_k^2 \mathbb{E} \|\hat{g}_k\|^2$$

$$\alpha_k (F(\bar{x}_k) - F(\bar{x}_*)) \leq \frac{1}{2} \mathbb{E} \|\bar{x}_k - \bar{x}_*\|^2 - \frac{1}{2} \mathbb{E} \|\bar{x}_{k+1} - \bar{x}_*\|^2 + \frac{1}{2} \alpha_k^2 \mathbb{E} \|\hat{g}_k\|^2$$

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

$$E F\left(\frac{\sum \alpha_i \bar{x}_i}{\sum \alpha_i}\right) \leq \frac{\sum \alpha_i F(\bar{x}_i) - \sum \alpha_i F(\bar{x}_*)}{\sum \alpha_i}$$

$$E F\left(\frac{\sum \alpha_i \bar{x}_i}{\sum \alpha_i}\right) - F(\bar{x}_*) \leq \frac{\sum \alpha_i (F(\bar{x}_i) - F(\bar{x}_*))}{\sum \alpha_i} \leq \frac{1}{\sum \alpha_i}$$

$$\|\bar{x}_0 - \bar{x}_*\| \leq R, \quad E \|\hat{g}_k\|^2 \leq G^2$$

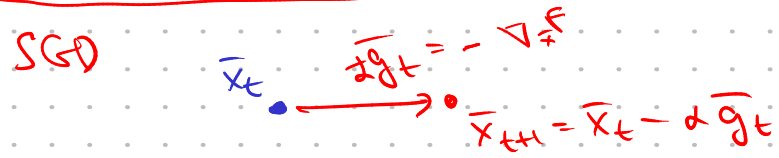
$$E F(\hat{\bar{x}}_k) - F(\bar{x}_*) \leq \frac{R^2 + G^2 \sum_{i=0}^k \alpha_i^2}{2 \sum_{i=0}^k \alpha_i} \leq \text{Konst}$$

[SGD] → ∞

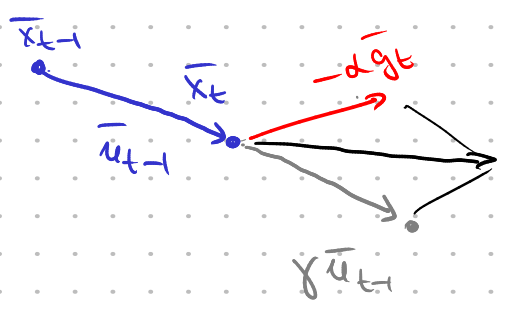
$\alpha_i = h$

$$\dots \leq \frac{R^2}{2h \cdot k} + \frac{G^2 h}{2}$$

$\xrightarrow[k \rightarrow \infty]{} 0$



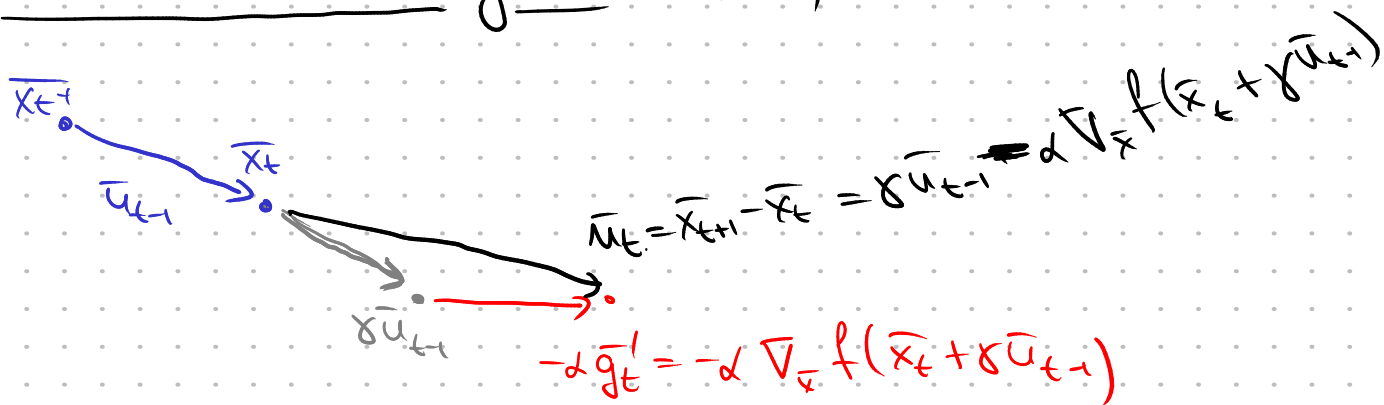
SGD with momentum



$$\bar{x}_{t+1} = \bar{x}_t + \gamma \bar{u}_{t-1} - \alpha \bar{g}_t$$

$\gamma = 0.9, 0.99, 0.995$

Nesterov accelerated gradient (NAG)



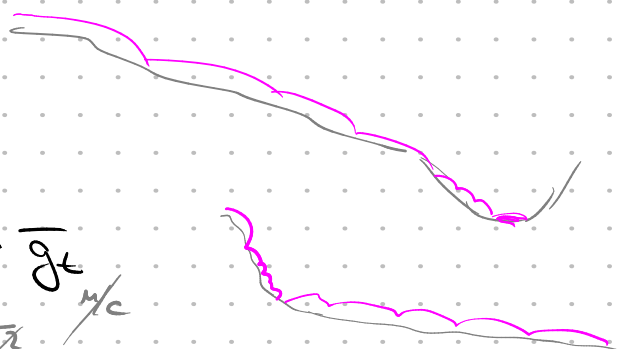
Adaptive SGD

Adagrad

$$G_0 = 0$$

$$G_{t,i} = G_{t-1,i} + \frac{g_{t,i}^2}{m/c}$$

$$\bar{x}_{t+1} = \bar{x}_t - \frac{\alpha}{\sqrt{G_{t+1}}} \cdot \frac{g_t}{m/c}$$



RMSprop

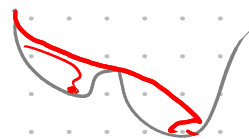
$$G_{t,i} = \gamma G_{t-1,i} + (1-\gamma) g_{t,i}^2$$

Adadelta

$$\bar{x}_{t+1} = \dots \frac{1}{(m/c)^2} \cdot \frac{g_t}{m/c}$$

$$\bar{x}_{t+1} = \bar{x}_t - \alpha \cdot \frac{\sqrt{R_{t+1}}}{\sqrt{G_{t+1}}} \cdot \frac{g_t}{(m/c)^2}$$

$$R_{t,i} = \beta R_{t-1,i} + (1-\beta) u_{t,i}^2$$



Adam

$$\bar{x}_{t+1} = \bar{x}_t - \frac{\alpha}{\sqrt{G_{t+1}}} \cdot \bar{m}_t$$

$$\beta_1 = 0.9$$

Nadan

AdamW

$$G_{t,i} = \beta_2 G_{t-1,i} + (1-\beta_2) g_{t,i}^2$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

$$m_{t,i} = \beta_1 m_{t-1,i} + (1-\beta_1) g_{t,i}$$