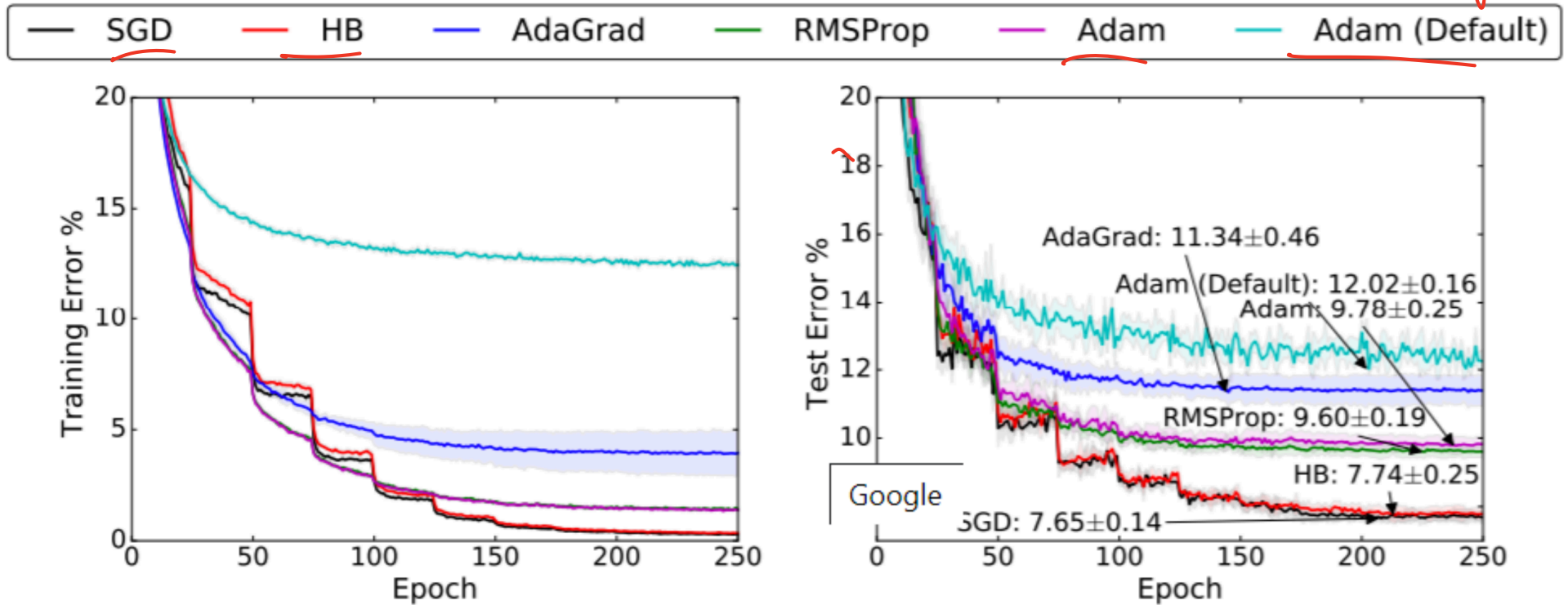


# Are these actually useful

hyper

↓



**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

---



# Gradient Explosion / Vanishing

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

$$f(x, W_1, \dots, W_{H+1}) = W_{H+1} \sigma(W_H \dots \sigma(W_1 x) \dots)$$

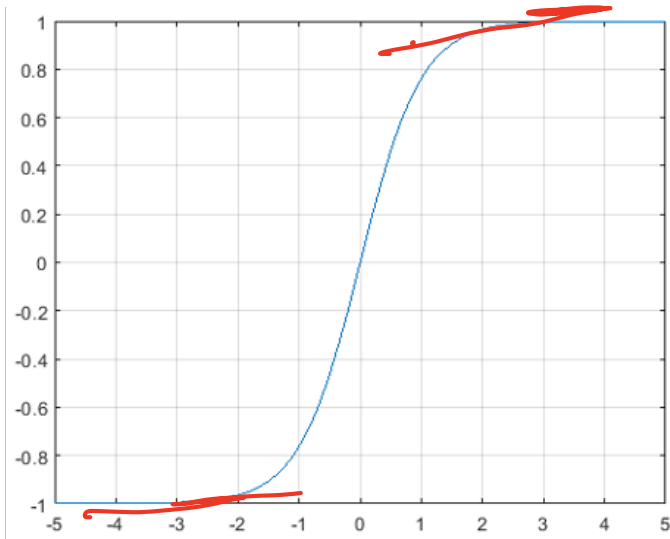
$$\frac{\partial J}{\partial W_n} = (W_{H+1} A_H \dots W_{n+1} A_n)^T (A_{n-1} W_{n-1} \dots W_1 x)$$

$$A_n = \text{diag}(\sigma'(W_n) \sigma(W_{n-1}) \dots \sigma(W_1 x) \dots)$$

$(W_n, A_n)$   $\cup$  magnitude small  $\rightarrow$  exp small  
large  $\rightarrow$  exp large

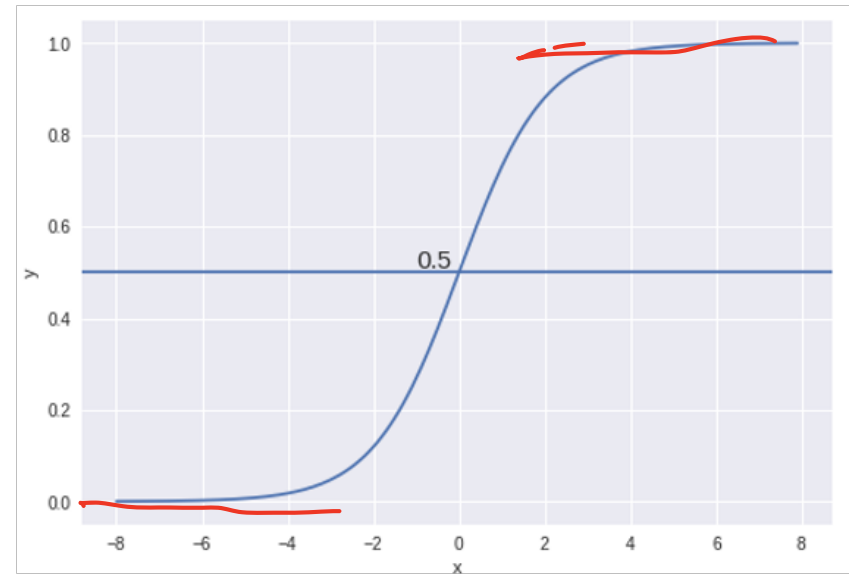
(2) sparse not align  $\rightarrow$  product exp small

# Activation Functions



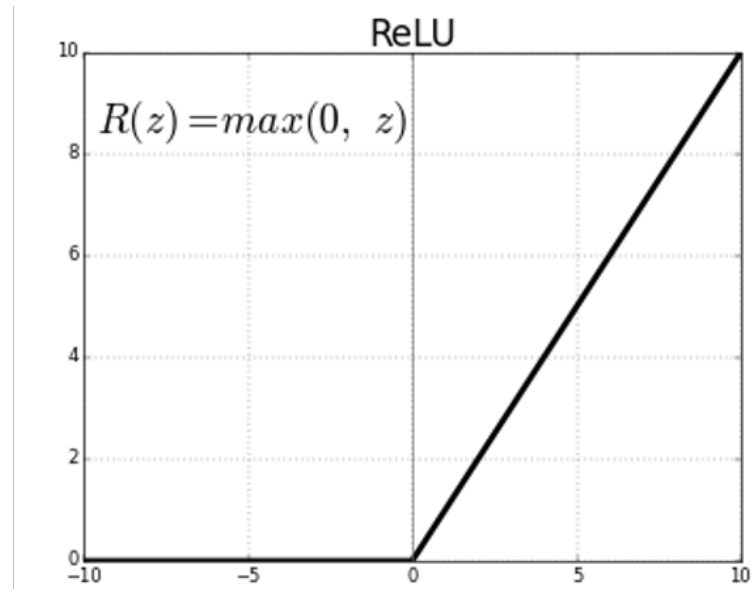
tanh

$$\sigma'(z) \rightarrow 0$$



sigmoid

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

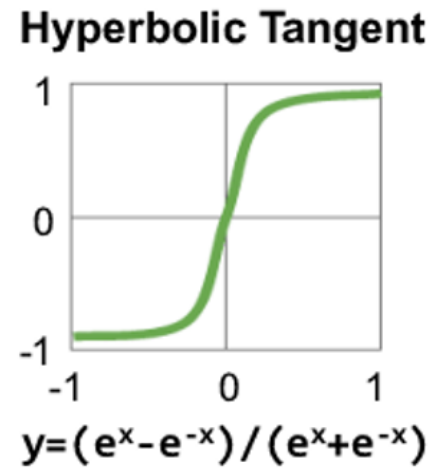
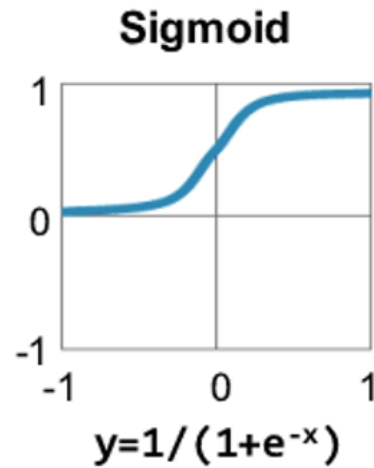


Rectified Linear Unit

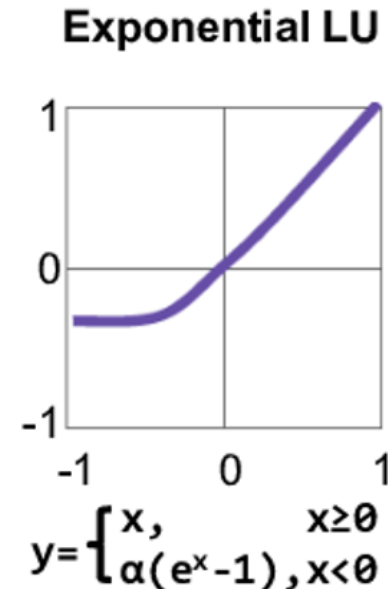
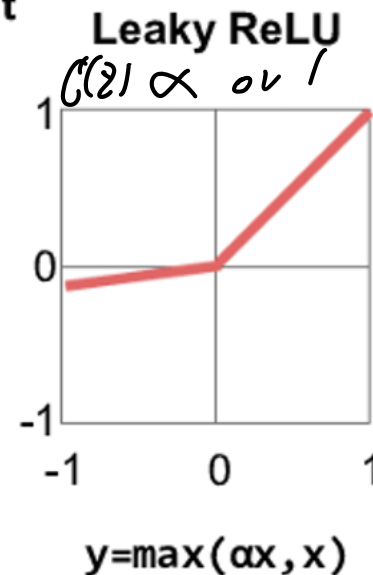
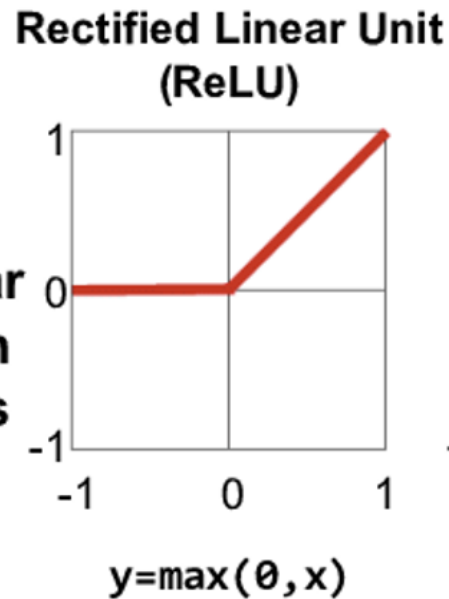
$$\sigma'(z) = 0 \text{ or } 1$$

# Activation Function

Traditional  
Non-Linear  
Activation  
Functions



Modern  
Non-Linear  
Activation  
Functions



$\alpha = \text{small const. (e.g. 0.1)}$

*avoid dying neurons*

conver:  $\|x_T - x^*\|_2 \leq e^{-T} \|x_0 - x^*\|$

# Initialization

- Zero-initialization  $\rightarrow$  all gradient  $\rightarrow 0$
- Large initialization  $\rightarrow$  scaling  $\rightarrow$  exp large
- Small initialization  $\rightarrow$  scaling  $\rightarrow$  exp small
- Design principles: random / distribution  $D_n$  for  $W_n$ 
  - Zero activation mean (no prior knowledge)
  - Activation variance remains same across layers

# Xavier Initialization (Glorot & Bengio, '10) $d_{in} \times d_{out}$

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

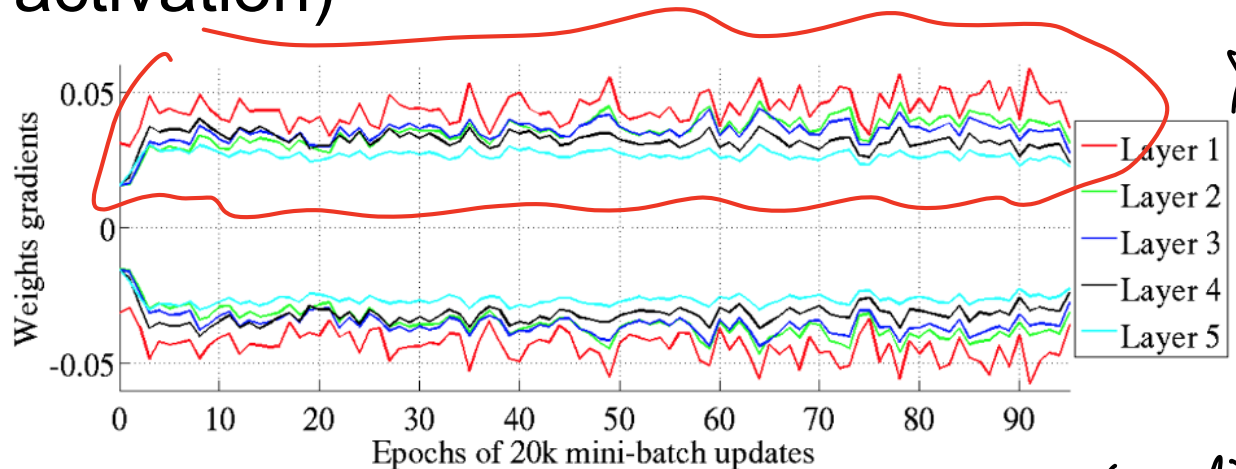
$W^{(h)} \in \mathbb{R}$   
 $d_h$ : fan-in  
 $d_{h+1}$ : fan-out

- $b^{(h)} = 0$

$\mathcal{N}(0, \frac{2}{d_{in} + d_{out}})$

- Experiments (tanh activation)

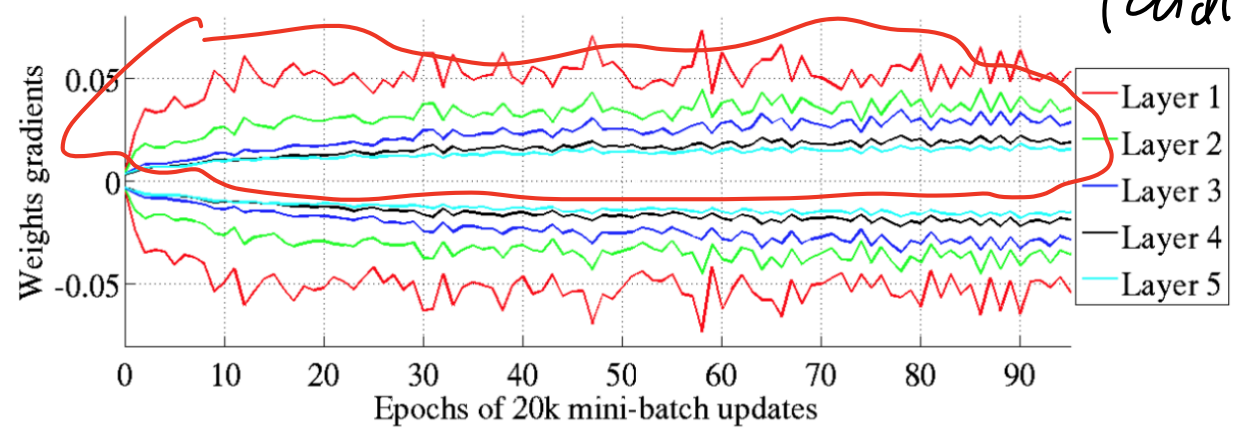
$\text{Var}(W_{ij}^{(h)})$   
 $= \frac{1}{12} \cdot \frac{(2\sqrt{6})^2}{d_{in} + d_{out}} = \frac{2}{d_{in} + d_{out}}$



Xavier init

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

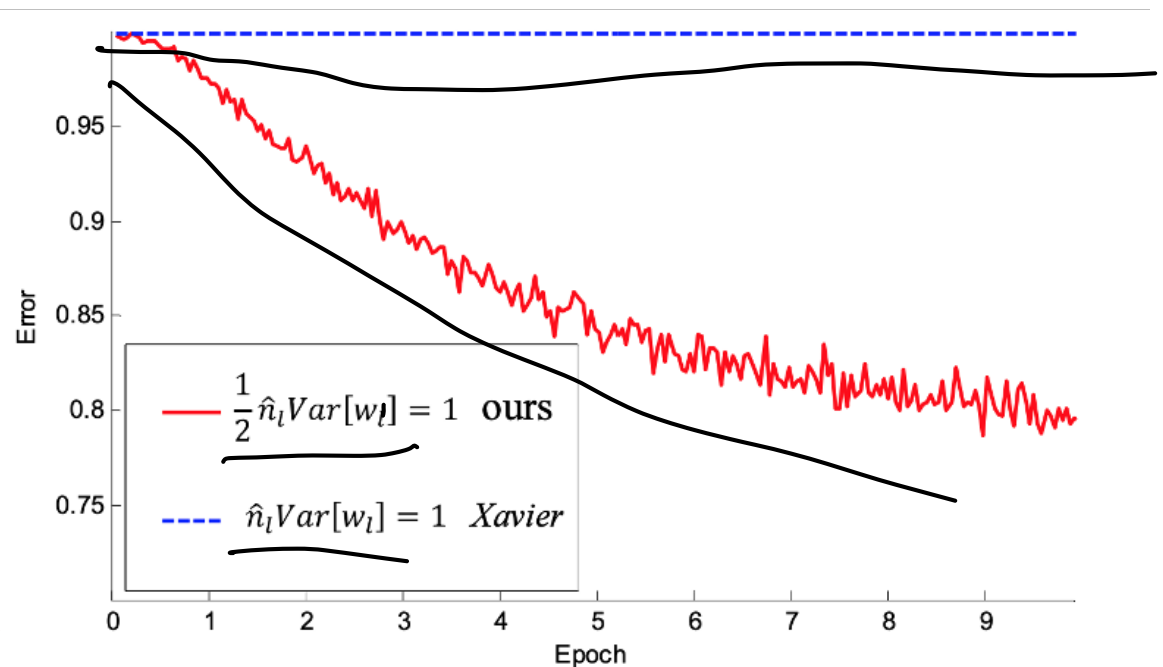
$\text{Var}(u_{ij}) = \frac{1}{3 d_{in}}$



tandem init

# 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





# Kaiming Initialization (He et al. '15)

Each  $z^h = W^h x^h \in \mathcal{R}^{d_{h+1}}$

$$x^{h+1} = \sigma(z^h)$$

$$z_i^h = \sum_{j=1}^{d_h} w_{ij}^h x_j^h$$

$$\mathbb{E}[w_{ij}^h] = 0 \Rightarrow \mathbb{E}[z_i^h] = 0$$

$$\text{Var}(z_i^h) = d_h \cdot \text{Var}(w_{ij}^h x_j^h)$$

$$= d_h (\text{Var}(w_{ij}^h) - \text{Var}(x_j^h))$$

$$+ \underbrace{(\mathbb{E} w_{ij}^h)^2}_{=0} \cdot \text{Var}(x_j^h) + \underbrace{\text{Var}(w_{ij}^h)}_{=0} \cdot (\mathbb{E}(x_j^h)^2)$$

$$= d_h \cdot \text{Var}(w_{ij}^h) - \mathbb{E}[(x_j^h)^2]$$

## Kaiming Initialization (He et al. '15)

$$\mathbb{E} \left[ (X_j^h)^2 \right] = \int_{-\infty}^{\infty} (x_j^h)^2 p(x_j^h) dx_j^h$$

$$= \int_{-\infty}^{\infty} \max(0, z_j^{h-1})^2 p(z_j^{h-1}) dz_j^{h-1}$$

(ReLU)

$$= \int_0^{\infty} (z_j^{h-1})^2 p(z_j^{h-1}) dz_j^{h-1}$$

(Symmetric)

$$= \frac{1}{2} \int_{-\infty}^{\infty} (z_j^{h-1})^2 p(z_j^{h-1}) dz_j^{h-1}$$

$$= \frac{1}{2} \text{var}(z_j^{h-1})$$

## Kaiming Initialization (He et al. '15)

We want  $\text{Var}(z_i^h) = \text{Var}(z_i^{h+1})$

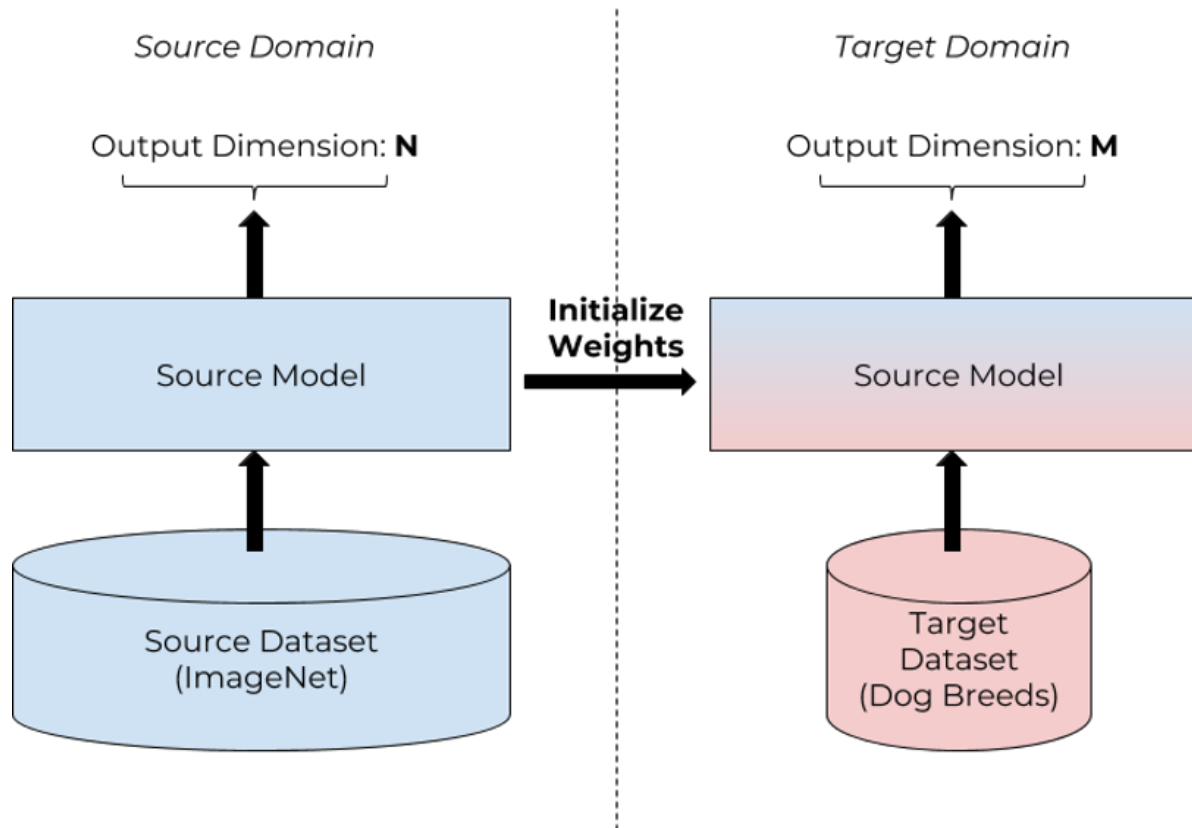
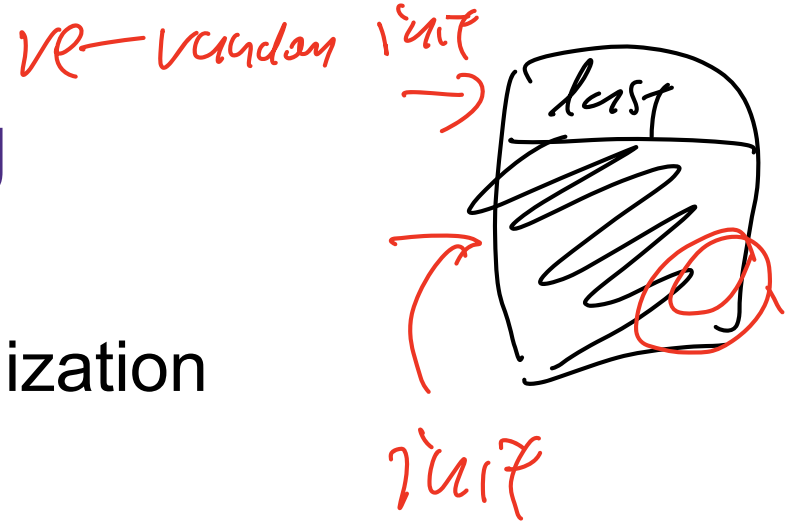
$$d_y \text{Var}(W_{ij}^y) \cdot \frac{1}{2} \text{Var}(z_i^{h+1}) = \text{Var}(z_i^{h+1})$$

$$\text{Var}(W_{ij}^y) = \frac{2}{d_y}$$

$$\text{Var}(z^{H+1}) = \text{Var}(z^1) \cdot \left( \prod_{h=1}^{H-1} \frac{d_y \text{Var}(W_{ij}^h)}{2} \right)$$

# Initialization by Pre-training

- Use a pre-trained network as initialization
- And then fine-tuning



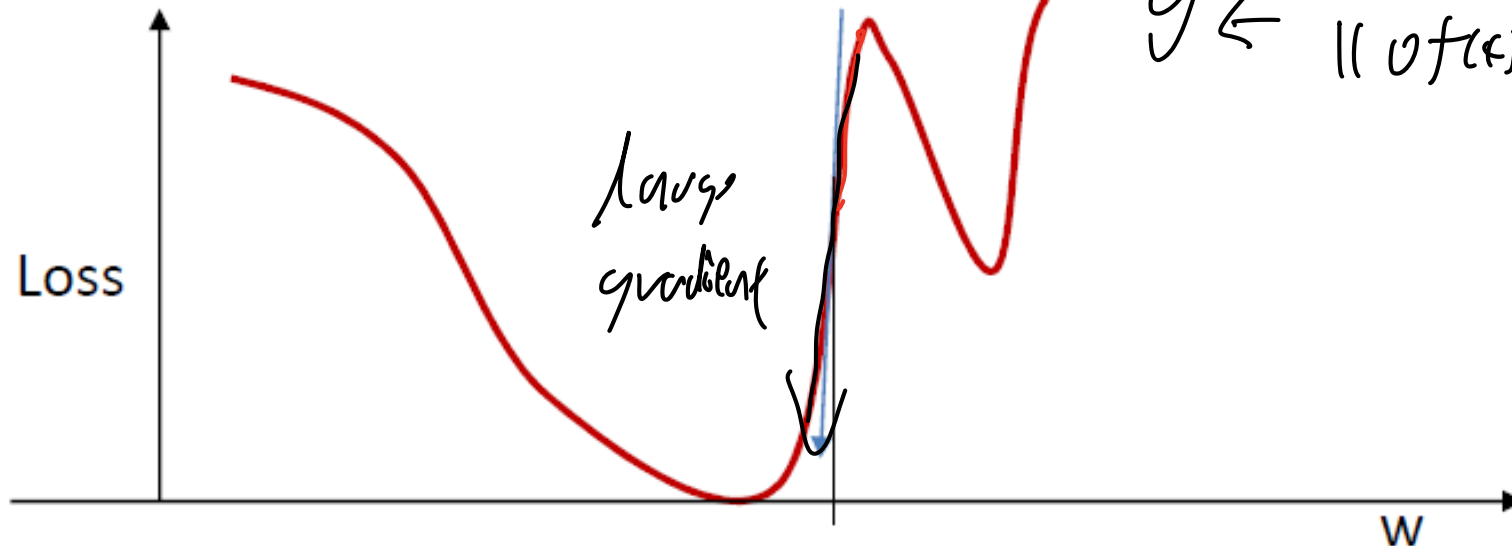
- (1) fix *vep* *paup*
- (2) fine ~~step~~

# Gradient Clipping

- The loss can occasionally lead to a steep descent
- This result in immediate instability
- If gradient norm bigger than a threshold, set the gradient to the threshold.

$$\frac{\partial f(x)}{\partial x}, \text{ if } \|\frac{\partial f(x)}{\partial x}\| \geq \text{threshold}$$

$$g \leftarrow \frac{\frac{\partial f(x)}{\partial x}}{\|\frac{\partial f(x)}{\partial x}\|} \cdot \text{threshold}$$



# Batch Normalization (Ioffe & Szegedy, '14)

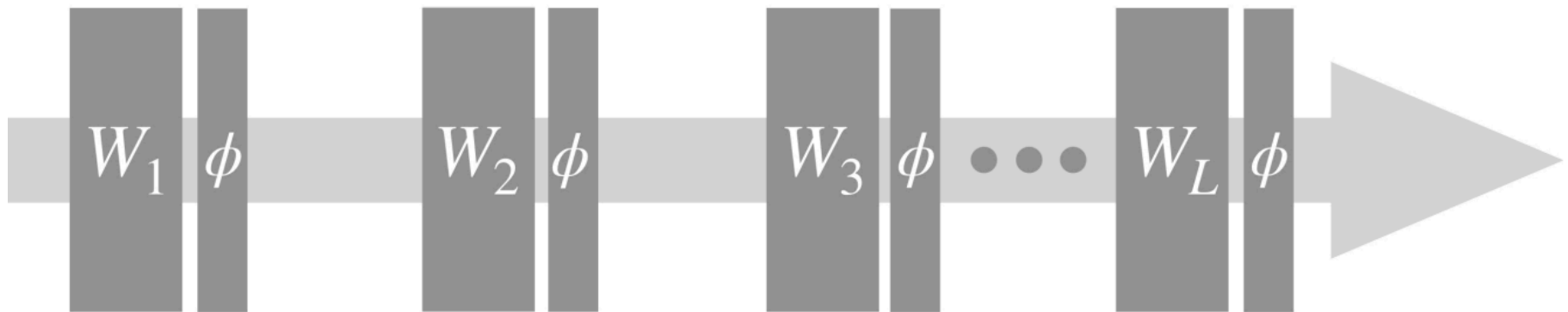
---

- **Normalizing/whitening** (mean = 0, variance = 1) the inputs is generally useful in machine learning.
  - Could normalization be useful at the level of hidden layers?
  - **Internal covariate shift**: the calculations of the neural networks change the distribution in hidden layers even if the inputs are normalized
- **Batch normalization** is an attempt to do that:
  - Each unit's **pre-activation** is normalized (mean subtraction, std division)
  - During training, mean and std is computed for each minibatch (can be backproped!)

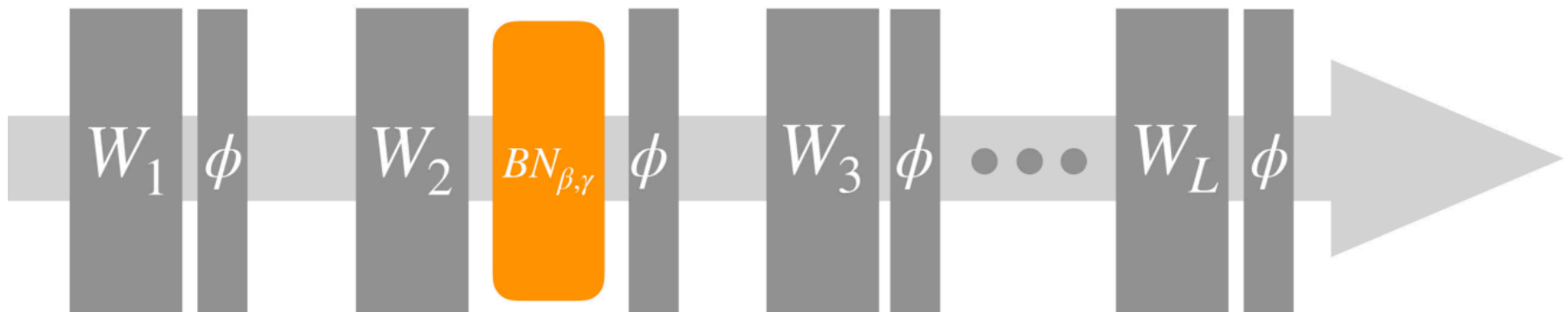
# Batch Normalization (Ioffe & Szegedy, '14)

$\phi$ : activation

## Standard Network



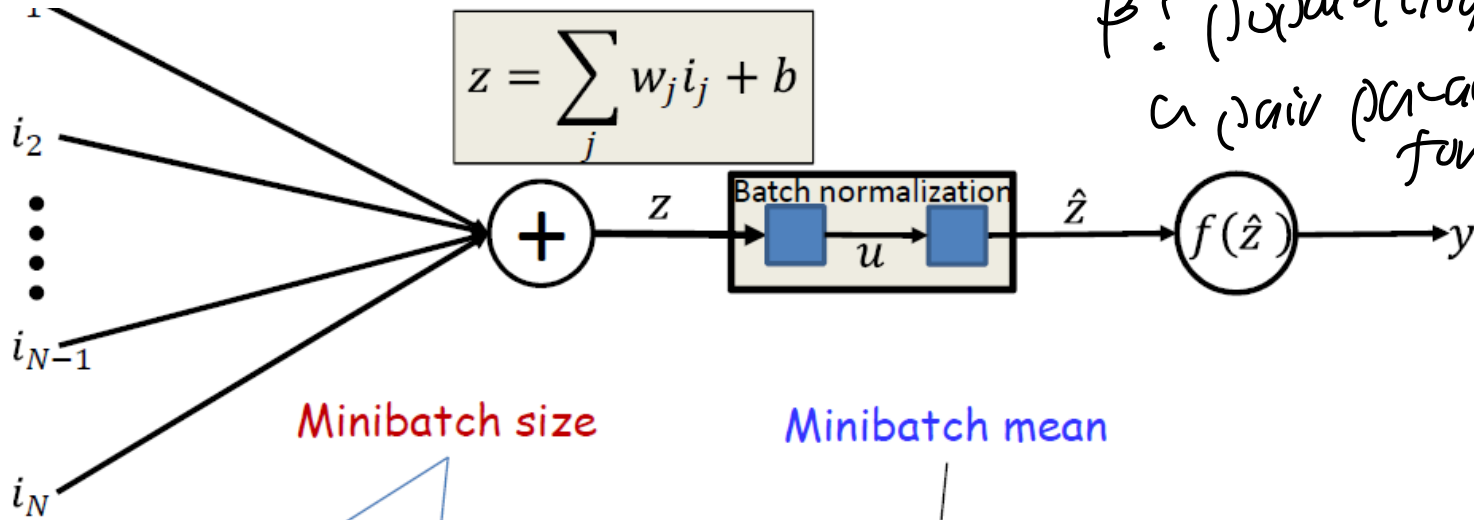
## Adding a BatchNorm layer (between weights and activation function)



BN

# Batch Normalization (Ioffe & Szegedy, '14)

$\gamma$ : population std  
 $\beta$ : population mean  
 a pair parameter forward BN layer



Minibatch size

Minibatch mean

Minibatch standard deviation

$$\mu_B = \frac{1}{B} \sum_{i=1}^B z_i$$

$$\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (z_i - \mu_B)^2$$

$$u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\hat{z}_i = \gamma u_i + \beta$$

rescale

prevent 0-div

$z_1, \dots, z_B$

$B$ : batch size

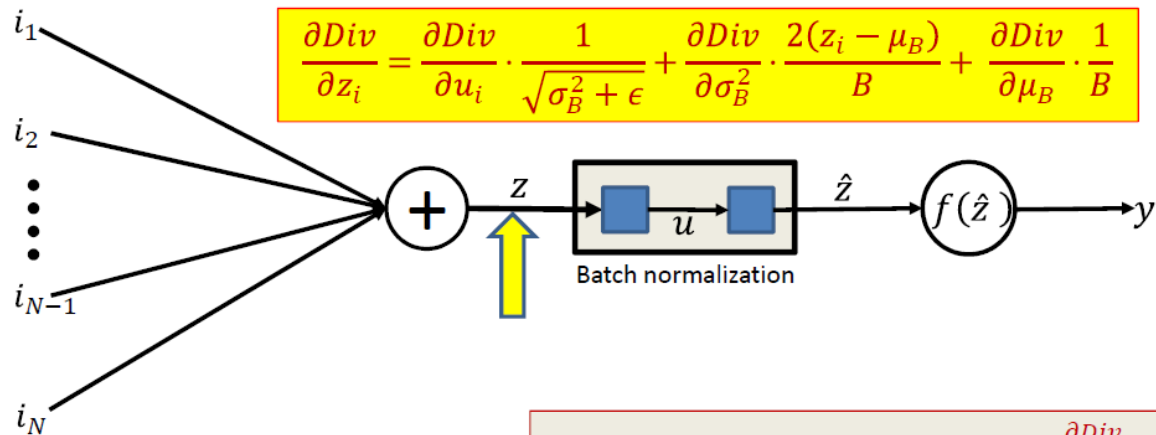


# Batch Normalization (Ioffe & Szegedy, '14)

- BatchNorm at training time
  - Standard backprop performed for each single training data
  - Now backprop is performed over entire batch.

$$\frac{\partial \text{Div}}{\partial \sigma_B^2} = \frac{-1}{2} (\sigma_B^2 + \epsilon)^{-3/2} \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i}$$

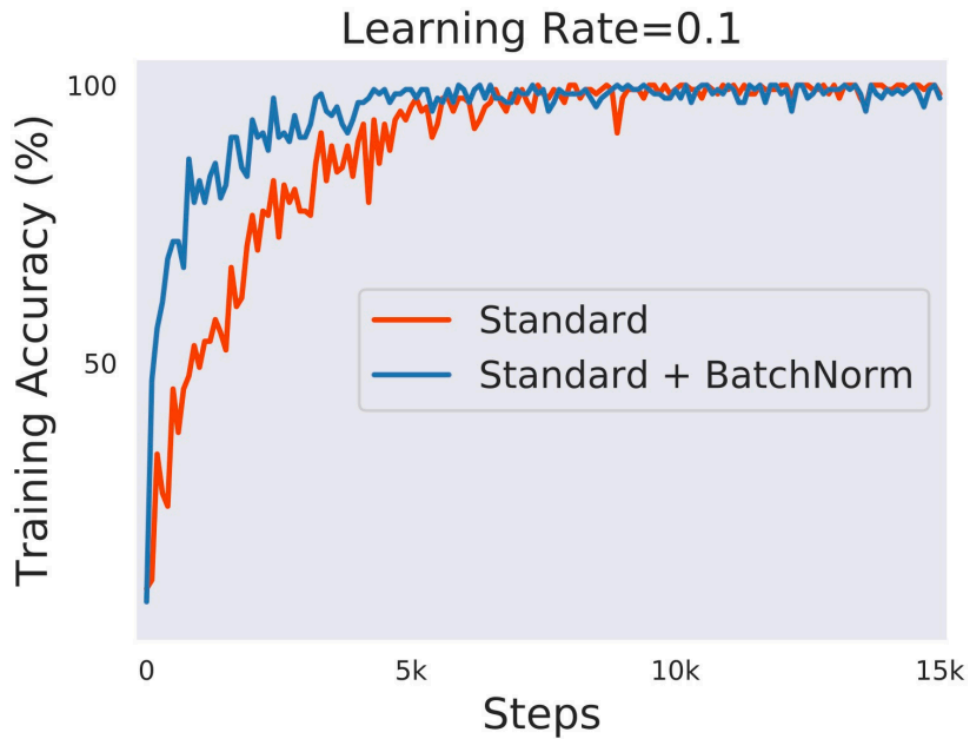
$$\frac{\partial \text{Div}}{\partial \mu_B} = \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i}$$



$$\frac{\partial \text{Div}}{\partial z_i} = \frac{\partial \text{Div}}{\partial u_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \text{Div}}{\partial \sigma_B^2} \cdot \frac{2(z_i - \mu_B)}{B} + \frac{\partial \text{Div}}{\partial \mu_B} \cdot \frac{1}{B}$$

The rest of backprop continues from  $\frac{\partial \text{Div}}{\partial z_i}$

# Batch Normalization (Ioffe & Szegedy, '14)



# What is BatchNorm actually doing?

---

- May not be due to covariate shift (Santurkar et al. '18):
  - Inject non-zero mean, non-standard covariance Gaussian noise after BN layer: removes the whitening effect
  - Still performs well.
- Only training  $\beta, \gamma$  with random convolution kernels gives non-trivial performance (Frankle et al. '20)
- BN can use exponentially increasing learning rate! (Li & Arora '19)  

---

# More normalizations

---

- Layer normalization (Ba, Kiros, Hinton, '16)
  - Batch-independent
  - Suitable for RNN, MLP
- Weight normalization (Salimans, Kingma, '16)
  - Suitable for meta-learning (higher order gradients are needed)
- Instance normalization (Ulyanov, Vedaldi, Lempitsky, '16)
  - Batch-independent, suitable for generation tasks
- Group normalization (Wu & He, '18)
  - Batch-independent, improve BatchNorm for small batch size

# Non-convex Optimization Landscape

---

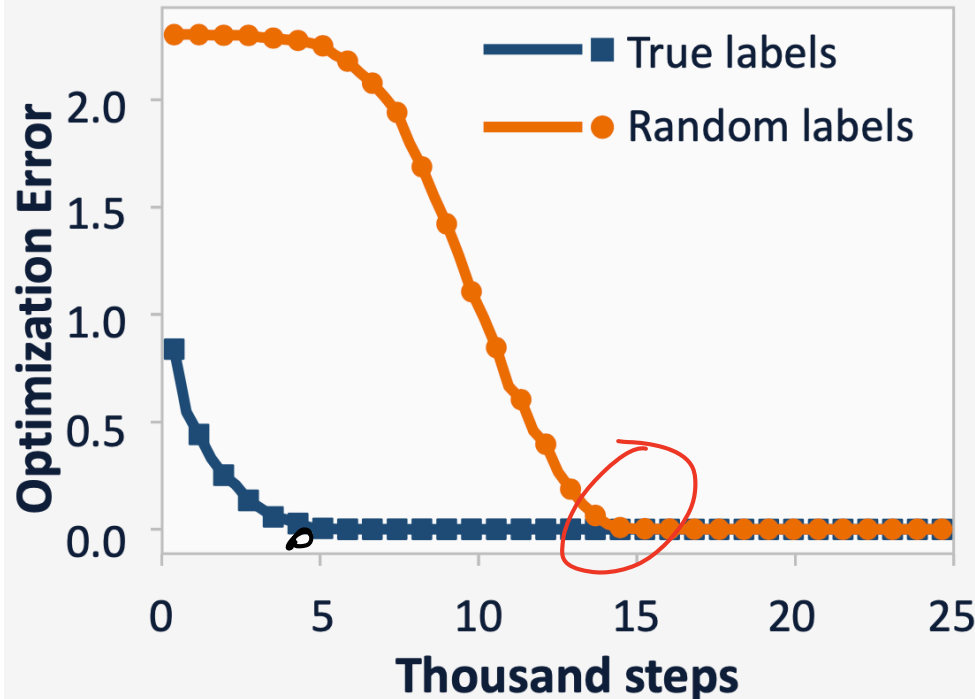


# Gradient descent finds global minima $\left\{ (x_i, y_i) \right\}_{i=1}^n$

Practice: gradient descent

$$\theta(t+1) \leftarrow \theta(t) - \eta \frac{\partial L(\theta(t))}{\partial \theta(t)}$$

$$\left\{ (x_i, \tilde{y}_i) \right\}_{i=1}^n$$
$$\tilde{y}_i \sim \mathcal{N}(0, 1)$$



Optimization error  $\rightarrow 0$  for both **true labels** and **random labels** !

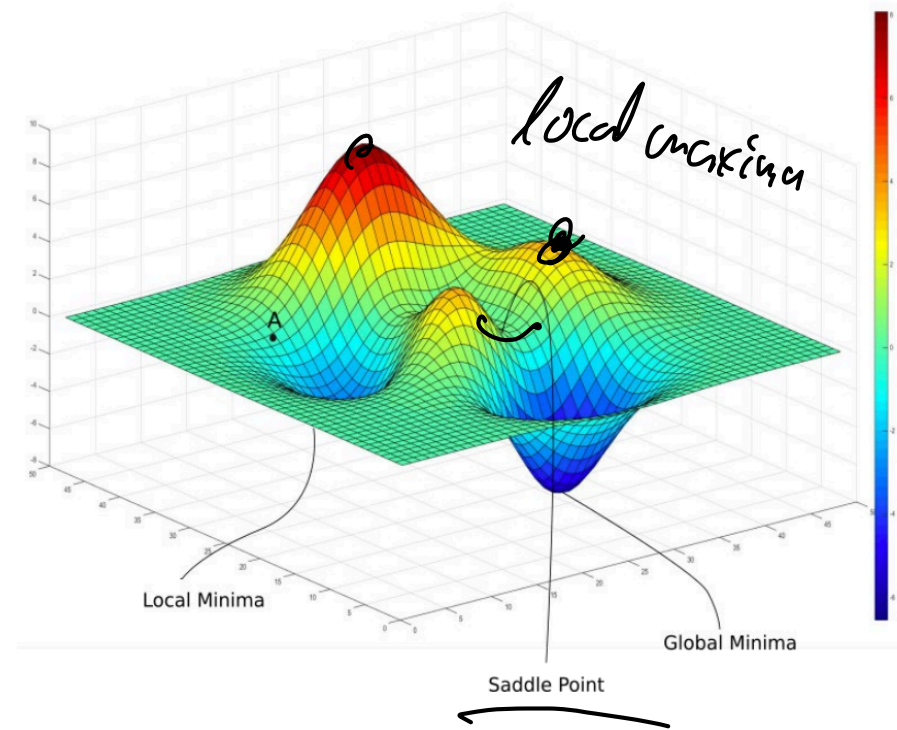
Zhang Bengio Hardt Recht Vinyals 2017

Understanding DL Requires Rethinking Generalization

# Types of stationary points

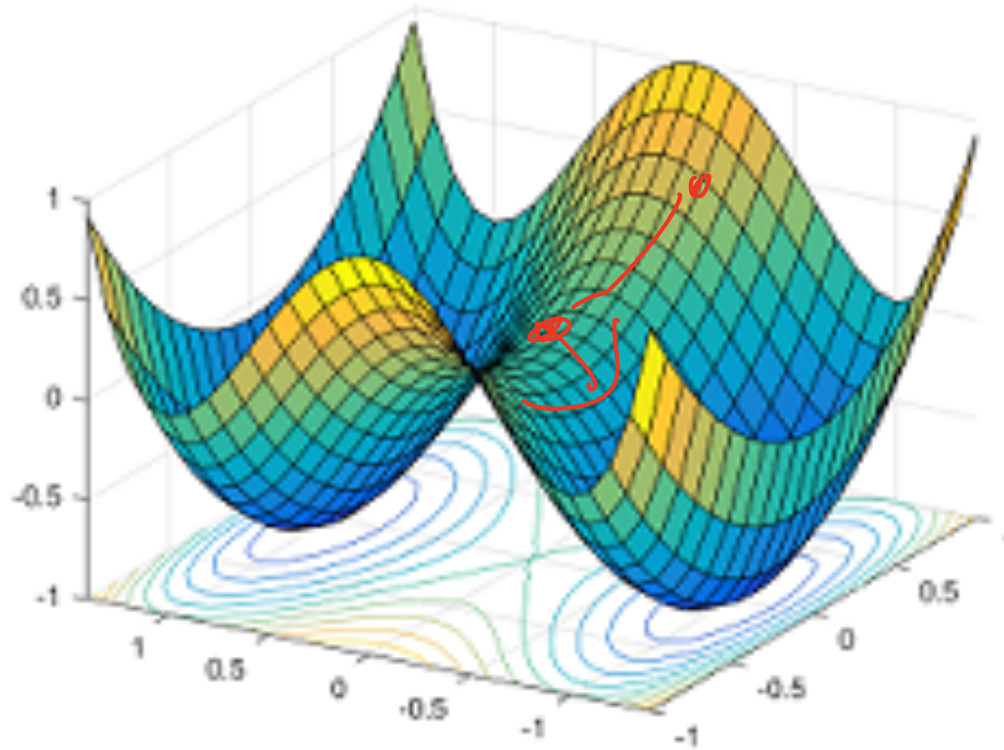
$$\min f(x)$$
$$f(x) = (NN(x) - y)^2$$

- Stationary points:  $x : \nabla f(x) = 0$
- Global minimum:  
 $x : f(x) \leq f(x') \forall x' \in \mathbb{R}^d$
- Local minimum:  
 $x : f(x) \leq f(x') \forall x' : \|x - x'\| \leq \epsilon$
- Local maximum:  
 $x : f(x) \geq f(x') \forall x' : \|x - x'\| \leq \epsilon$
- Saddle points: stationary points that are not a local min/max



# Landscape Analysis

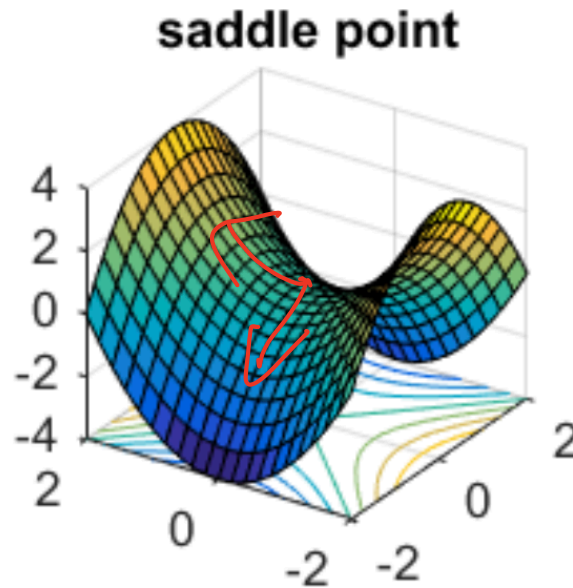
---



- All local minima are global!
  - Gradient descent can escape saddle points.
- $f$



# Strict Saddle Points (Ge et al. '15, Sun et al. '15)



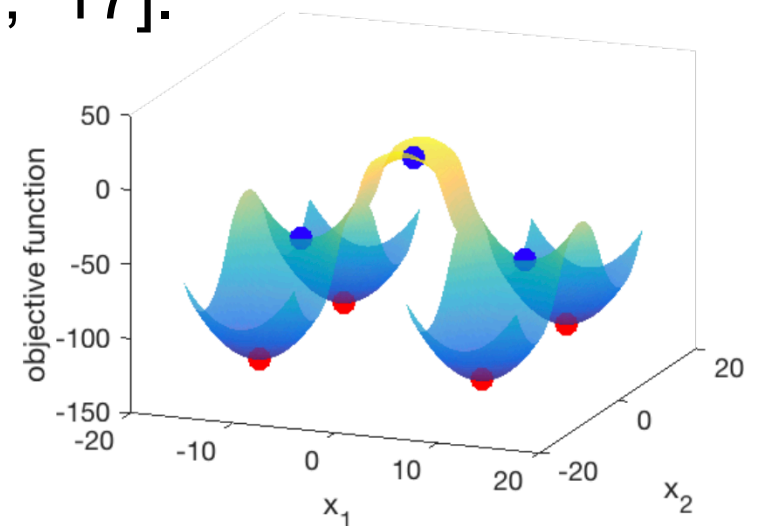
$$\nabla^2 f(x) = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \\ = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix},$$

- Strict saddle point: a saddle point and  $\lambda_{\min}(\nabla^2 f(x)) < 0$

# Escaping Strict Saddle Points

- **Noise-injected** gradient descent can escape strict saddle points in polynomial time [Ge et al., '15, Jin et al., '17].
- Randomly initialized gradient descent can escape all strict saddle points asymptotically [Lee et al., '15].
  - Stable manifold theorem.
- Randomly initialized gradient descent can take exponential time to escape strict saddle points [Du et al., '17].

If 1) all local minima are global, and 2) are saddle points are strict, then noise-injected (stochastic) gradient descent finds a global minimum in polynomial time



# What problems satisfy these two conditions

---

- Matrix factorization
- Matrix sensing
- Matrix completion
- Tensor factorization
- Two-layer neural network with quadratic activation

# What about neural networks?

---

- Linear networks (neural networks with linear activation functions): **all local minima are global, but there exists saddle points that are not strict** [Kawaguchi '16].
  - Non-linear neural networks with:
    - Virtually any non-linearity,
    - Even with Gaussian inputs,
    - Labels are generated by a neural network of the same architecture,
- There are many bad local minima** [Safran-Shamir '18, Yun-Sra-Jadbaie '19].