

$$-\log p(\bar{x} | \bar{w}) - \log p(\bar{w})$$

$$L(\bar{w}) = \sum_{\bar{x} \in D} [\log p(\bar{x} | \bar{w}) + \log p(\bar{w})]$$

$$= N \cdot \mathbb{E}_{\text{unif}(D)} [\log p(\bar{x} | \bar{w}) + \log p(\bar{w})]$$

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

$$\bar{w}^{(0)}, \bar{w}^{(k+1)} = \bar{w}^{(k)} - \eta \cdot \nabla_{\bar{w}} F(\bar{w}^{(k)}) \quad \text{GD - gradient descent}$$

$$F(\bar{w}) \approx \bar{w}^{(k)} + (\bar{w} - \bar{w}^{(k)}) \cdot \bar{g}_k + \frac{1}{2} (\bar{w} - \bar{w}^{(k)})^T H_k (\bar{w} - \bar{w}^{(k)})$$

$$-\bar{g}_k + H_k \cdot \bar{w} - H_k \cdot \bar{w}^{(k)} = 0$$

$$\bar{w} = \bar{w}^{(k)} + H_k^{-1} \cdot \bar{g}_k \quad \left(\frac{\partial F}{\partial w_i} \frac{\partial w_j}{\partial w_j} \right)$$

quasi-Newton algorithms L-BFGS $\bar{g}_k \rightarrow \text{low rank} \approx H^{-1}$

SGD - stochastic gradient descent

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

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

$$\begin{cases} q(\bar{y}) = \text{unif}(D) \\ q(\bar{y}) = p(\bar{y} | D) \end{cases}$$

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

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

ye $\bar{y}_i \sim q(\bar{y})$

$$\mathbb{E} \hat{g}_k = \bar{g}_k$$

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

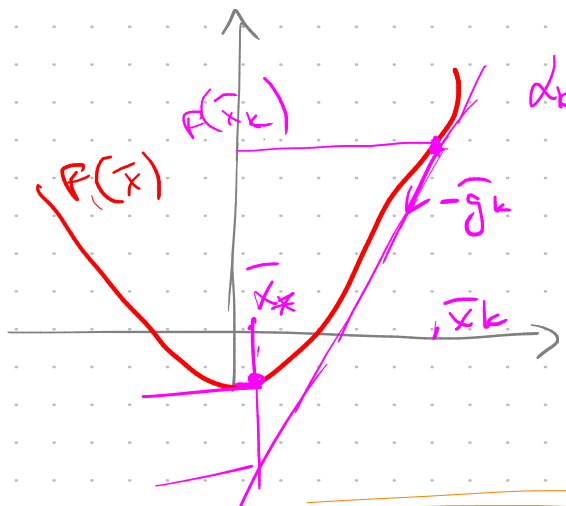
$$\text{GD} \left[F(\bar{x}) \rightarrow \min \quad \bar{g}_k \quad F(\bar{x}_k - \alpha \cdot \bar{g}_k) \rightarrow \min \quad (?) \right]$$

\bar{x}_* residue

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

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

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



$$\alpha_k \cdot F(\bar{x}_*) \approx \alpha_k (F(\bar{x}_k) - (\bar{x}_* - \bar{x}_k)^T \hat{g}_k)$$

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

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

$$P\left(\frac{\sum \alpha_i \bar{x}_i}{\sum \alpha_i}\right) \leq \frac{\sum \alpha_i F(\bar{x}_i)}{\sum \alpha_i}$$

$$F(\bar{x}_k) - F(\bar{x}_*) \leq \sum \alpha_i (F(\bar{x}_i) - F(\bar{x}_*))$$

$$\sum \alpha_i \cdot (E F(\bar{x}_k) - F(\bar{x}_*)) \leq E \sum \alpha_i (F(\bar{x}_i) - F(\bar{x}_*)) \leq$$

$$\leq \frac{1}{2} \left(\dots \right)$$

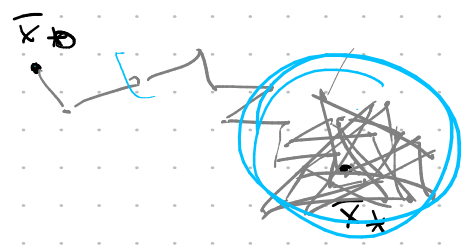
$$E F(\bar{x}_k) - F(\bar{x}_*) \leq \frac{E \|\bar{x}_0 - \bar{x}_*\|^2 + \sum_{i=0}^k \alpha_i^2 E \|\hat{g}_i\|^2}{2 - \sum_{i=0}^k \alpha_i}$$

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

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

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

$$\alpha_i = \alpha \Rightarrow \frac{R^2 + G^2 \cdot \alpha^2 \cdot k}{2 - \alpha k} = \frac{R^2}{2\alpha k} + \frac{G^2 \cdot \alpha}{2}$$



- Adam

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

$$x_{t+1,i} = x_{t,i} - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} \cdot m_{t,i}$$

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

$$\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-7.8}$$

Weight decay

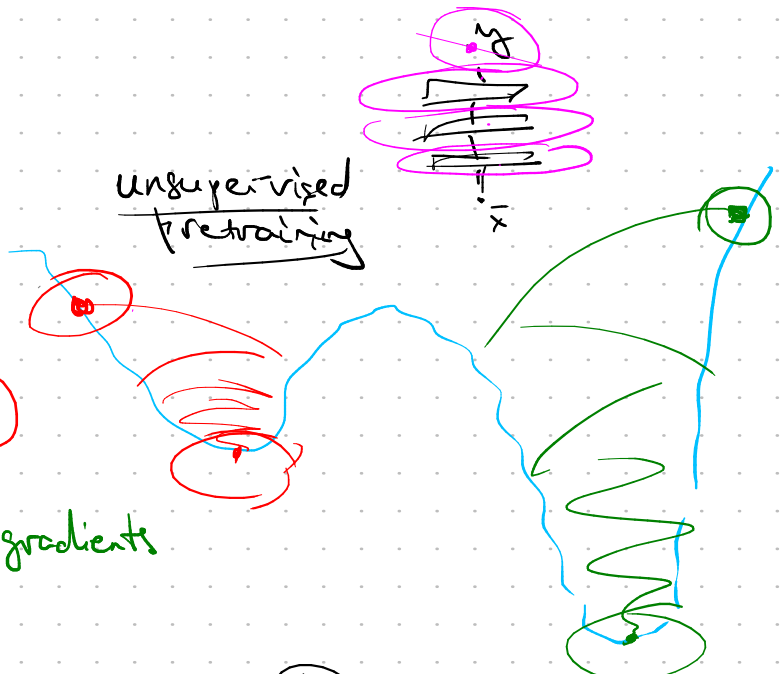
$$L_2: L(\bar{w}) + \lambda \|\bar{w}\|^2$$

Weight decay:

$$\bar{w} := \left(\frac{1-\eta}{1-\eta}\right) \cdot \bar{w} - \alpha \cdot \nabla_{\bar{w}} L(\bar{w}) = \bar{w} - \alpha \cdot \nabla_{\bar{w}} (L(\bar{w}) + \frac{\lambda}{2\alpha} \|\bar{w}\|^2)$$



unsupervised pretraining



Weight initialization

$$w^{(k)} = w^{(k-1)} - \eta \nabla_{w^{(k-1)}} L(w^{(k-1)})$$

exploding
vanishing gradients

$$y = \sum_{i=1}^n w_i x_i = \sum_{i=1}^n w_i \cdot x_i$$

$$\text{Var}[y_i] = \text{Var}[w_i x_i] =$$

$$= \mathbb{E}[w_i^2 x_i^2] - \mathbb{E}[w_i x_i]^2 = \mathbb{E}[x_i^2] \cdot \text{Var}[w_i] + \mathbb{E}[w_i^2] \cdot \text{Var}[x_i] + \text{Var}[w_i] \cdot \text{Var}[x_i]$$

$$\text{Var}[y_i] = \text{Var}[w_i] \cdot \text{Var}[x_i]$$

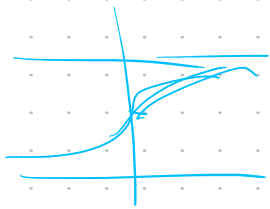
$$\text{Var}[y_i] = n \cdot \text{Var}[w_i] \cdot \text{Var}[x_i] \approx 1 = O(1/n)$$

Xavier init

$$w_i \sim \text{Unif}\left[-\frac{\sqrt{3}}{\sqrt{n}}, \frac{\sqrt{3}}{\sqrt{n}}\right]$$

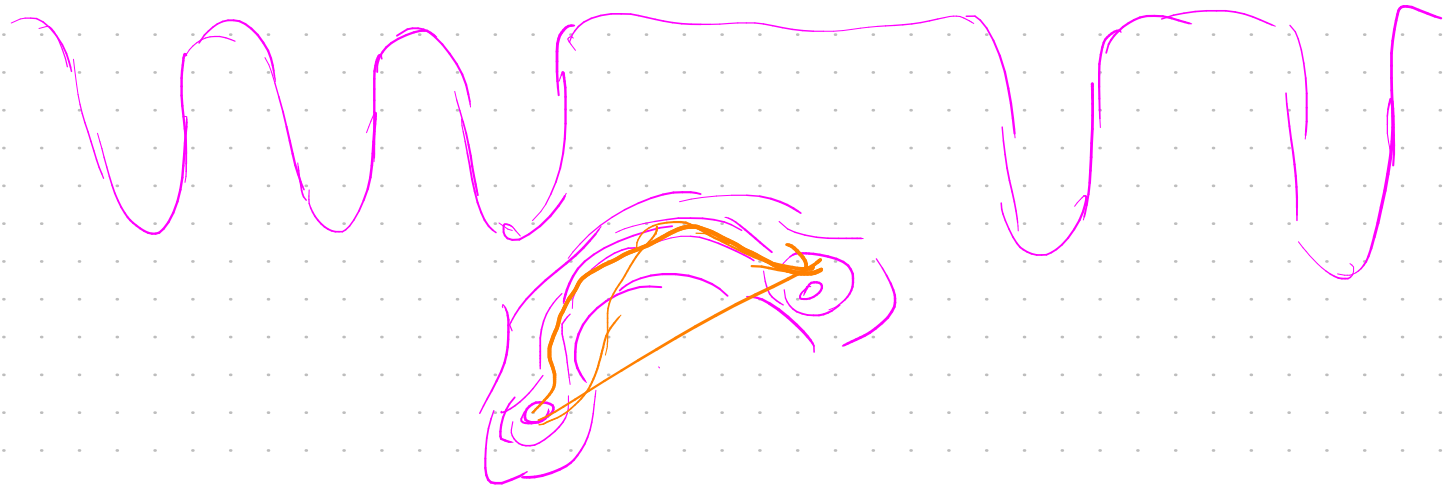
$$\text{Var}(\text{Unif}(a,b)) = \frac{(b-a)^2}{12}$$

$$\text{Var}(\text{Unif}[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}]) = \frac{(2/\sqrt{n})^2}{12} = \frac{1}{3n}$$



$$\text{Var}(w_i) = \frac{2}{n} \quad \text{ReLU}$$

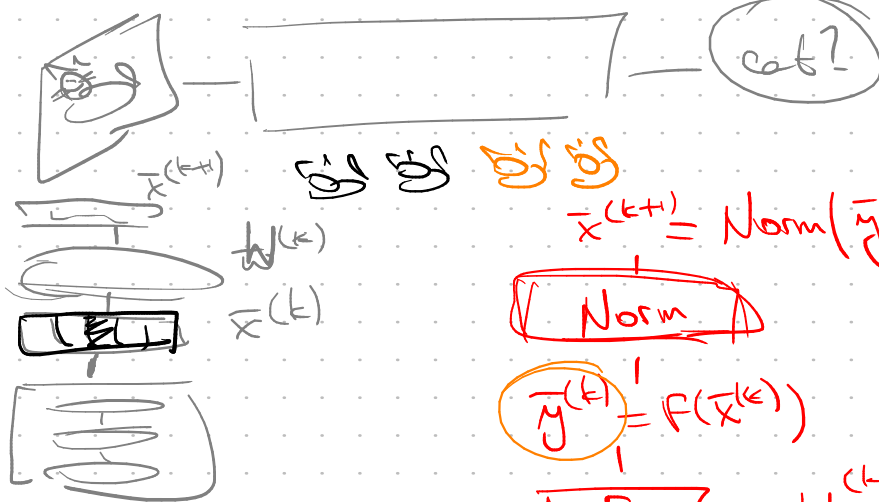
He init



Batch normalization

Covariate shift

Internal covariate shift

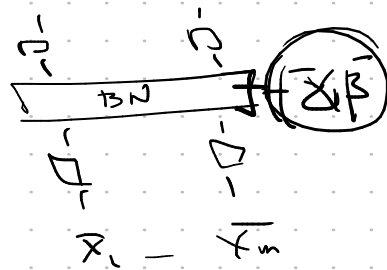


$$\bar{x}^{(k+1)} = \text{Norm}(\bar{y}^{(k)})$$

$$\bar{y}^{(k)} = F(\bar{x}^{(k)})$$

Layer normalization

$$\bar{x}^{(k+1)} = \frac{\bar{y}^{(k)} - E_n[\bar{y}^{(k)}]}{\sqrt{\text{Var}_n[\bar{y}^{(k)}]}}$$



$$\text{BN}(\bar{y}^{(k)}) = \frac{\bar{y}^{(k)} - E_n[\bar{y}^{(k)}]}{\sqrt{\text{Var}_n[\bar{y}^{(k)}]}} \odot \bar{\gamma} + \bar{\beta}$$

