

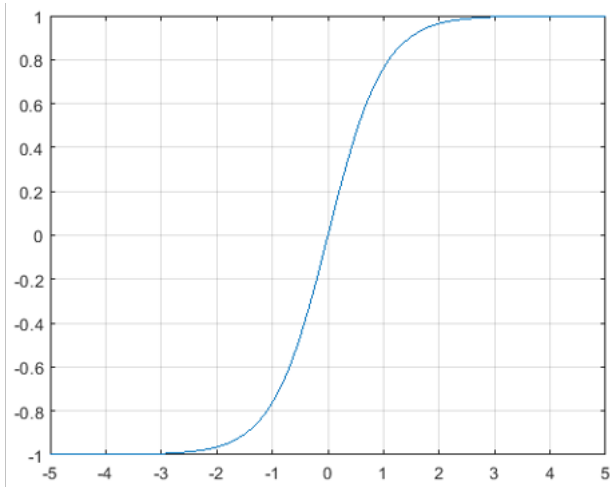
Important Techniques in Neural Network Training



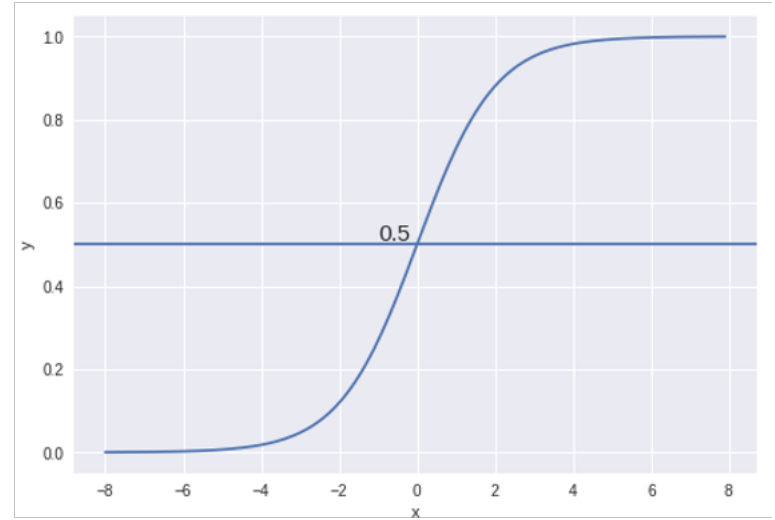
Gradient Explosion / Vanishing

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

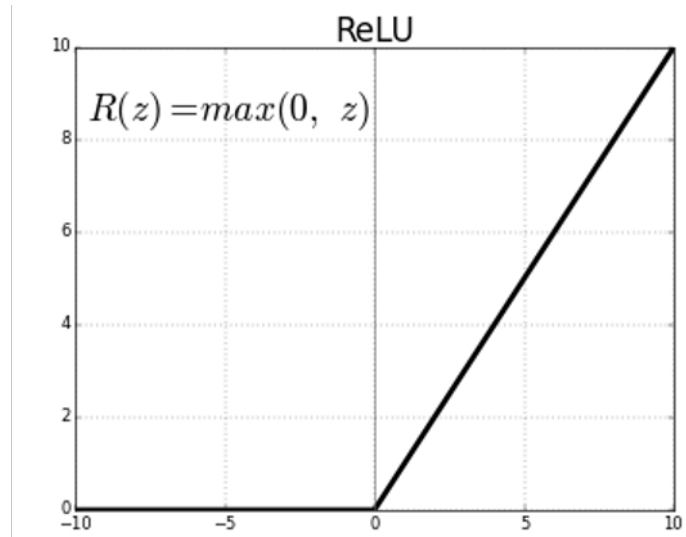
Activation Functions



tanh



sigmoid

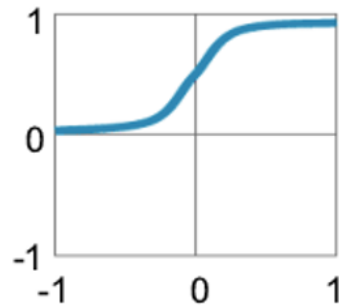


Rectified Linear Unit

Activation Function

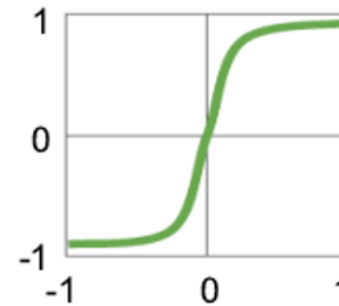
Traditional Non-Linear Activation Functions

Sigmoid



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

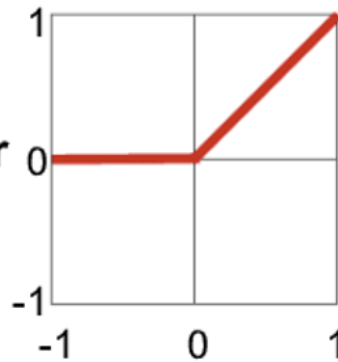
Hyperbolic Tangent



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

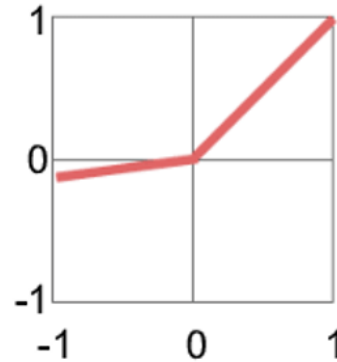
Modern Non-Linear Activation Functions

Rectified Linear Unit (ReLU)



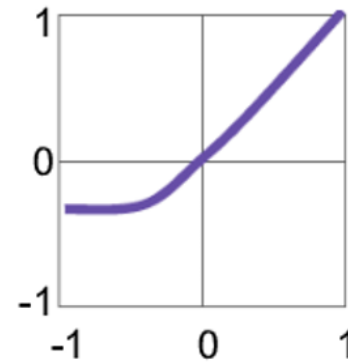
$$y = \max(0, x)$$

Leaky ReLU



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

Exponential LU



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

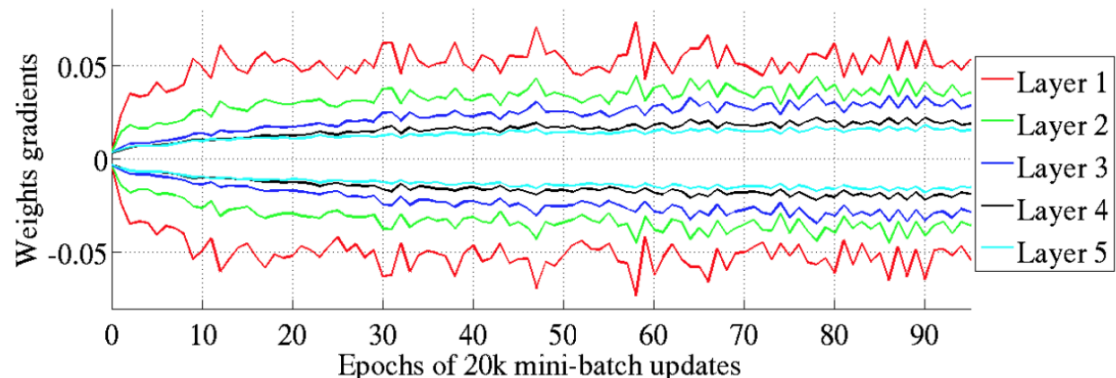
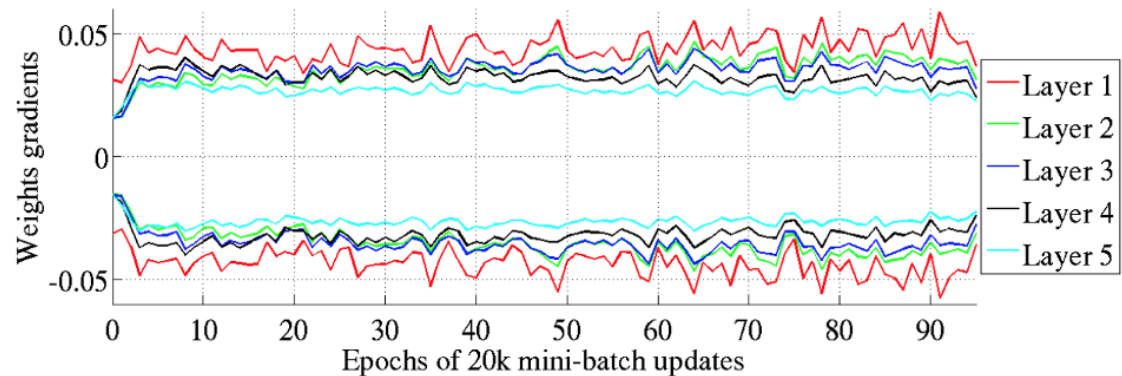
α = small const. (e.g. 0.1)

Initialization

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

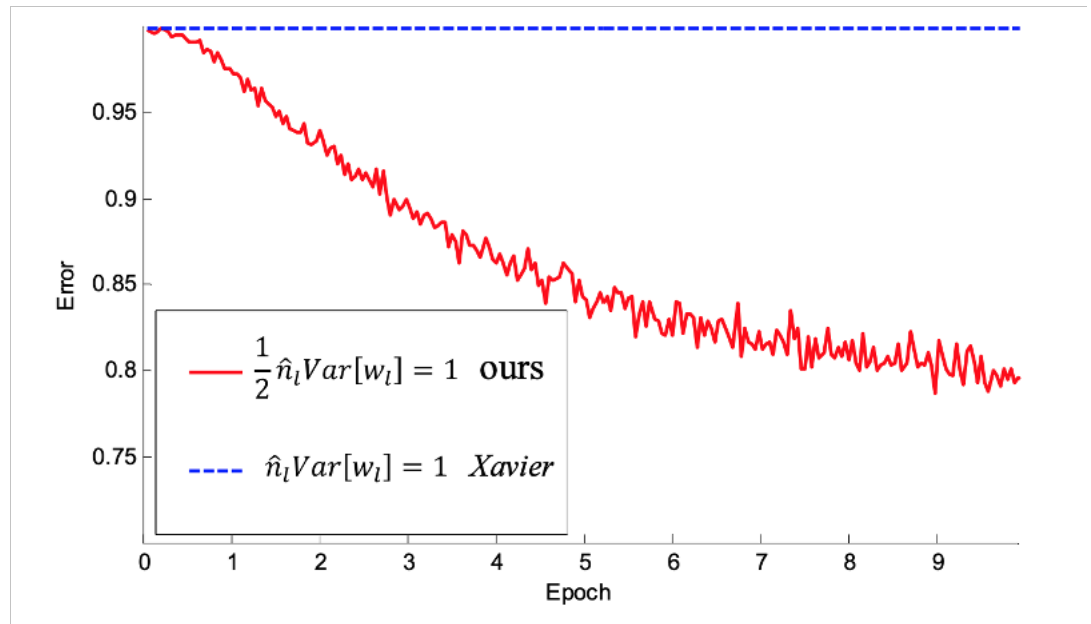
Xavier Initialization (Glorot & Bengio, '10)

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



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

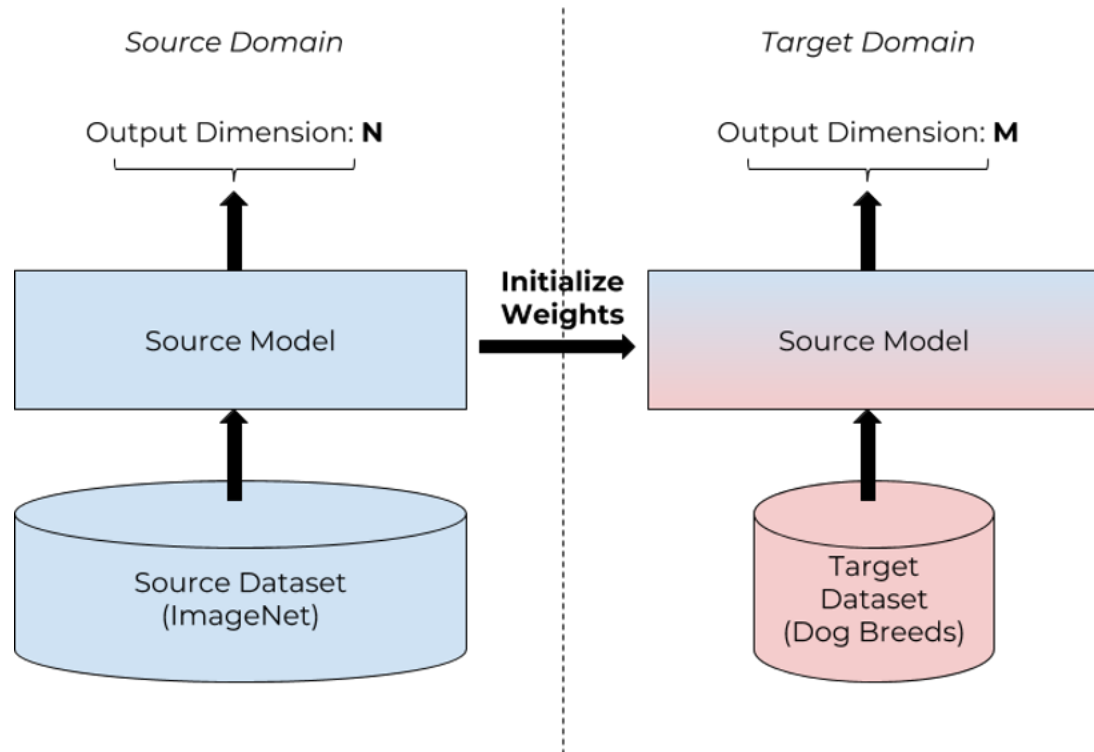


Kaiming Initialization (He et al. '15)

Kaiming Initialization (He et al. '15)

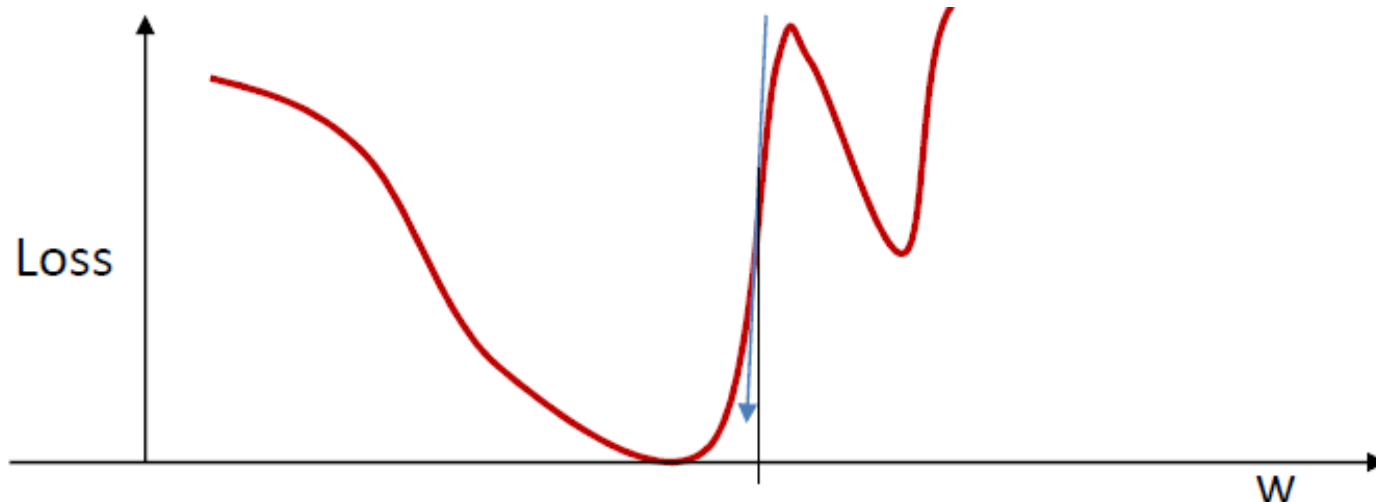
Initialization by Pre-training

- 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.

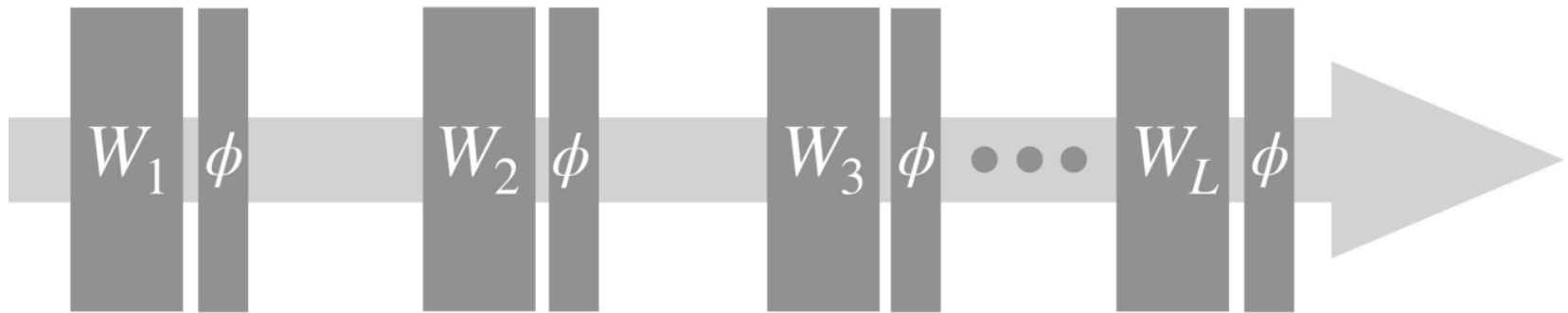


Batch Normalization (Ioffe & Szegedy, '14)

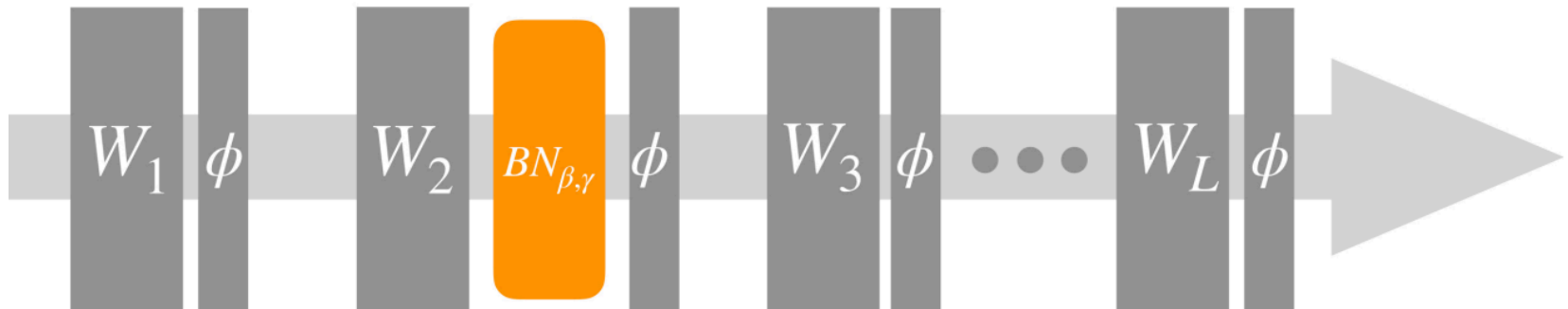
- **Normalizing/whitening** (mean = 0, variance = 1) the inputs is generally useful in machine learning.
 - Could normalization be useful at the level of hidden layers?
 - **Internal covariate shift**: the calculations of the neural networks change the distribution in hidden layers even if the inputs 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)

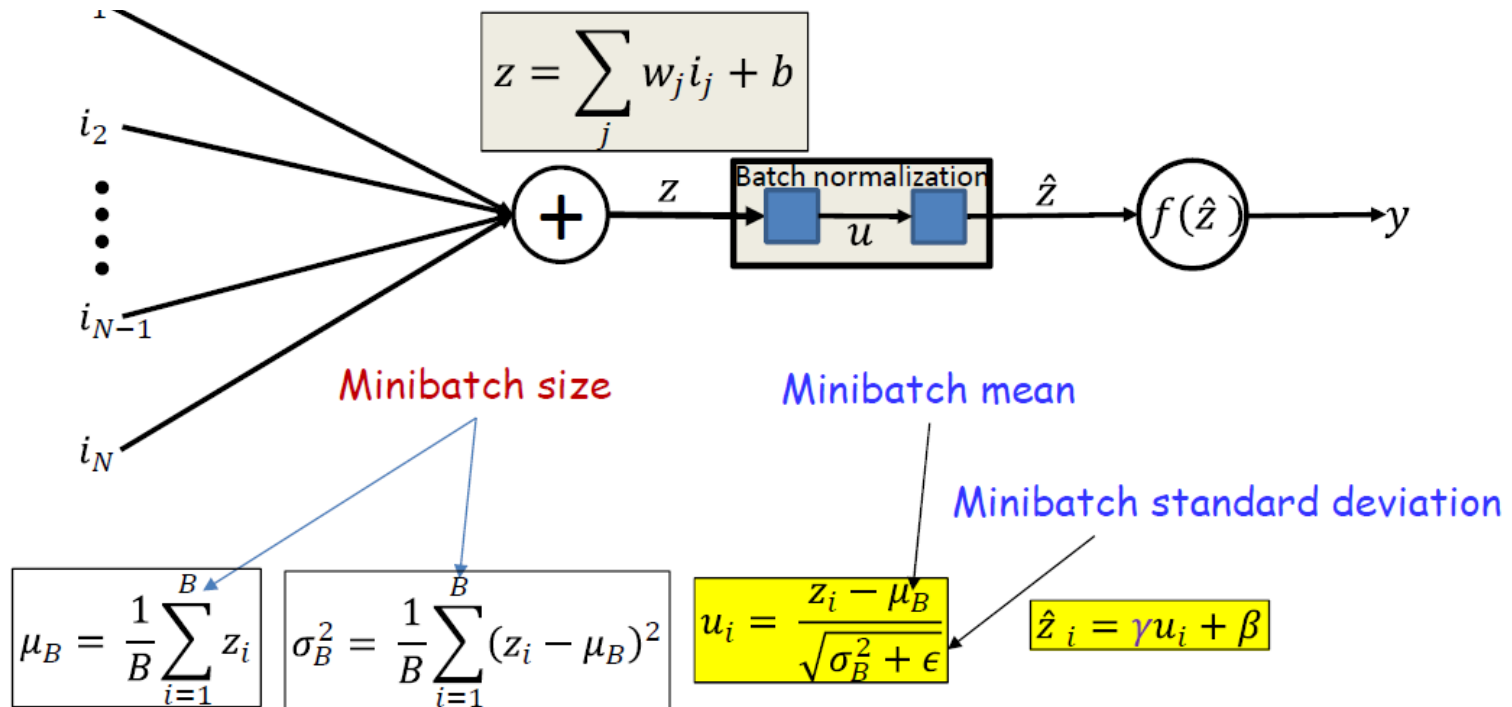
Standard Network



Adding a BatchNorm layer (between weights and activation function)



Batch Normalization (Ioffe & Szegedy, '14)

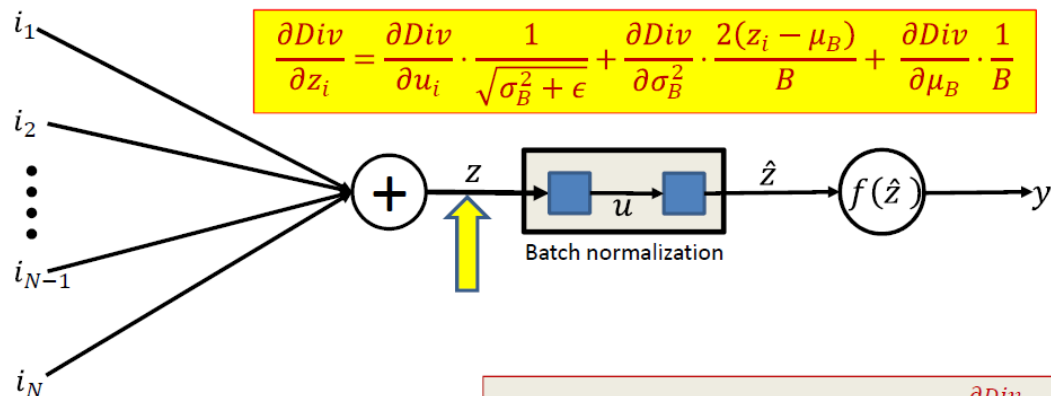


Batch Normalization (Ioffe & Szegedy, '14)

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

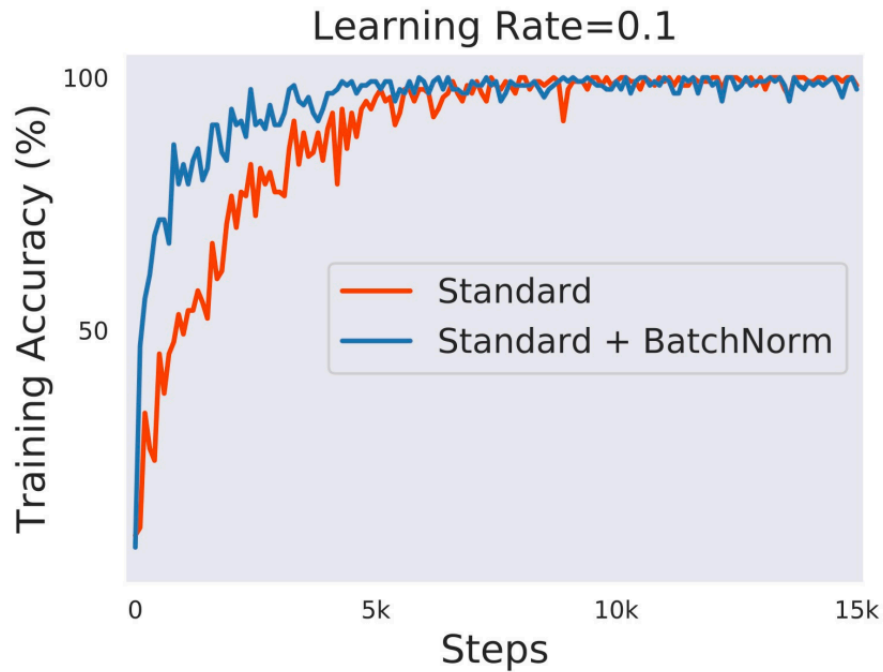
$$\frac{\partial Div}{\partial \sigma_B^2} = \frac{-1}{2} (\sigma_B^2 + \epsilon)^{-3/2} \sum_{i=1}^B \frac{\partial Div}{\partial u_i}$$

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



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 β, γ 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

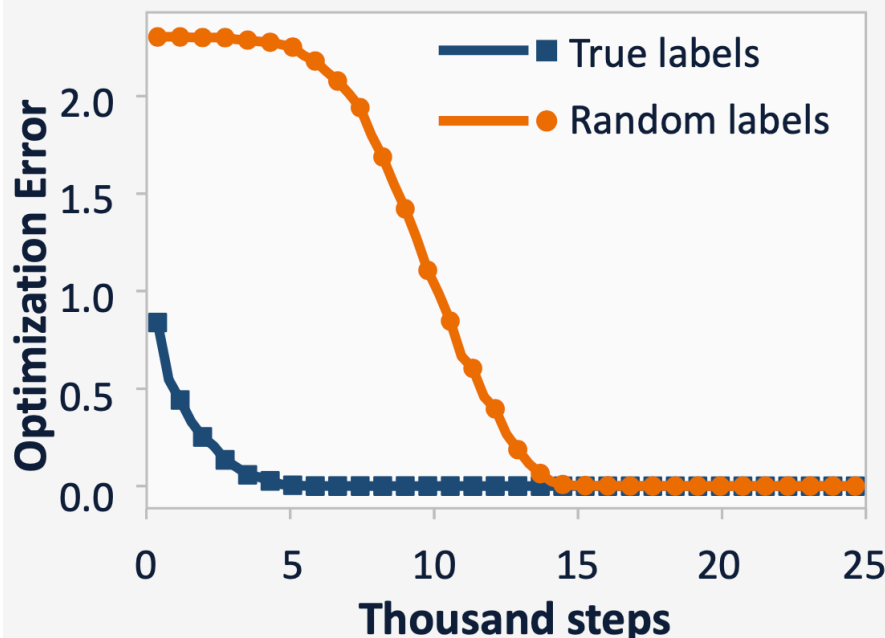
Non-convex Optimization Landscape

W

Gradient descent finds global minima

Practice: gradient descent

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



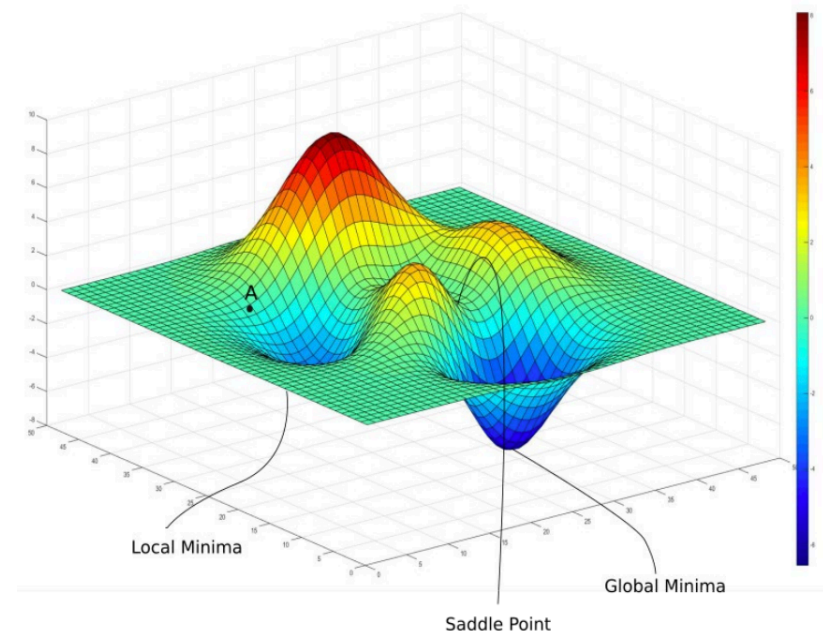
Optimization error $\rightarrow 0$ for both **true labels** and **random labels** !

Zhang Bengio Hardt Recht Vinyals 2017

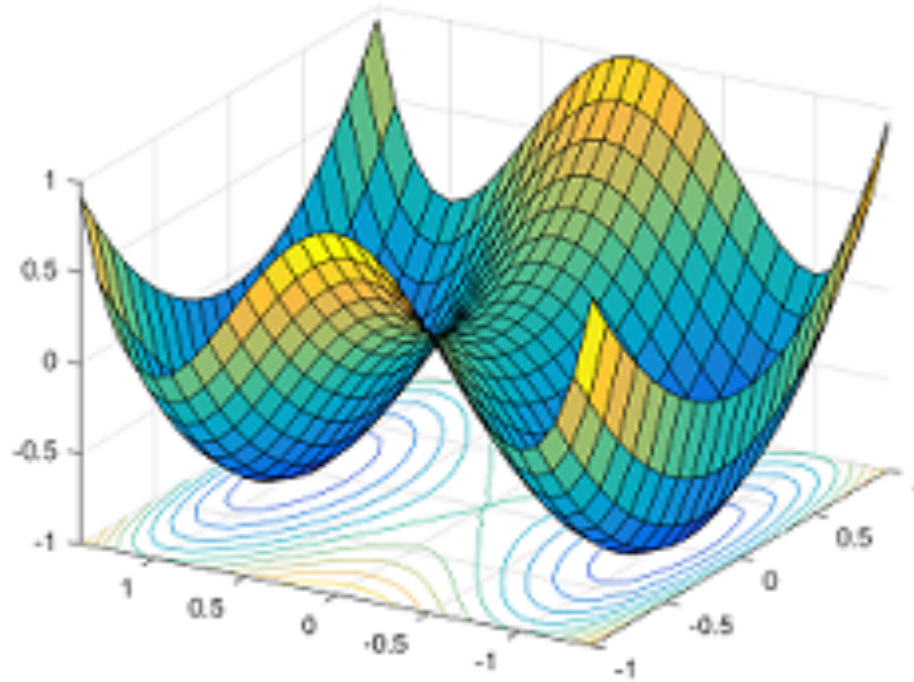
Understanding DL Requires Rethinking Generalization

Types of stationary points

- 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

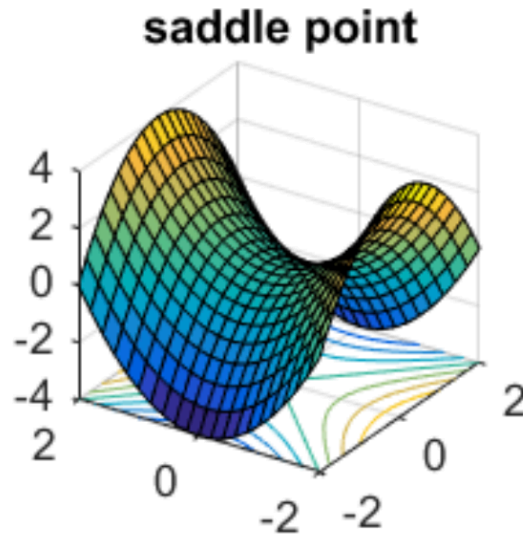


Landscape Analysis



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

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

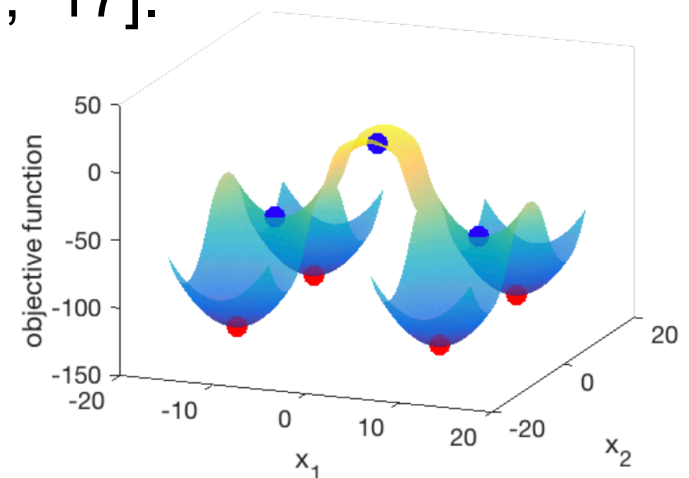


- Strict saddle point: a saddle point and $\lambda_{\min}(\nabla^2 f(x)) < 0$

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 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].