# Important Techniques in Neural Network Training

W

# Gradient Explosion / Vanishing

$\sigma: ReLU$

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

$$f(X, W_1, \ldots, W_H) = W_{H+1} \, \sigma(W_H \cdots \sigma(W_1 X) \cdots)$$

$$\frac{\partial f}{\partial W_n} = \left(W_{H+1} A_{H+1} \cdots W_{n+1} A_n\right)^T \left(A_{n-1} W_{n-1} \cdots W_1 X\right)$$

$$A_n = diag\left(\sigma'\left(W_n \, \sigma(\cdots \sigma(W_1 X) \cdots)\right)\right)$$

① magnitude  if  $W_n \cdots W_{H+1}$ small $\Rightarrow$ exp small

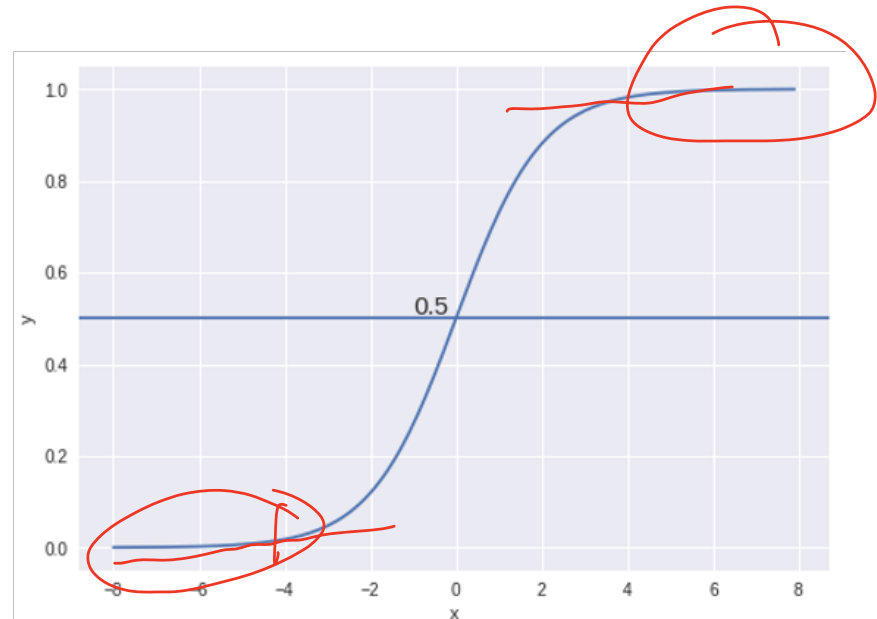$W_1, \ldots W_{H+1}$ large $\Rightarrow$ exp large

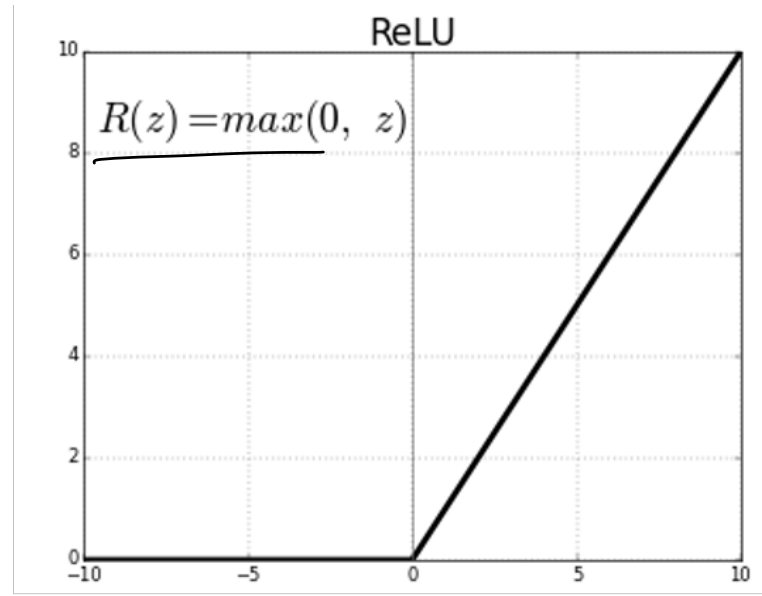② spale may not align $\longrightarrow$ multiplication small

# Activation Functions



tanh

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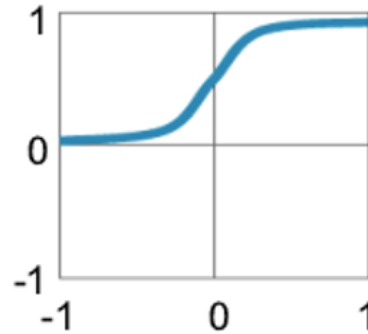
sigmoid

$$\frac{1}{1 + \exp(-z)}$$

ReLU

$$R(z) = max(0, \ z)$$

Rectified Linear United

$$\sigma'(z) = 0 \ or \ 1$$

# Activation Function



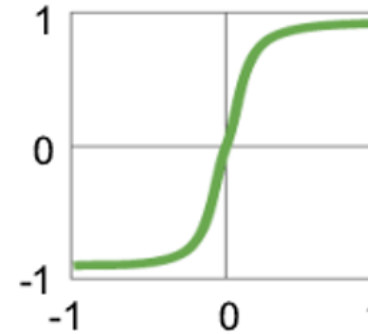**Traditional Non-Linear Activation Functions**

**Sigmoid**

$y=1/(1+e^{-x})$

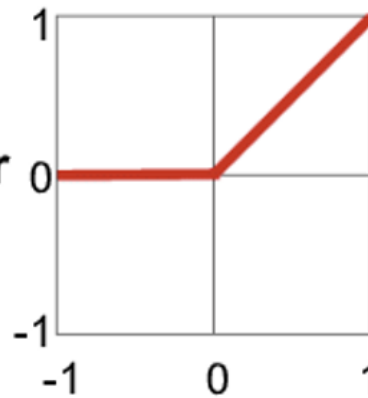**Hyperbolic Tangent**

$y=(e^x-e^{-x})/(e^x+e^{-x})$

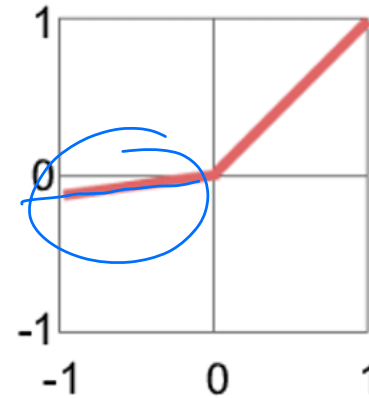**Modern Non-Linear Activation Functions**

**Rectified Linear Unit (ReLU)**

$y=\max(0,x)$

**Leaky ReLU**

$y=\max(\alpha x,x)$

**Exponential LU**

$y=\begin{cases} x, & x\geq 0 \\ \alpha(e^x-1), & x<0 \end{cases}$

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

# Initialization

$$\text{convex:} \quad \frac{f(x_0) - f(x_*)}{T}$$

$W_n^{ij} \sim \text{Gaussian or Unit}$

- Zero-initialization
- Large initialization $\rightarrow$ exp large grad
- Small initialization $\rightarrow$ exp small gradient

- Design principles:
  - Zero activation mean     (no prior knowledge)

  - Activation variance remains same across layers

# Kaiming Initialization (He et al. '15)
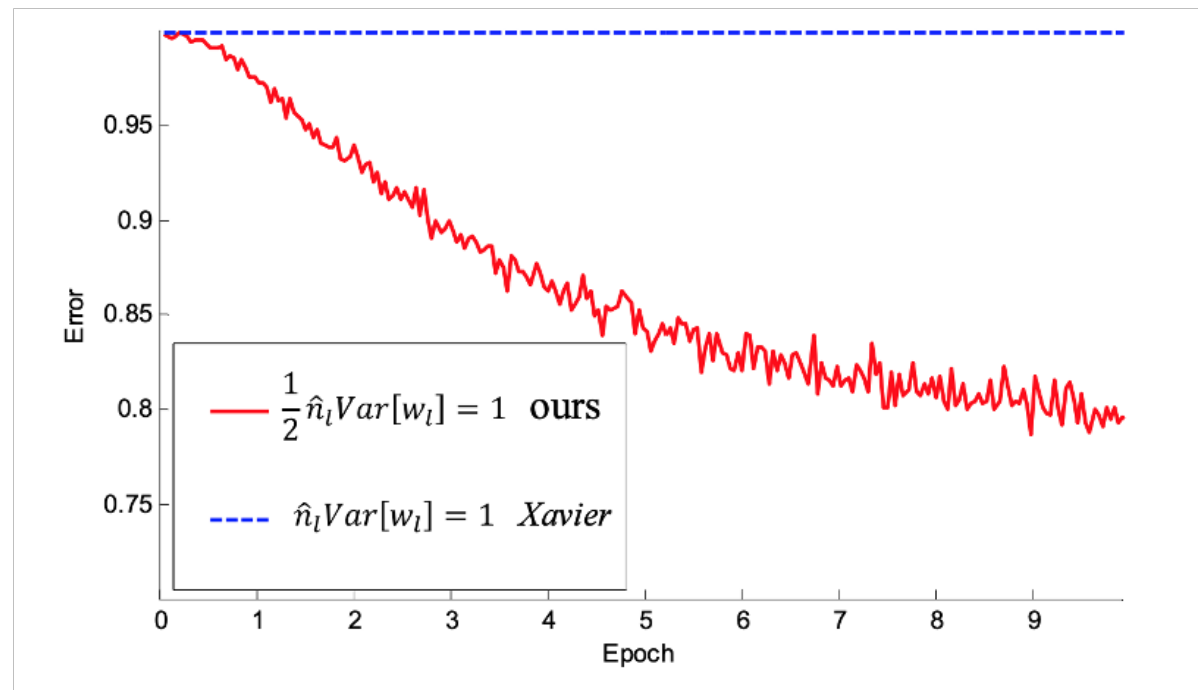
$d_h : \text{fan-in}$
$d_{h+1} : \text{fan-out}$

$W_h : \mathbb{R}^{d_{h+1} \times d_h}$

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

# Kaiming Initialization (He et al. '15)

Each layer

$$Z^h = W^h \cdot X^h \qquad \text{\color{red}pre-activation}$$

$$X^{h+1} = \sigma(Z^h)$$

$$Z_i^h = \sum_{j=1}^{d_h} W_{ij}^h X_j^h$$

**Goal**  $\qquad$ $Z^h$ : mean zero, same Var for all layers

$$\text{if } \mathbb{E}[W_{ij}^h] = 0 \implies \mathbb{E}[Z_i^h] = 0 ,$$

$$\text{Var}(Z_i^h) = d_h \cdot \text{Var}(W_{ij}^h X_j^h)$$

$$= d_h \left( \text{Var}(W_{ij}^h) \cdot \text{Var}(X_j^h) \right.$$

$$\left. + \underbrace{(\mathbb{E}[W_{ij}^h])^2}_{=0} \cdot \text{Var}(X_j^h) + \text{Var}(W_{ij}^h) \cdot \mathbb{E}[X_j^h]^2 \right)$$

$$= d_h \cdot \text{Var}(W_{ij}^h) \cdot \mathbb{E}[(X_j^h)^2]$$

# Kaiming Initialization (He et al. '15)

$$\mathbb{E}\left[(x_j^h)^2\right] = \int_{-\infty}^{\infty} \left(x_j^h\right)^2 P\left(x_j^h\right) d\,x_j^h$$

$$= \int_{-\infty}^{\infty} \max\left(0, z_j^{h-1}\right)^2 \cdot P\left(z_j^{h-1}\right) d\,z_j^{h-1}$$

$$= \int_{0}^{\infty} \left(z_j^{h-1}\right)^2 P\left(z_j^{h-1}\right) d\,z_j^{h-1}$$

$$= \frac{1}{2} \int_{-\infty}^{\infty} \left(z_j^{h-1}\right)^2 P\left(z_j^{h-1}\right) d\,z_j^{h-1}$$

$$= \frac{1}{2} \text{Var}\left(z_j^{h-1}\right)$$

(ReLU)

$\left(\begin{array}{c}\text{symmetry}\\\text{of init}\\\text{around } 0\end{array}\right)$

# Kaiming Initialization (He et al. '15)

We want $\text{Var}\left(z_i^h\right) = \text{Var}\left(z_j^{h-1}\right)$

$$d_h \cdot \text{Var}\left(w_{ij}^h\right) \cdot \frac{1}{2} \cdot \text{Var}\left(z_j^{h-1}\right) = \text{Var}\left(z_j^{h-1}\right)$$

$$\implies \text{Var}\left(w_{ij}^h\right) = \frac{2}{d_h}$$

$$\text{Var}\left(z_j^h\right) = \text{Var}\left(z^0\right)\left(\prod_{h'=1}^{h} \frac{d_{h'}}{2} \text{Var}\left(w_{ij}^{h'}\right)\right)$$

$$\underbrace{\qquad}_{O(1)}$$

# Initialization by Pre-training

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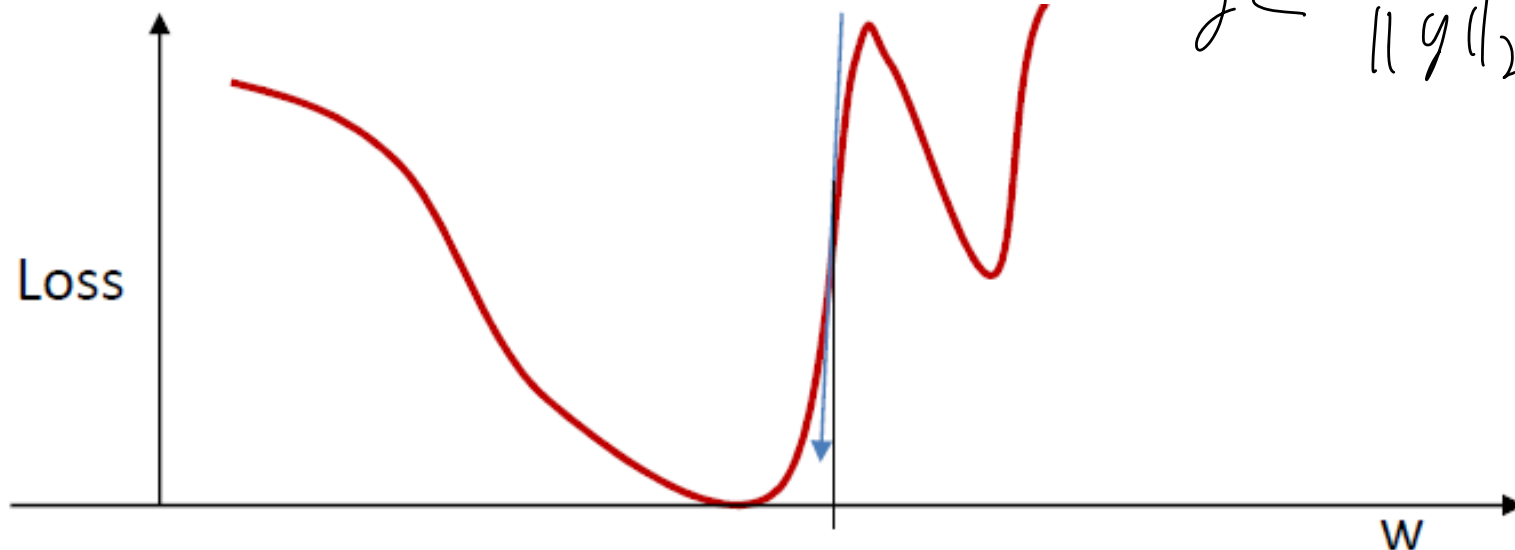- And then fine-tuning

# 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.

$$g = \nabla f(x_t) \quad \text{if } \|g\|_2 > \text{threshold}$$

$$g \leftarrow \frac{g}{\|g\|_2} \cdot \text{threshold}$$
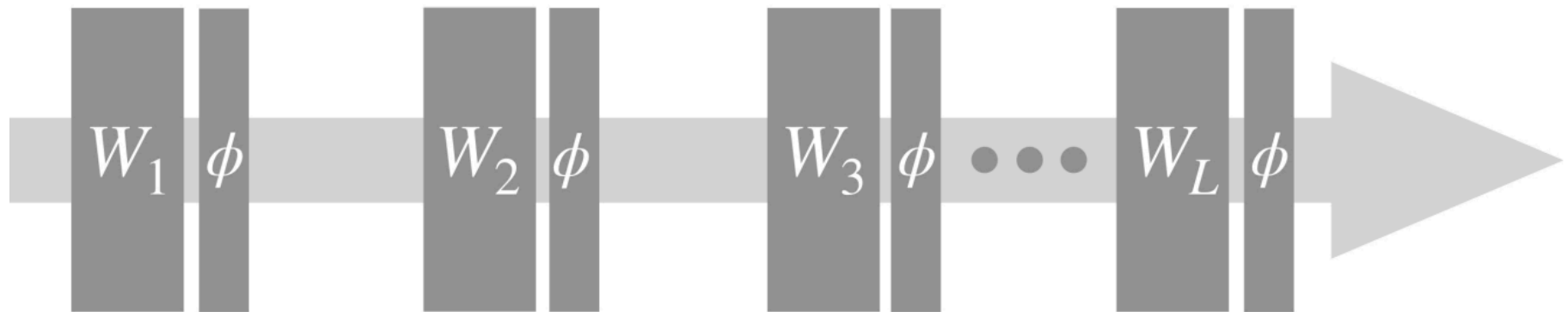
Loss

W

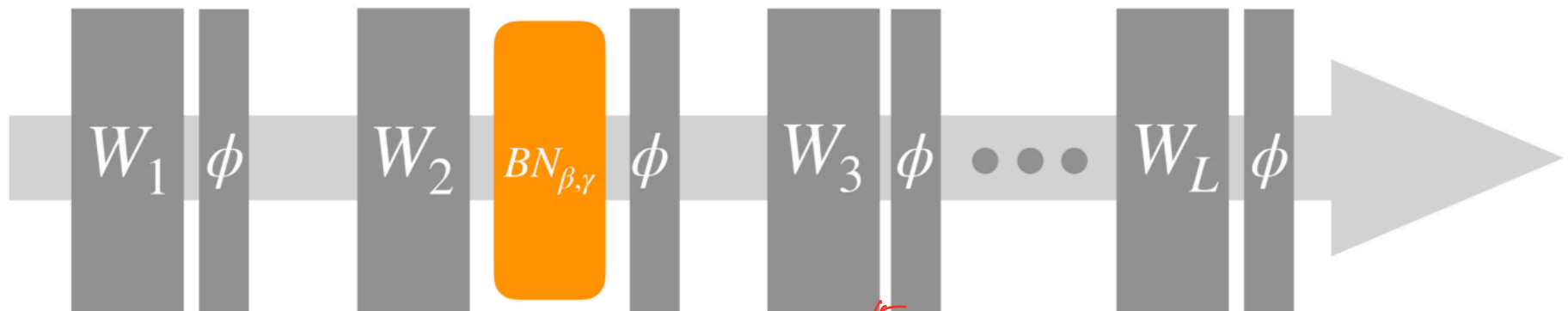# Batch Normalization (Ioffe & Szegedy, '14)

- <span style="color:red">Normalizing/whitening</span> (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!

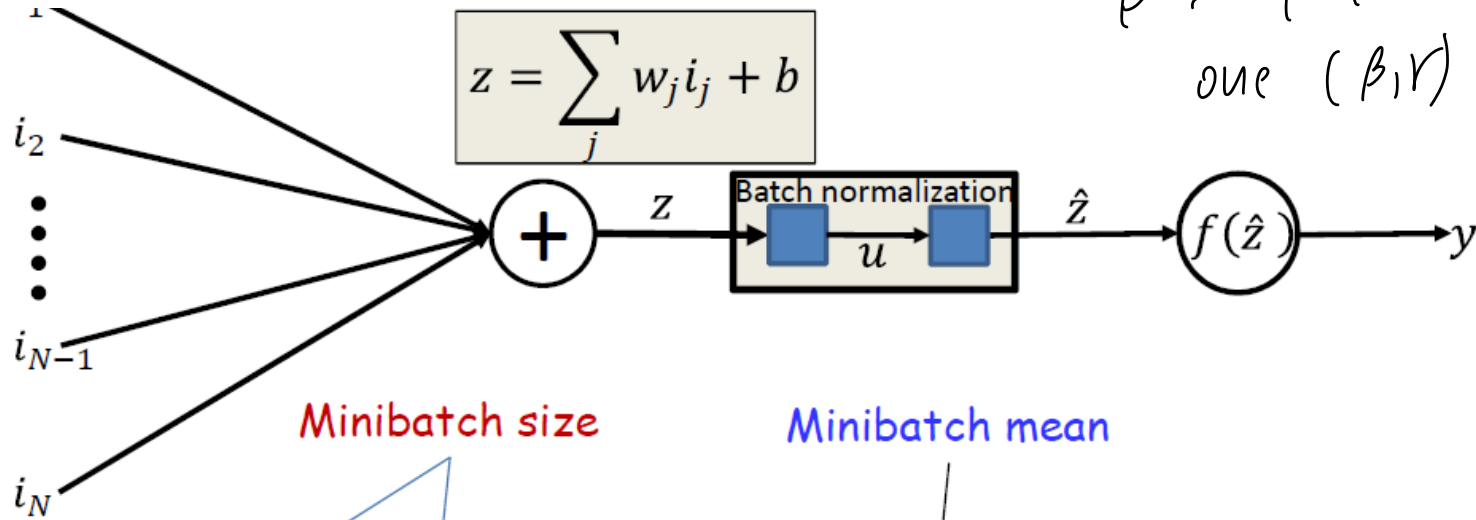# Batch Normalization (Ioffe & Szegedy, '14)

**Standard Network**



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

# Batch Normalization (Ioffe & Szegedy, '14)

$\gamma$ : population std

$\beta$ : population mean

one $(\beta, \gamma)$ for each BN layer

$$z = \sum_j w_j i_j + b$$

$i_2$

$\vdots$

$i_{N-1}$

$i_N$

$z$  →  **Batch normalization** $u$  →  $\hat{z}$  →  $f(\hat{z})$  →  $y$

**Minibatch size**

**Minibatch mean**

**Minibatch standard deviation**

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

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

$$u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$
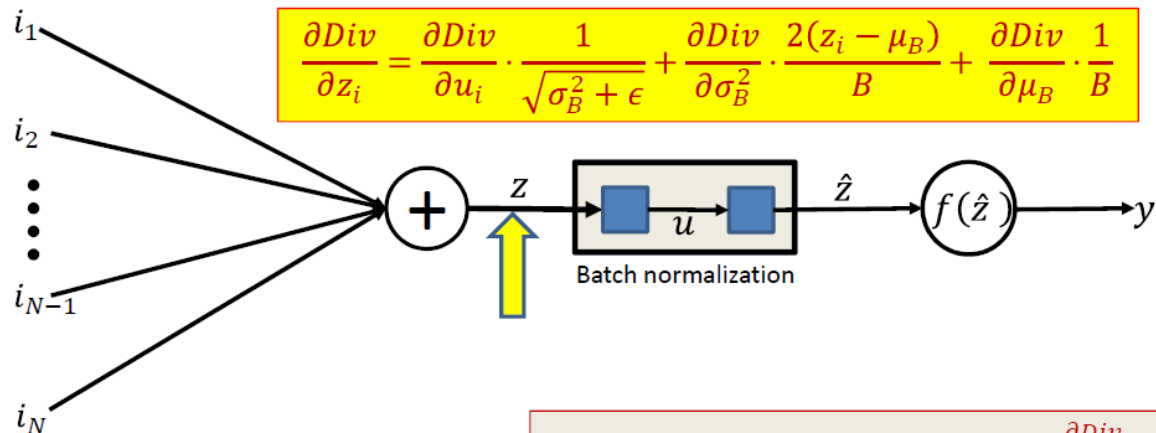
$$\hat{z}_i = \gamma u_i + \beta$$

$\hat{z}_i$ : mean $\beta$

std : $\gamma$

# Batch Normalization (Ioffe & Szegedy, '14)

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

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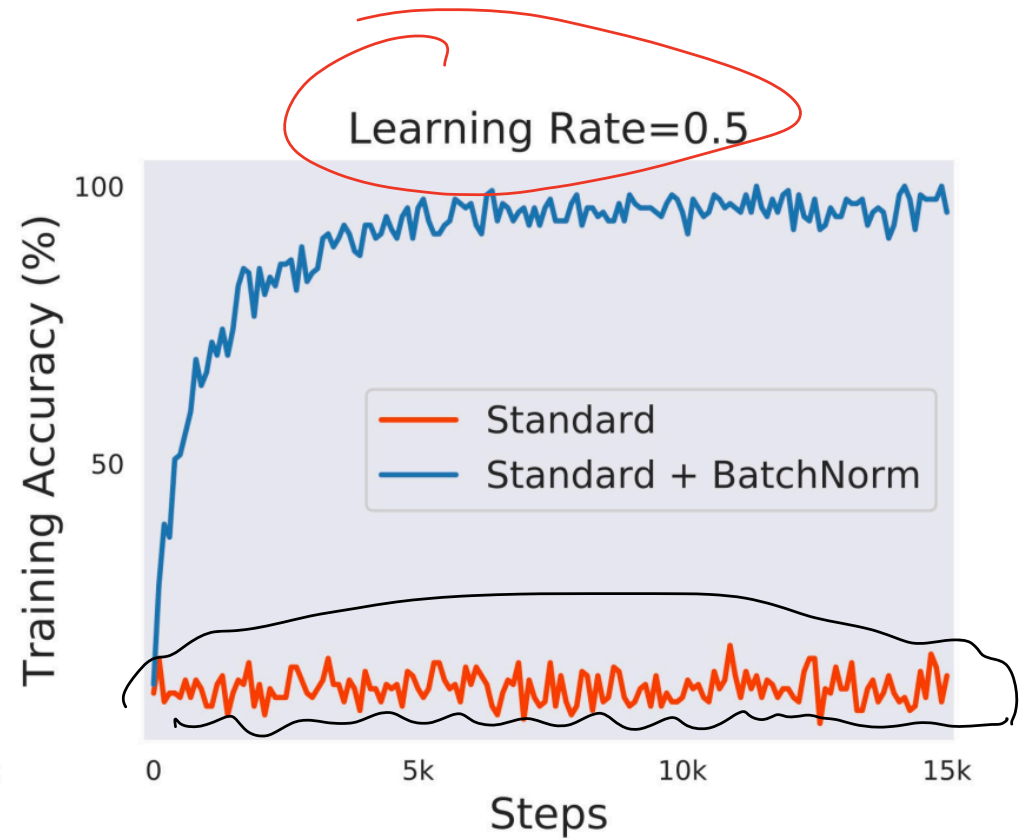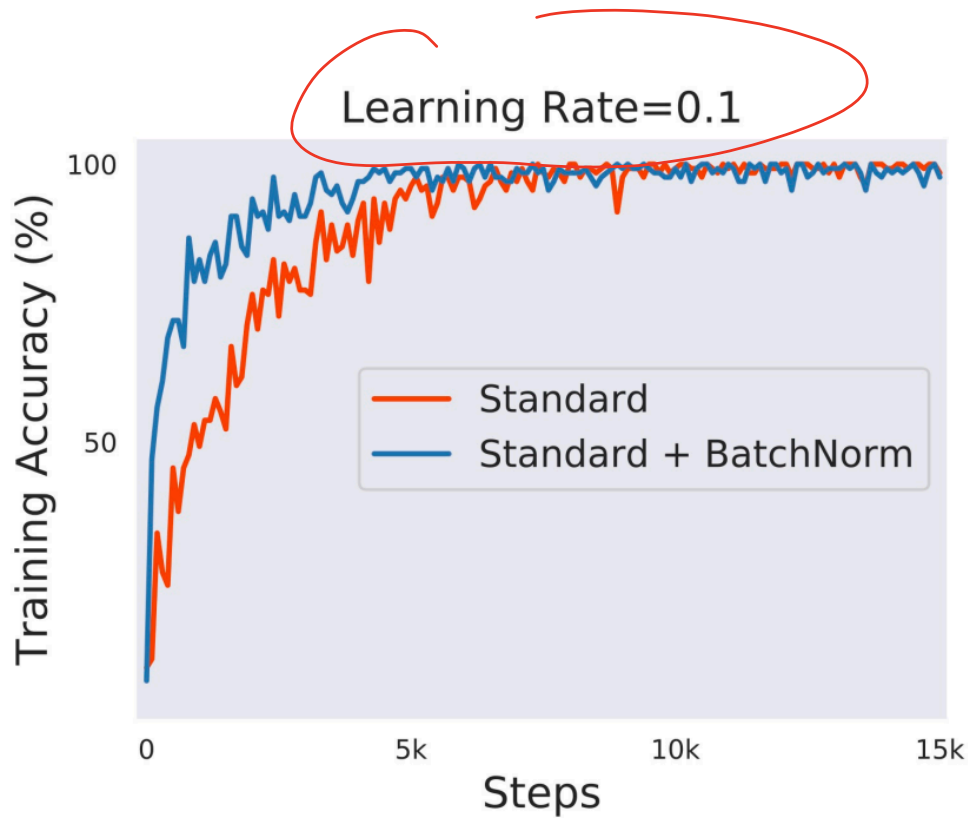
$$\frac{\partial Div}{\partial \mu_B} = \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \sum_{i=1}^{B} \frac{\partial Div}{\partial u_i}$$

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



Batch normalization

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

# Batch Normalization (Ioffe & Szegedy, '14)

# What is BatchNorm actually doing?

- May not 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

$$X^y, \quad \frac{X^y}{\|X^y\|_2}$$

- 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)
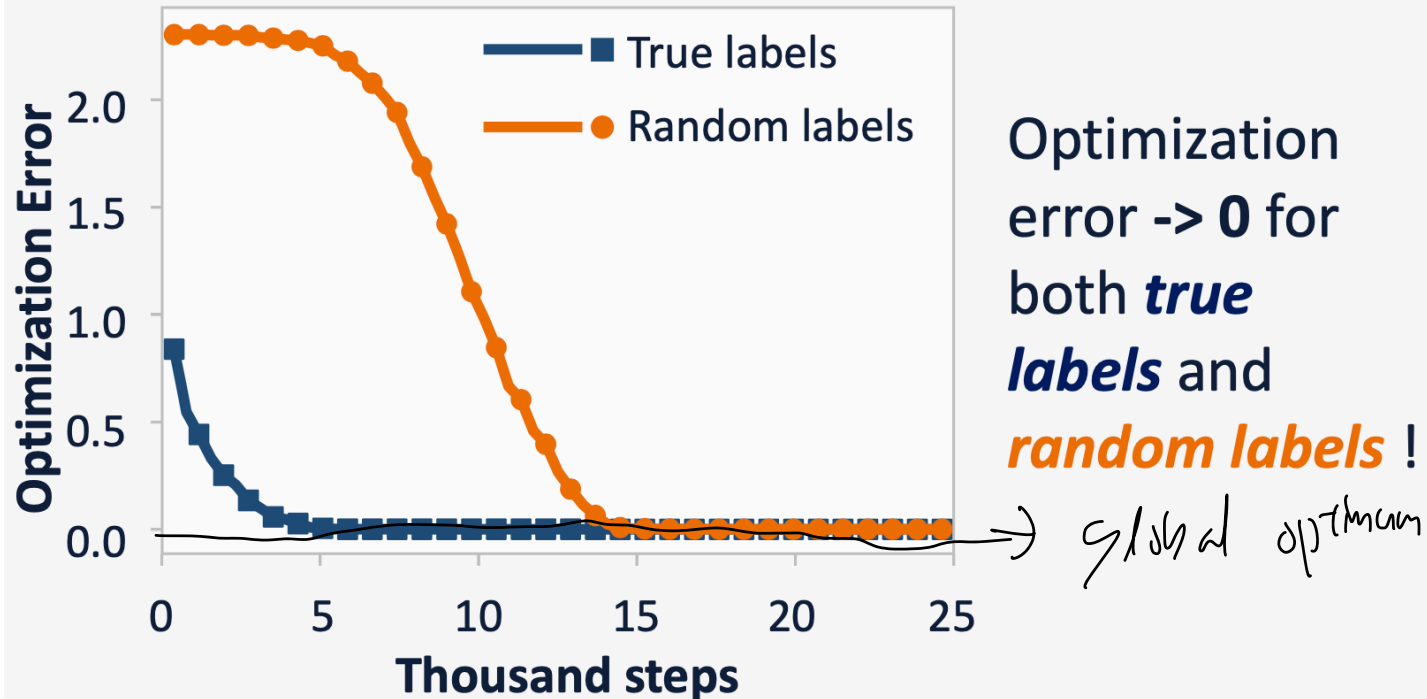- ….

# Non-convex Optimization Landscape

W

# Gradient descent finds global minima

over-parameterized



**Practice:** gradient descent

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

Optimization
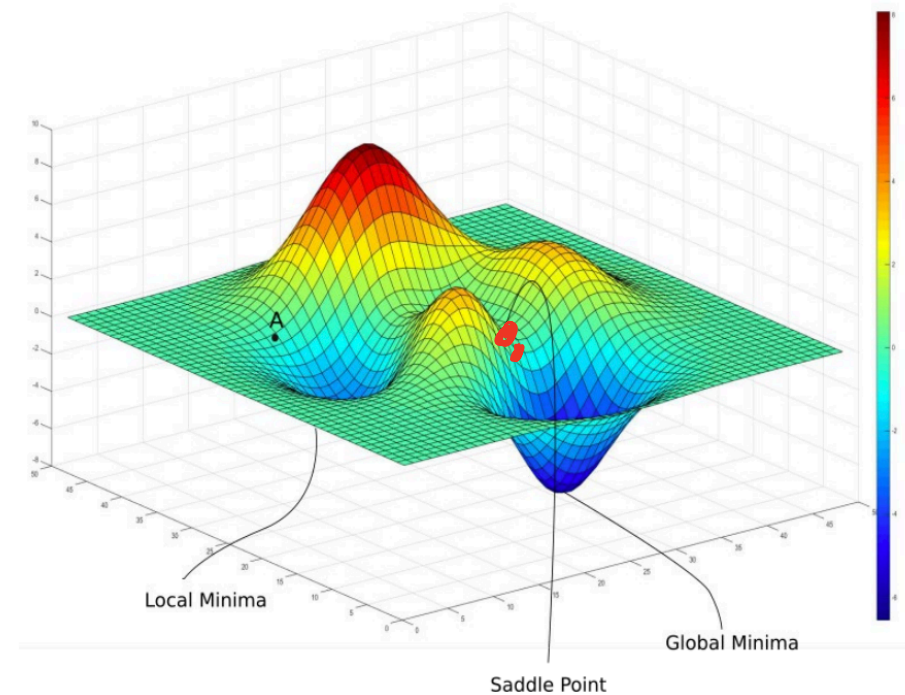error -> **0** for
both ***true
labels*** and
***random labels*** !

→ Global optimum
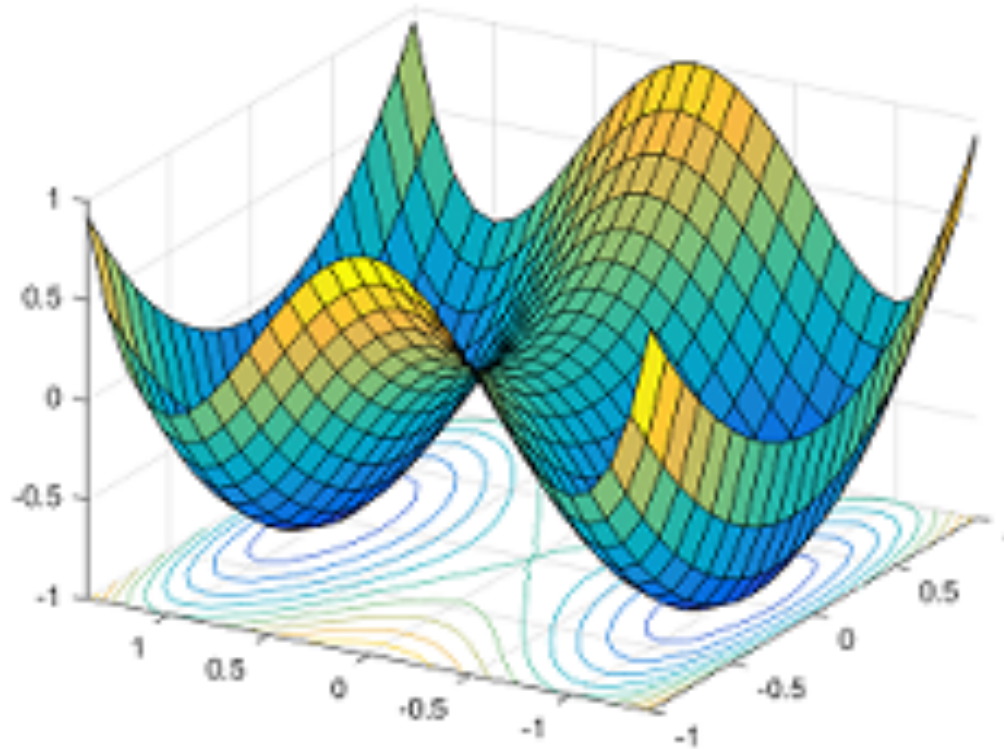
Zhang Bengio Hardt Recht Vinyals 2017
Understanding DL Requires Rethinking Generalization

# Types of stationary points

- 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:
$x : f(x) \geq f(x') \forall x' : \|x - x'\| \leq \epsilon$
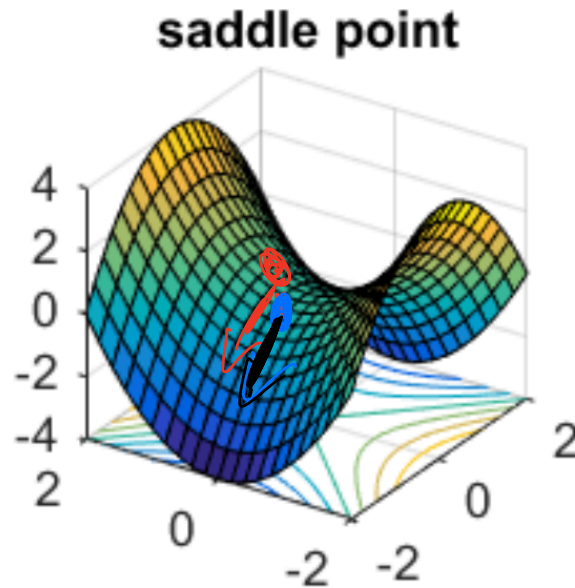- Saddle points: stationary points that are not a local min/max



Local Minima

Global Minima

Saddle Point

# Landscape Analysis



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

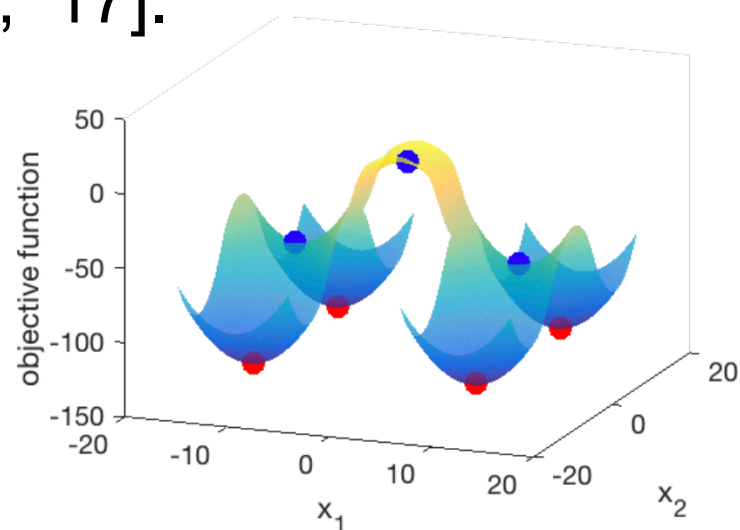# Strict Saddle Points (Ge et al. '15, Sun et al. '15)



saddle point

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

$$\nabla^2 f(x) \cdot \nabla f(x)$$

# 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 activations functions): all local minima are global, but there exists saddle points that are not strict [Kawaguchi '16].

$$W_{H+1} W_H \cdots W_1 X$$

- 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].