

Simon OH Today 10:30A

Optimization Methods for Deep Learning

W

Gradient descent for non-convex optimization

Decsent Lemma: Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be twice differentiable, and $\|\nabla^2 f\|_2 \leq \beta$. Then setting the learning rate $\eta = 1/\beta$, and applying gradient descent, $x_{t+1} = x_t - \eta \nabla f(x_t)$, we have:

$$f(x_t) - f(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla f(x_t)\|_2^2.$$

$f(x_t) \downarrow$
 \Downarrow
 $\curvearrowleft \curvearrowright$

$$X: \nabla f(x) = 0$$

Converging to stationary points

$$\beta = \frac{1}{\beta}$$

ϵ - approximate stationary points

Theorem: In $T = O(\frac{\beta}{\epsilon^2})$ iterations, we have $\|\nabla f(x)\|_2 \leq \epsilon$.

$$Pf: f(x_{t+1}) \leq f(x_t) - \frac{\gamma}{2} \|\nabla f(x_t)\|_2^2$$

Sum over $t = 0, 1, \dots, T$

$$\sum_{t=1}^T f(x_t) \leq \sum_{t=0}^{T-1} f(x_t) - \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow f(x_T) \leq f(x_0) - \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \leq f(x_0) - f(x_T) \leq f(x_0) - \min_x f(x)$$

$$\Rightarrow T \cdot \frac{\gamma}{2} \min_{0 \leq t \leq T-1} \|\nabla f(x_t)\|_2^2 \leq f(x_0) - \min_x f(x) \leq \epsilon$$

$$\Rightarrow \min_{0 \leq t \leq T-1} \|\nabla f(x_t)\|_2 = \sqrt{\frac{2\beta(f(x_0) - \min_x f(x))}{T}} = \frac{\sqrt{2\beta}(f(x_0) - \min_x f(x))}{\sqrt{T}} = \frac{\sqrt{2\beta}}{\sqrt{T}} \epsilon$$

Gradient Descent for Quadratic Functions

Opt: $x=0$

Symmetric, full-rank, $\lambda_{\min}(A) > 0$

Problem: $\min_x \frac{1}{2} x^T A x$ with $A \in \mathbb{R}^{d \times d}$ being positive-definite.

Theorem: Let λ_{\max} and λ_{\min} be the largest and the smallest eigenvalues of A . If we set $\eta \leq \frac{1}{\lambda_{\max}}$, we have

$$\|x_t\|_2 \leq (1 - \eta \lambda_{\min})^t \|x_0\|_2$$

$$\|x_{t+1}\|_2 = \|x_t - \eta A x_t\|_2$$

$$= \|(I - \eta A)x_t\|_2$$

$$\leq \|I - \eta A\|_2 \cdot \|x_t\|_2$$

$$= (1 - \eta \lambda_{\max}) \cdot \|x_t\|_2$$

$$\leq (1 - \eta \lambda_{\max})^t \|x_0\|_2$$

to make $\rho \cdot \|x_t\|_2 \leq \epsilon$

$$\eta = \frac{1}{\lambda_{\max}}$$

$$\Rightarrow O\left(\frac{\lambda_{\max}}{\lambda_{\min}} \log\left(\frac{1}{\epsilon}\right)\right)$$

$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$: condition number

can be generalized for non-convex

e.g. $A = \begin{pmatrix} 10 & 0 \\ 0 & 1 \end{pmatrix}$

Momentum: Heavy-Ball Method (Polyak '64)

$$x_{t+1} = x_t - \gamma_f \cdot \nabla f(x_t)$$

Problem: $\min_x f(x)$

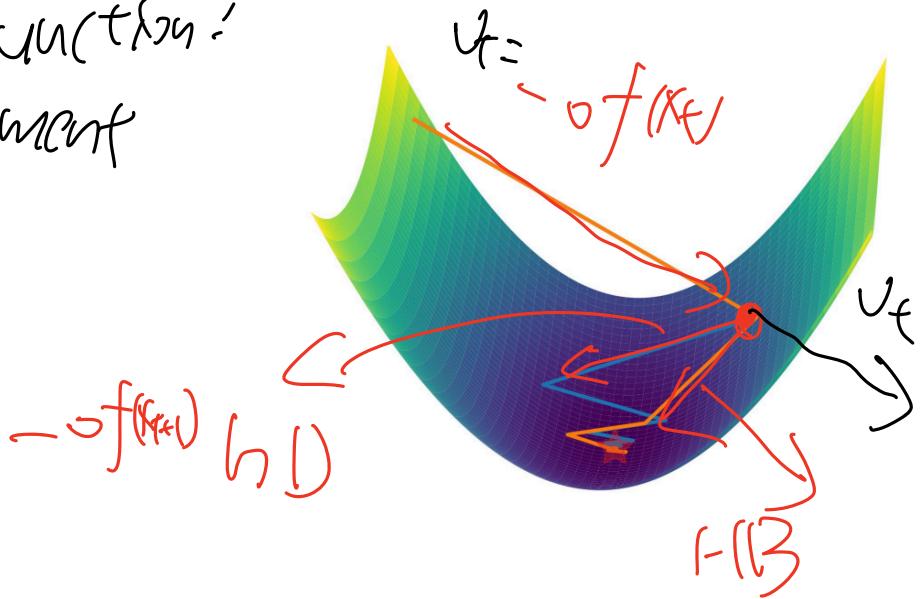
Method: $v_{t+1} = -\nabla f(x_t) + \beta v_t$
 $x_{t+1} = x_t + \eta v_{t+1}$

For quadratic function:
provable improvement

$$\mathcal{O}\left(\sqrt{K} \cdot \log\left(\frac{1}{\epsilon}\right)\right)$$

$$\text{vs. } \mathcal{O}(K \cdot \log\left(\frac{1}{\epsilon}\right))$$

• slow not work
for strongly convex



Momentum: Nesterov Acceleration (Nesterov '89)

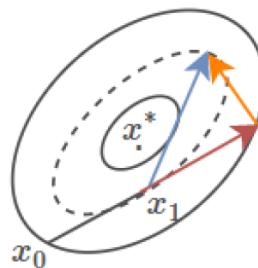
Problem: $\min_x f(x) - \nabla f(x_t)^T$ for $H\beta$

Method: $v_{t+1} = -\nabla f(x_t + \beta v_t) + \beta v_t$
 $x_{t+1} = x_t + \eta v_{t+1}$ *lookhead*

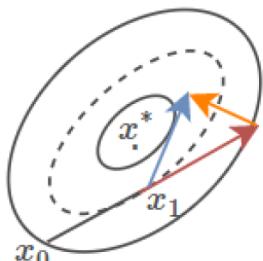
For general strongly convex functions: $O(\sqrt{\frac{L}{\alpha}})$

- use ODE

Polyak's Momentum



Nesterov Momentum



Newton's Method

2nd order method
 $\mathcal{O}(d^3)$

Newton's Method: $x_{t+1} = x_t - \eta (\nabla^2 f(x_t))^{-1} \nabla f(x_t)$

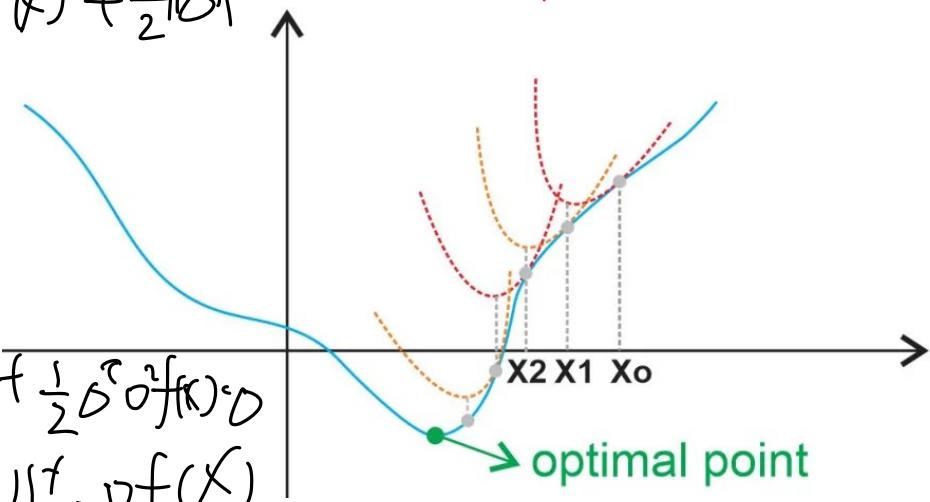
- G1: $x_{t+1} = x_t - \eta \cdot \nabla f(x_t)$
- $f(x+\delta) \approx f(x) + \delta^T \nabla f(x) + \frac{1}{2} \|\delta\|^2$ → isotropic quadratic
- $\Delta \leftarrow -\nabla f(x)$

• Newton:

$$f(x+\delta) \approx f(x) + \delta^T \nabla f(x) + \frac{1}{2} \delta^T \nabla^2 f(x) \delta$$

$$\Rightarrow \Delta \leftarrow -(\nabla^2 f(x))^{-1} \cdot \nabla f(x)$$

$$T = \mathcal{O}(\log \log (\frac{1}{\epsilon}))$$



AdaGrad (Duchi et al. '11)

- (1) diagonal
- (2) build inverse Hessian
in an online

Newton Method: $x_{t+1} = x_t - \eta(\nabla^2 f(x_t))^{-1} \nabla f(x_t)$

AdaGrad: separate learning rate for every parameter

initialization

$$x_{t+1} = x_t - \eta(G_t + \epsilon I)^{-1} \nabla f(x_t), \quad (G_t)_{ii} = \sqrt{\sum_{j=1}^{t-1} (\nabla f(x_t)_i)^2}$$

\$G_t\$ is always diagonal

- effective learning rate: $\eta \cdot (G_t + \epsilon I)^+$ small
- Default hyper parameters work well

RMSProp (Hinton et al. '12)

root mean Squared Propagability

AdaGrad: separate learning rate for every parameter

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1} \nabla f(x_t), (G_t)_{ii} = \sqrt{\sum_{j=1}^{t-1} (\nabla f(x_t)_i)^2}$$

RMSProp: exponential weighting of gradient norms

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1/2} \nabla f(x_t),$$
$$(G_{t+1})_{ii} = \beta(G_t)_{ii} + (1 - \beta)(\nabla f(x_t)_i)^2$$

$0 < \beta < 1$

$\text{if: } \beta^t \leq t$

$\beta^{t-1} \rightarrow 0$

AdaDelta (Zeiler '12)

RMSProp:

$$x_{t+1} = x_t - \frac{\eta(G_{t+1} + \epsilon I)^{-1/2}}{(G_{t+1})_{ii} = \beta(G_t)_{ii} + (1 - \beta)(\nabla f(x_t)_i)^2} \nabla f(x_t),$$

AdaDelta:

$$x_{t+1} = x_t - \eta \Delta x_t, \quad \text{unit of } x$$
$$\Delta x_t = \sqrt{u_t + \epsilon} \cdot (G_{t+1} + \epsilon I)^{-1/2} \nabla f(x_t) \quad \frac{\partial f}{\partial x} / \left(\frac{\partial f}{\partial x} \right) \propto \text{unit of } f$$
$$(G_{t+1})_{ii} = \rho(G_t)_{ii} + (1 - \rho)(\nabla f(x_t)_i)^2, \quad \text{unit}^2 \text{ of } x$$
$$u_{t+1} = \rho u_t + (1 - \rho) \|\Delta x_t\|_2^2 \quad \frac{1}{\|\Delta x_t\|_2}$$

GD : $\Delta x \propto \frac{\partial f}{\partial x} \propto \text{unit of } x$

Newton : $\Delta x \propto \frac{\partial f}{\partial x} / \left(\frac{\partial^2 f}{\partial x^2} \right) \propto \text{unit of } x$

Adam (Kingma & Ba '14)

Momentum:

$$v_{t+1} = -\nabla f(x_t) + \beta v_t, x_{t+1} = x_t + \eta v_{t+1}$$

RMSProp: exponential weighting of gradient norms

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1/2} \nabla f(x_t),$$

$$(G_t)_{ii} = \beta(G_t)_{ii} + (1 - \beta)(\nabla f(x_t)_i)^2$$

Adam:

$$v_{t+1} = \underbrace{\beta_1 v_t + (1 - \beta_1) \nabla f(x_t)}_{\text{Momentum}}$$

$$(G_{t+1})_{ii} = \beta_2(G_t)_{ii} + (1 - \beta_2)(\nabla f(x_t)_i)^2$$

$$x_{t+1} = x_t - \eta(G_{t+1} + \epsilon I)^{-1/2} v_{t+1}$$

Default choice nowadays.

∃ convex problem s.t. Adam does not converge

Other Optimizers

- AdamW
- NAdam
- RAdam
- GGT
- K-FAC
- ...

adaptive

Are these actually useful

Heavy-tail

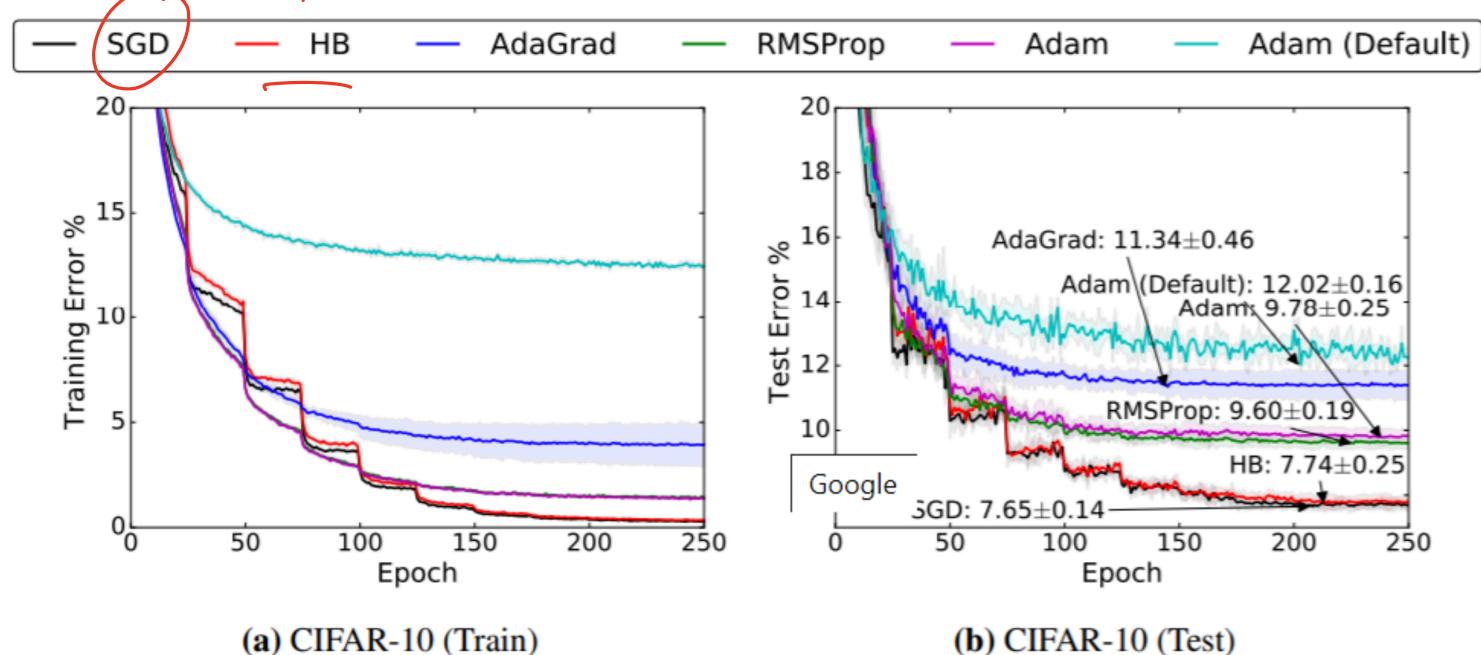


Figure 1: Training (left) and top-1 test error (right) on CIFAR-10. The annotations indicate where the best performance is attained for each method. The shading represents \pm one standard deviation computed across five runs from random initial starting points. In all cases, adaptive methods are performing worse on both train and test than non-adaptive methods.

Wilson, Roelofs, Stern, Srebro, Recht '18

Important Techniques in Neural Network Training

W

Gradient Explosion / Vanishing

$f: \mathbb{R}^d \rightarrow \mathbb{C}$

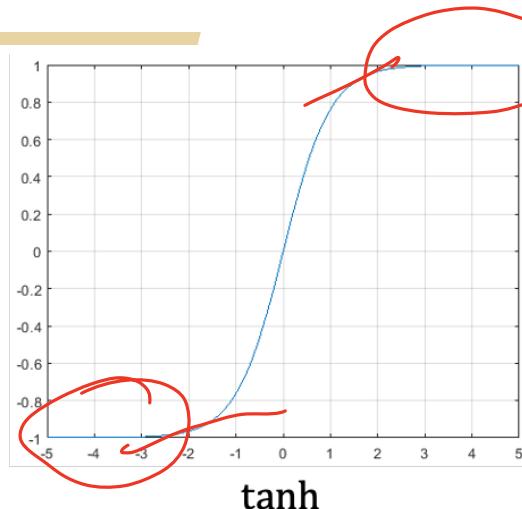
- Deeper networks are harder to train:
 - Intuition: gradients are products over layers
 - Hard to control the learning rate

$$f(x, w_1, \dots, w_{H+1}) = w_{H+1} \sigma(w_H \dots \sigma(w_1 x) \dots)$$

$$\frac{\partial f}{\partial w_n} \leq (w_{H+1} A_H \dots w_{n+1} A_n)^T (A_{n-1} w_n \dots w_1 x)$$

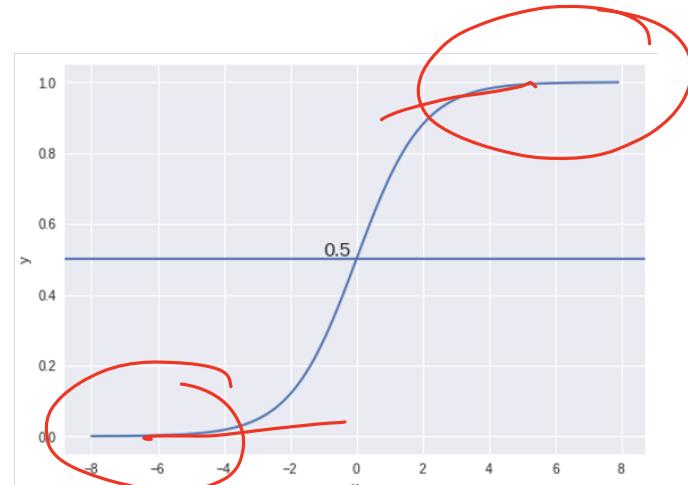
- magnitude $\sqrt{\quad}$ gradient exp large
- magnitude \downarrow / subspans don't align
gradient exp small

Activation Functions



gradient
vanishing

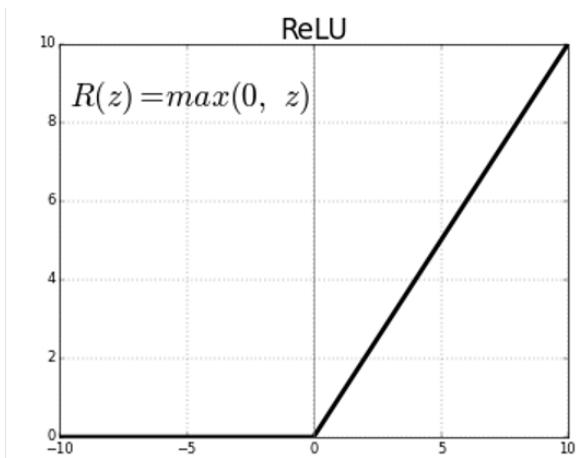
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



sigmoid (σ)

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

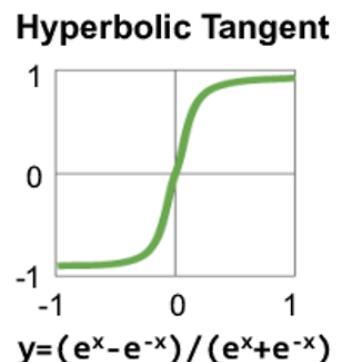
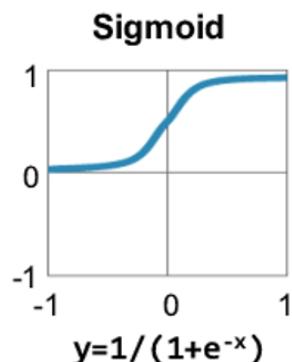
$f'(z) = 0$ or /



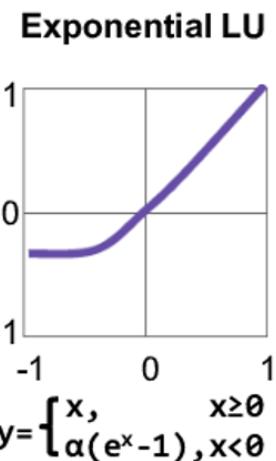
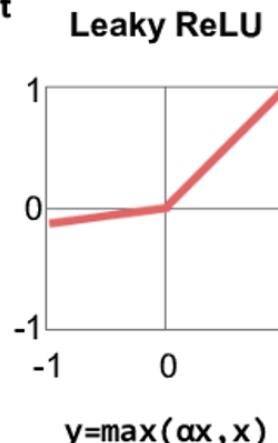
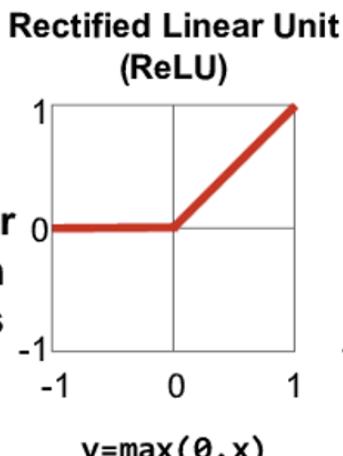
Rectified Linear United

Activation Function

**Traditional
Non-Linear
Activation
Functions**



**Modern
Non-Linear
Activation
Functions**

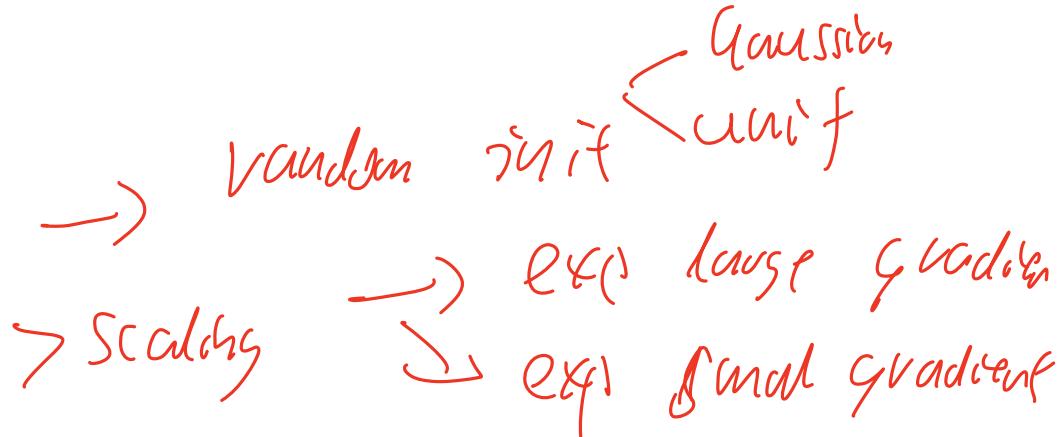


α = small const. (e.g. 0.1)

$\alpha \approx 0.1$

Initialization

- Zero-initialization
- Large initialization
- Small initialization
- Design principles:
 - Zero activation mean
 - Activation variance remains same across layers

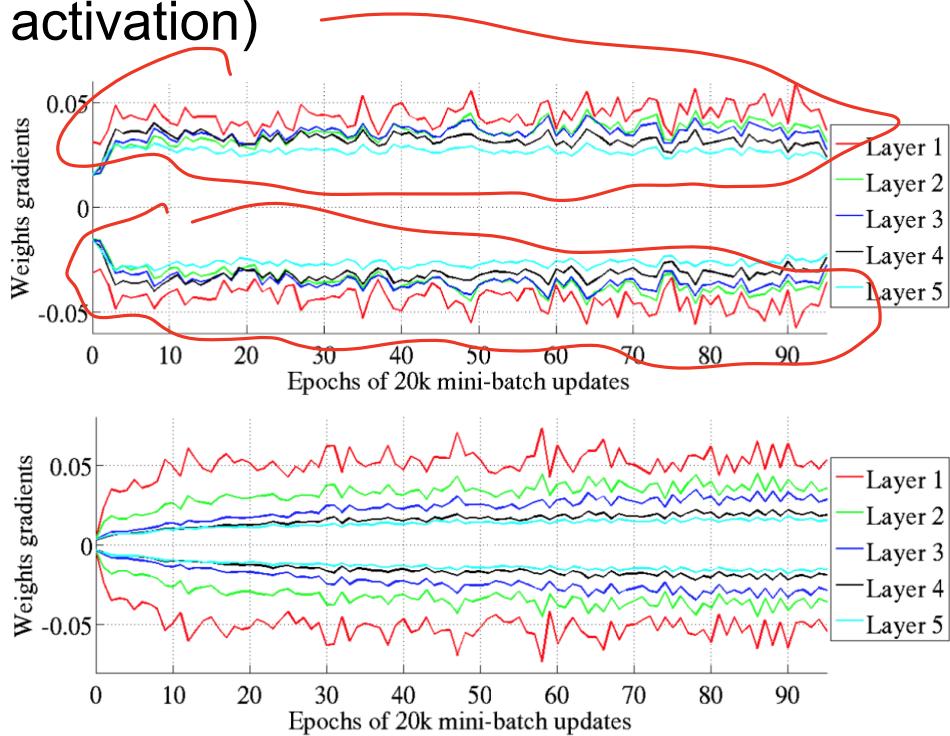


Xavier Initialization (Glorot & Bengio, '10)

- $W_{ij}^{(h)} \sim \text{Unif} \left[-\frac{\sqrt{6}}{\sqrt{d_h + d_{h+1}}}, \frac{\sqrt{6}}{\sqrt{d_h + d_{h+1}}} \right]$
- $b^{(h)} = 0$
- Experiments (tanh activation)

$$\begin{aligned} \text{Var}(W_{ij}^{(h)}) &= \frac{1}{12} \cdot \left(\frac{2\sqrt{6}}{d_{h+1} + d_h} \right)^2 \\ &= \frac{1}{d_{h+1} + d_h} \end{aligned}$$

$\text{Unif} \left[-\frac{1}{\sqrt{d_h + d_{h+1}}}, \frac{1}{\sqrt{d_h + d_{h+1}}} \right]$



$W^{(h)} \in \mathbb{R}^{d_{h+1} \times d_h}$

d_h : fan-in
 d_{h+1} : fan-out

Kaiming Initialization (He et al. '15)

- $W_{ij}^{(h)} \sim \mathcal{N}\left(0, \frac{2}{d_h}\right)$.
- $b^{(h)} = 0$
- Designed for ReLU activation
- 30-layer neural network

