

# Important Techniques in Neural Network Training

---



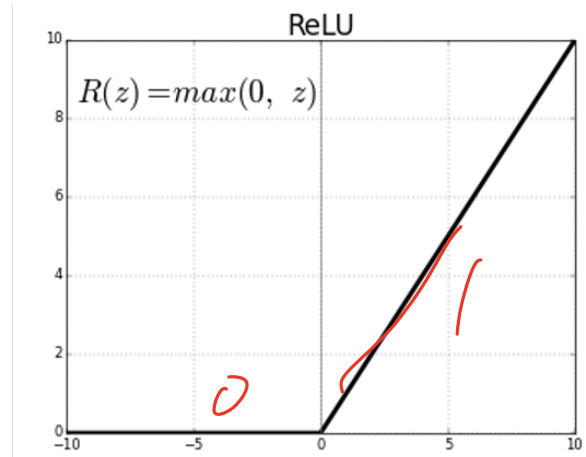
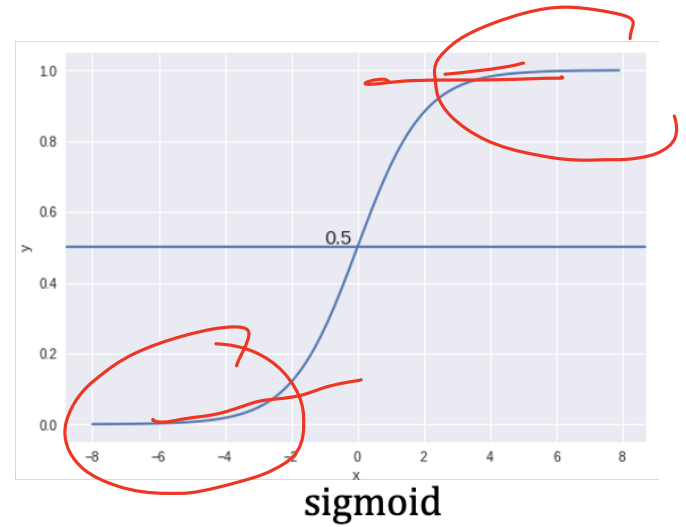
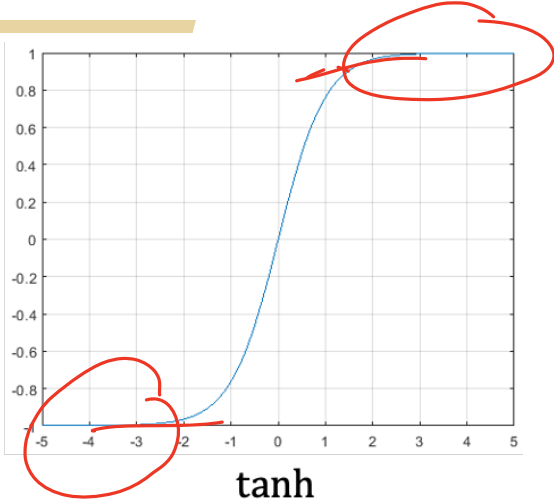
# Gradient Explosion / Vanishing

---

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

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

# Activation Functions

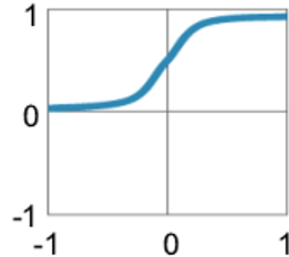


Rectified Linear Unit

# Activation Function

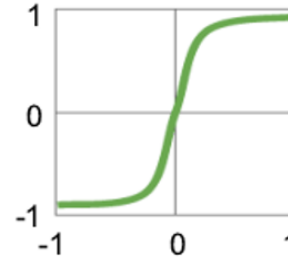
**Traditional  
Non-Linear  
Activation  
Functions**

**Sigmoid**



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

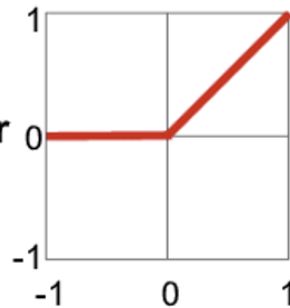
**Hyperbolic Tangent**



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

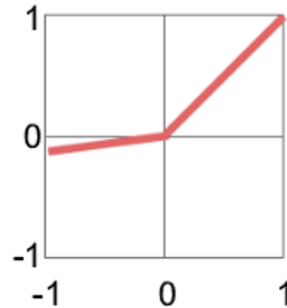
**Modern  
Non-Linear  
Activation  
Functions**

**Rectified Linear Unit  
(ReLU)**



$$y = \max(0, x)$$

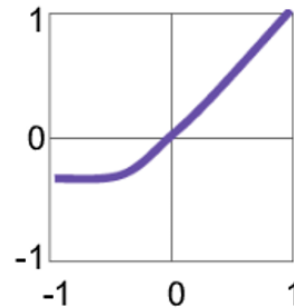
**Leaky ReLU**



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

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

**Exponential LU**



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

# Initialization

---

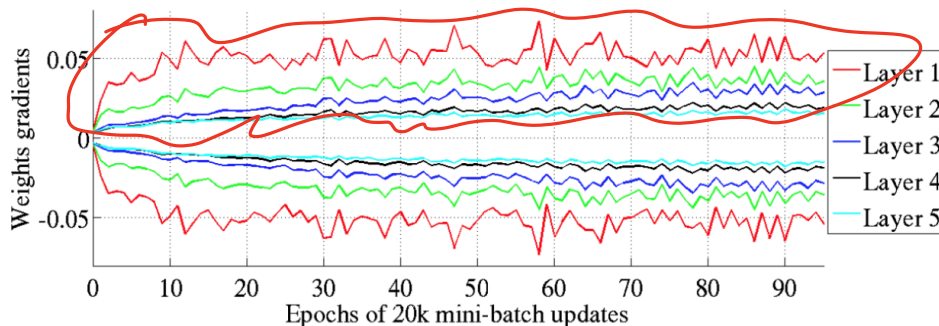
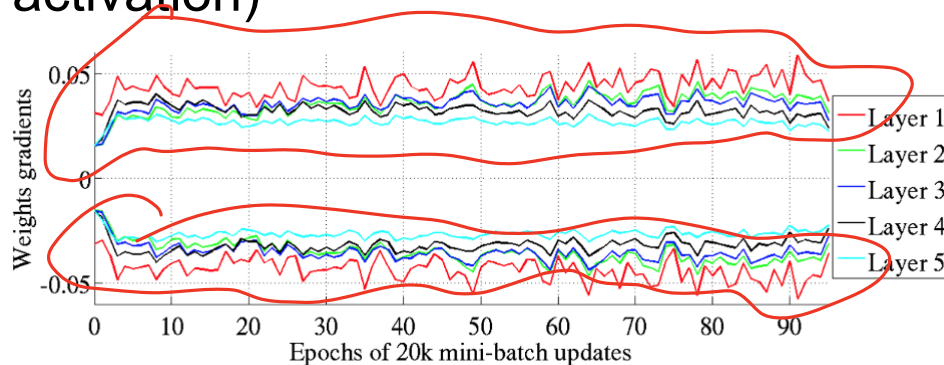
- Zero-initialization
  - Large initialization
  - Small initialization
- gradient 0
- > scaling
- 
- Design principles:
    - Zero activation mean
    - Activation variance remains same across layers
- balancing
- >

# 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)

$$\mathbb{E}(W_{ij}^{(h)}) = 0$$

$$\text{Var}(W_{ij}^{(h)}) = \frac{2}{d_h + d_{h+1}}$$



unif

$$\left( -\frac{1}{\sqrt{d_h}}, \frac{1}{\sqrt{d_h}} \right)$$



# Kaiming Initialization (He et al. '15)

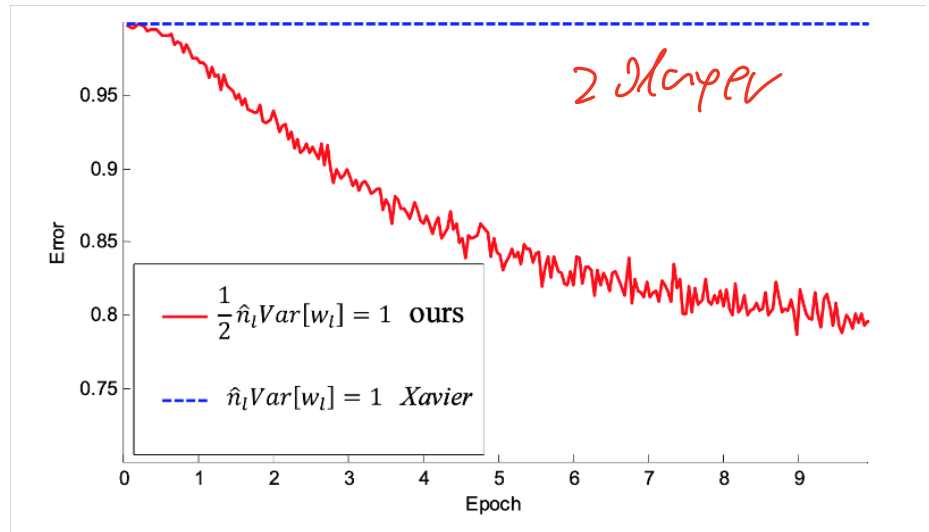
- $W_{ij}^{(h)} \sim \mathcal{N}\left(0, \frac{2}{d_h}\right)$ .

if  $dw = dw_{out}$   
v.s.  $var = \frac{1}{d_h}$

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

- Designed for ReLU activation

- 30-layer neural network



# Kaiming Initialization (He et al. '15)

$$Z^h = W^h \cdot X^h$$

$$X^1, W^1 X^h, \delta(W^1 X^h) \dots$$

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

$$Z^1, X^2,$$

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

$$\begin{bmatrix} \vdots \\ Z_i^h \end{bmatrix} = \begin{bmatrix} \vdots \\ W^h \end{bmatrix} \begin{bmatrix} \vdots \\ X^h \end{bmatrix}$$

$$\mathbb{E}[Z_i^h] = 0$$

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

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

$$+ (\mathbb{E} W_{ij}^h)^2 \cdot \text{Var}(X_j^h) + \text{Var}(W_{ij}^h) \cdot (\mathbb{E} (X_j^h))^2)$$

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



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

$$\begin{aligned} \mathbb{E}[(X_j^h)^2] &= \int_{-\infty}^{\infty} (X_j^h)^2 p(X_j^h) dX_j^h \\ &= \int_{-\infty}^{\infty} \max(0, z_j^{h-1})^2 p(z_j^{h-1}) dz_j^{h-1} \\ &\stackrel{\text{ReLU}}{=} \int_0^{\infty} (z_j^{h-1})^2 p(z_j^{h-1}) dz_j^{h-1} \\ &\stackrel{(\text{Symmetry of initialization})}{=} \frac{1}{2} \int_{-\infty}^{\infty} (z_j^{h-1})^2 p(z_j^{h-1}) dz_j^{h-1} \\ &= \frac{1}{2} \text{Var}(z_j^{h-1}) \end{aligned}$$

Want

$$\text{Var}(Z_j^h) = \text{Var}(Z_j^{h-1})$$

$$dh \cdot \text{Var}(W_{j,i}^h) \stackrel{!}{=} \frac{1}{2} \cancel{\text{Var}(Z_j^{h-1})} = \cancel{\text{Var}(Z_j^{h-1})}$$

$$\text{Var}(W_{j,i}^h) = \frac{2}{dh}$$

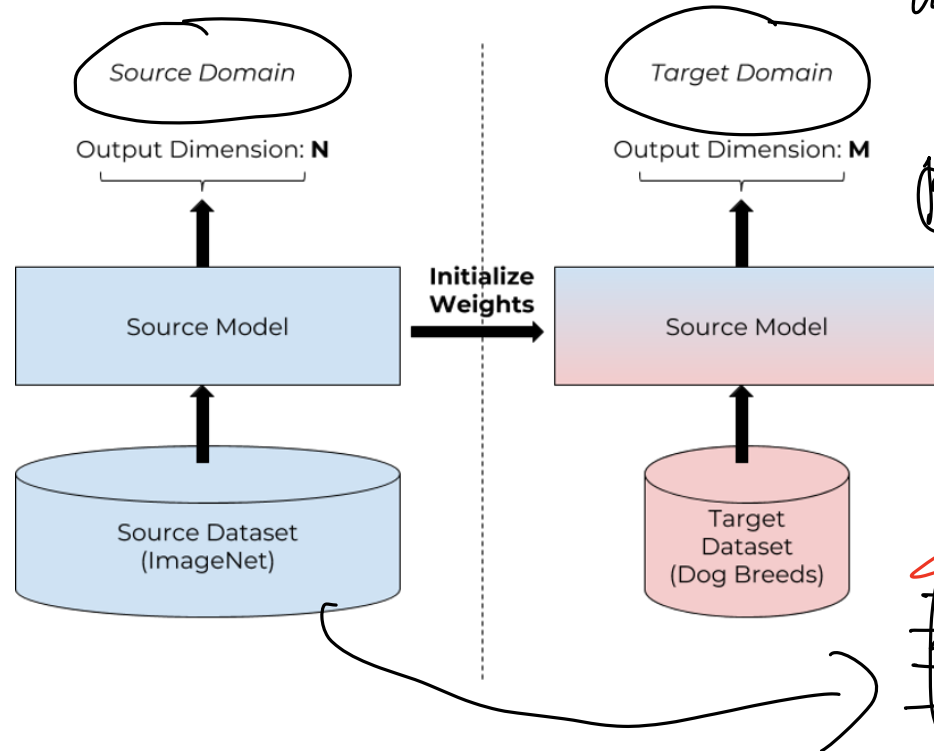
$$\text{Var}(Z^{t+1}) = \text{Var}(Z^1) \underbrace{\left( \prod_{h=1}^{t-1} \frac{dh}{2} \text{Var}(W_{j,i}^h) \right)}$$

$\mathcal{O}(C)$

# Initialization by Pre-training

BERT, GPT-3

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



train on  
wiki

predict  
missing  
word

MM

machine  
translation

only the  
last layer

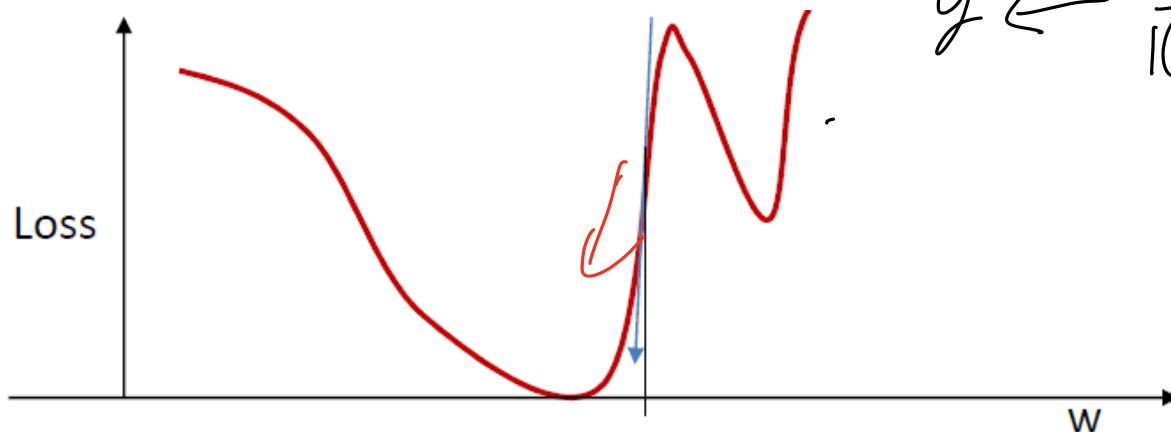
re-random  
initialization

re-random  
initialization

# 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.

$$g = \nabla f(x^t), \quad \|g\|_2 > \text{threshold}$$
$$g \leftarrow \frac{g}{\|g\|_2} \cdot \text{threshold}$$



# 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!)

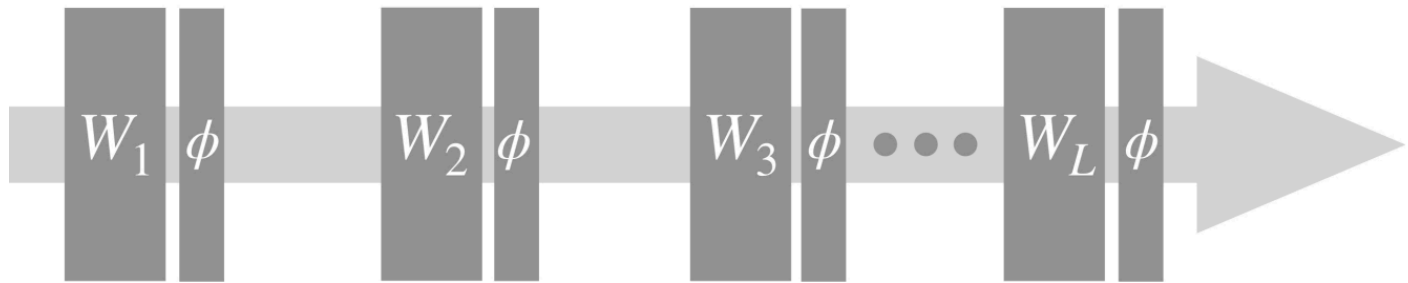
population mean  
variance

mean →  
variance →  
some  
mean  
var

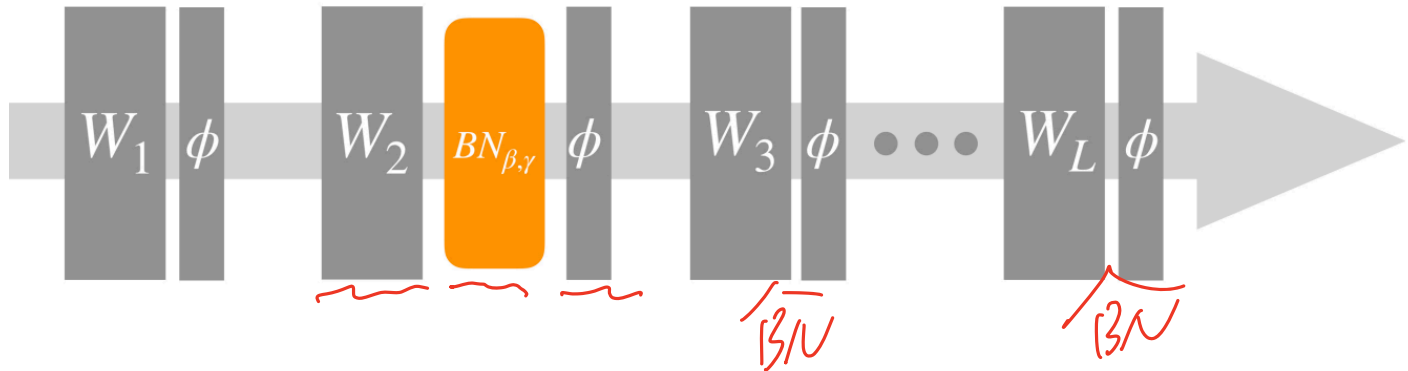
# Batch Normalization (Ioffe & Szegedy, '14)

*new layer*

Standard Network



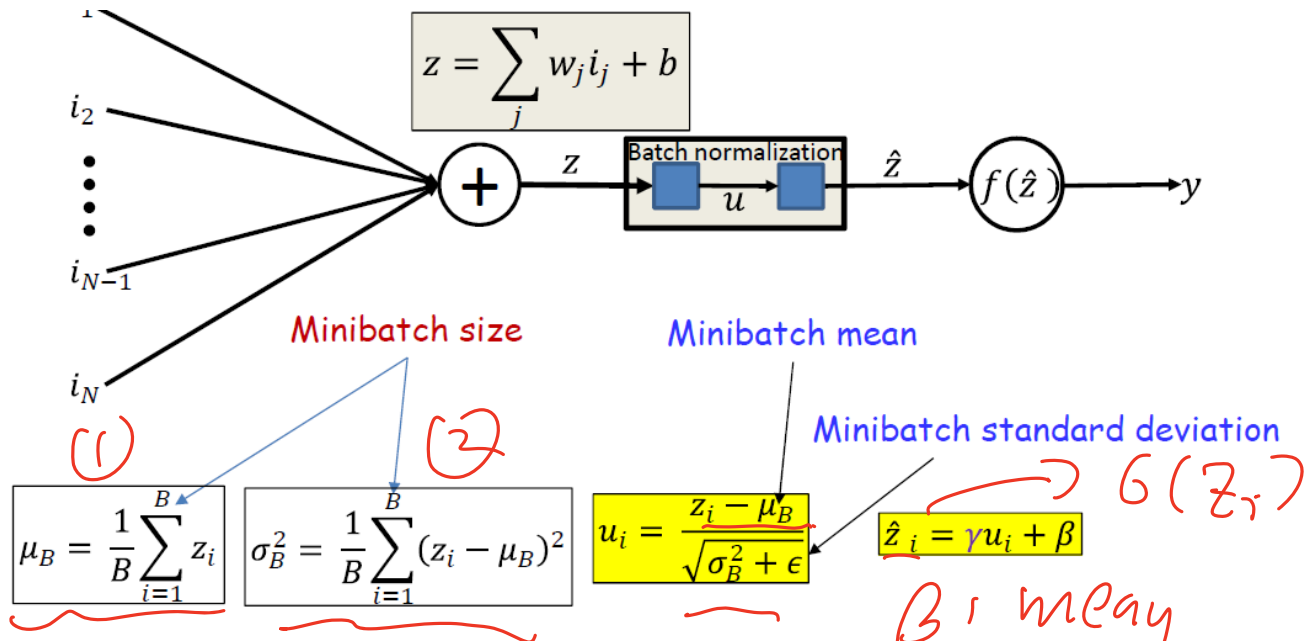
Adding a BatchNorm layer (between weights and activation function)



# Batch Normalization (Ioffe & Szegedy, '14)

$$z_i = w_i \cdot X$$

$\epsilon$ : hyper-parameter



$$z_1, \dots, z_B$$

$B$ : batch size

$\beta$ : mean

$\gamma$ : std

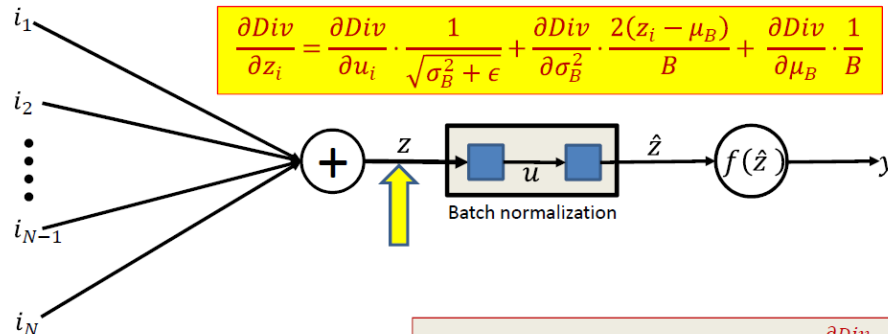
# Batch Normalization (Ioffe & Szegedy, '14)

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

32, 128  
1024

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

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

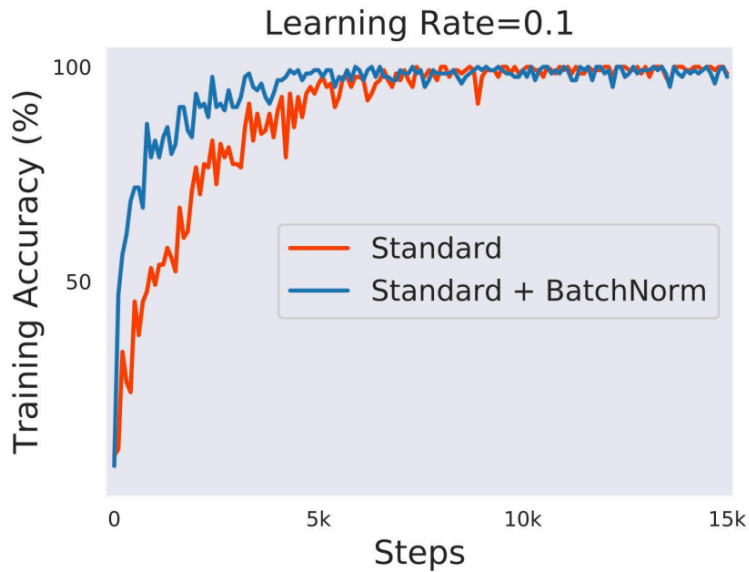


$$\frac{\partial \text{Div}}{\partial z_i} = \frac{\partial \text{Div}}{\partial u_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \text{Div}}{\partial \sigma_B^2} \cdot \frac{2(z_i - \mu_B)}{B} + \frac{\partial \text{Div}}{\partial \mu_B} \cdot \frac{1}{B}$$

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 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)
- Instant 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

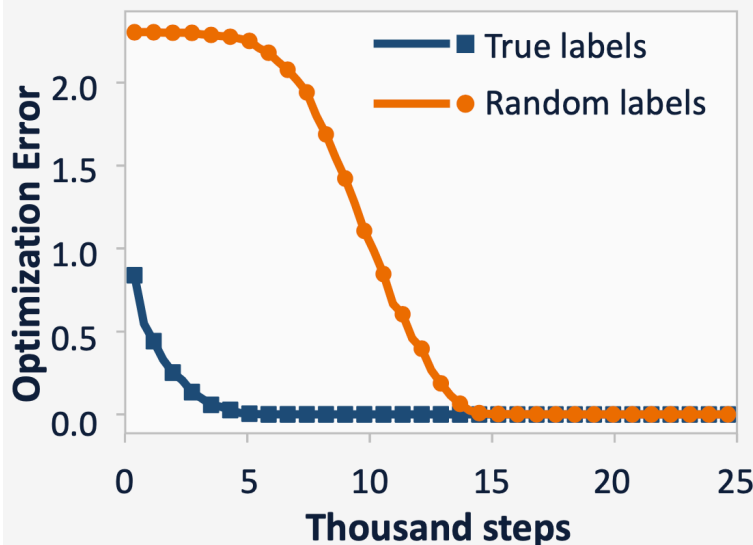
---



# Gradient descent finds global minima

**Practice:** gradient descent

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



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

Zhang Bengio Hardt Recht Vinyals 2017

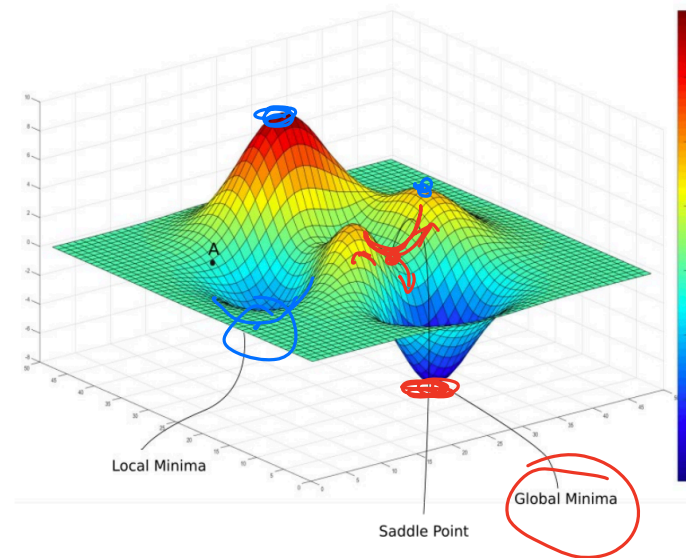
Understanding DL Requires Rethinking Generalization

# Types of stationary points

found by GN

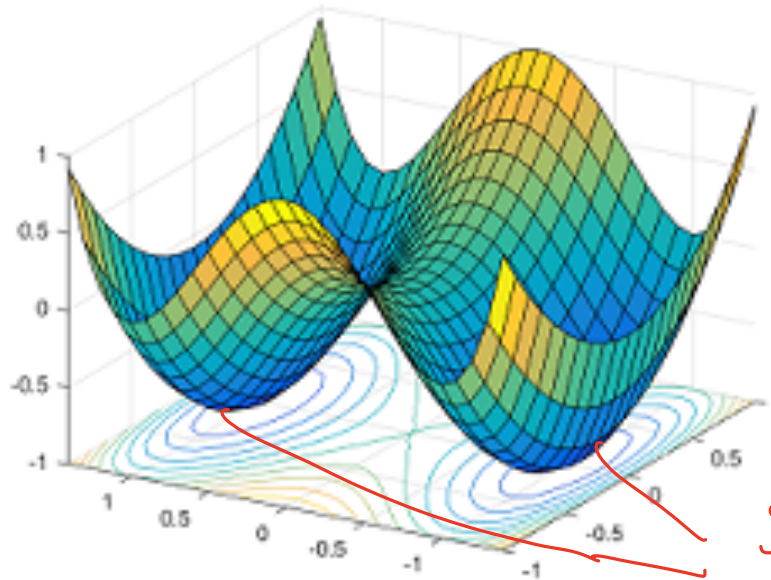
$$\min_x f(x)$$

- Stationary points:  $x : \nabla f(x) = 0$
- Global minimum:  
 $x : f(x) \leq f(x') \forall x' \in \mathbb{R}^d$
- Local minimum: *neighborhood*  
 $x : f(x) \leq f(x') \forall x' : \|x - x'\| \leq \epsilon$
- 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

---



same function  
value

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

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

$$\lambda_{\min}(A) < 0$$

choose  $v$

the eigenvector

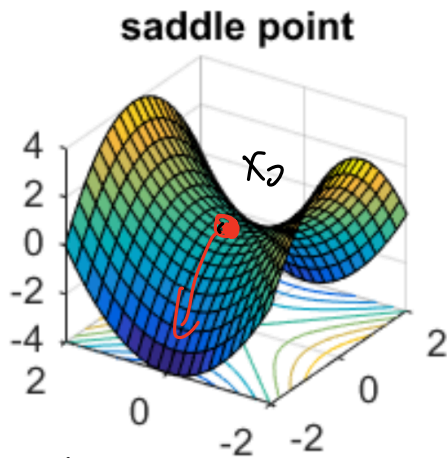
for  $\lambda_{\min}(A)$ ,  $\|v\|_2 = 1$

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

Quadratic

$$f(x) = \frac{1}{2} x^T A x$$

$x=0$  is a stationary point



$$f(x) \approx (x - x_0)^T \nabla f(x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0)$$

$$|v^T x_{t+1}| \geq (1 - \eta \lambda_{\min}(A))^t |v^T x_0|$$

$$x_0 \neq 0$$

$$\begin{aligned} & \|x_{t+1}\|_2 \geq \|v^T x_{t+1}\| \\ & |v^T x_{t+1}| = |v^T (x_t - \eta \nabla f(x_t))| \\ & = |v^T x_t - \eta \lambda_{\min}(A) v^T x_t| \\ & = |(1 - \eta \lambda_{\min}(A)) \cdot v^T x_t| \end{aligned}$$

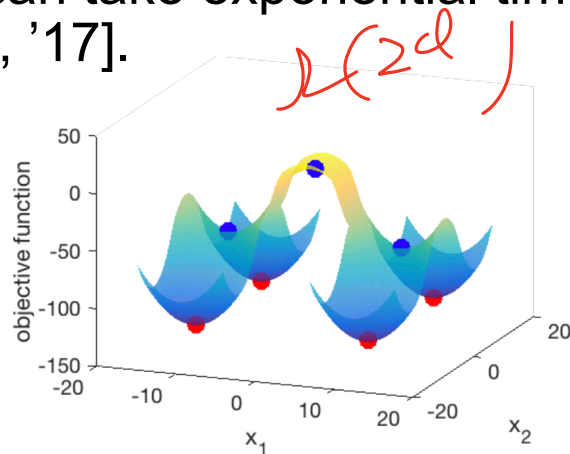


# Escaping Strict Saddle Points

$$X_{t+1} = X_t - \eta \nabla f(X_t) + \sqrt{\epsilon} \xi$$

- **Noise-injected** gradient descent can escape strict saddle points in polynomial time [Ge et al., '15, Jin et al., '17].  $\xi \sim \mathcal{N}(0, I)$
- Randomly initialized gradient descent can escape all strict saddle points asymptotically [Lee et al., '15].
  - Stable manifold theorem. *all points  $\rightarrow$  strict saddle*  
 $\mu(\downarrow) = 0$
- 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].
  - 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].