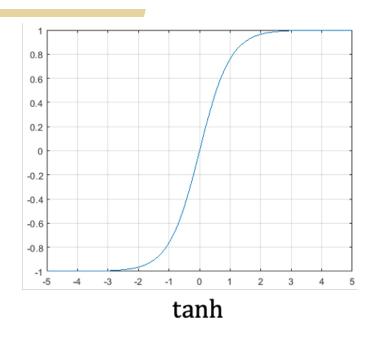
Important Techniques in Neural Network Training

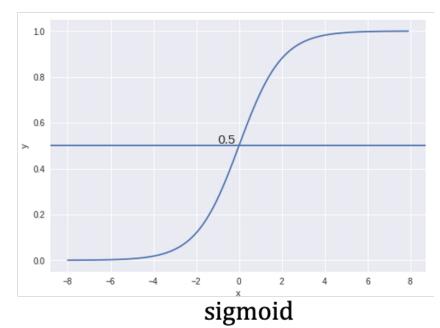


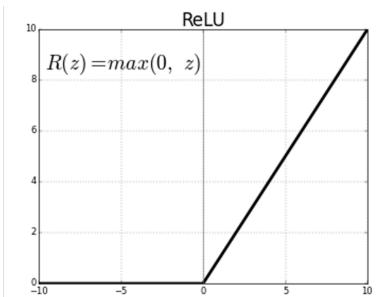
Gradient Explosion / Vanishing

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

Activation Functions

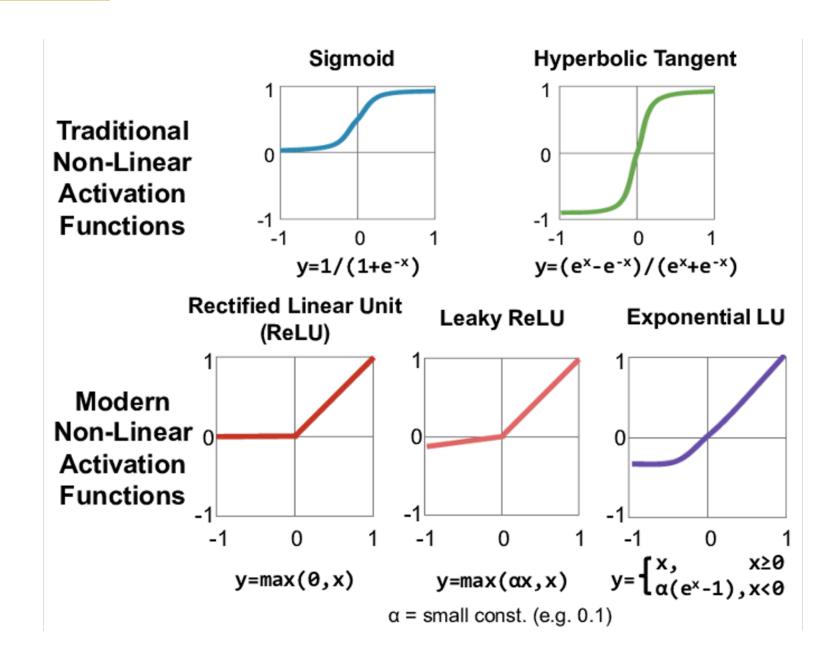






Rectified Linear United

Activation Function



Initialization

- Zero-initialization
- Large initialization
- Small initialization



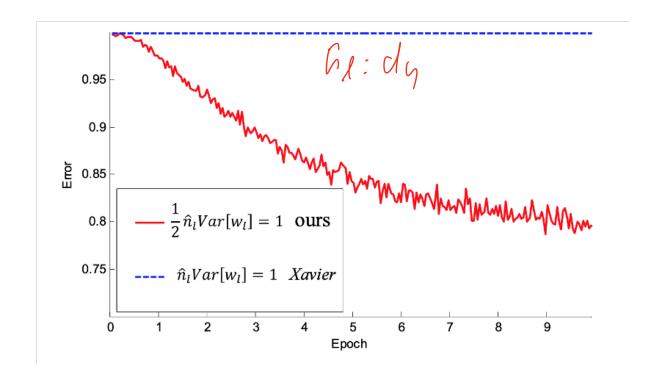
- Design principles:
 - Zero activation mean
 - Activation variance remains same across layers

Kaiming Initialization (He et al. '15) $W^{(b)} \subseteq \mathcal{D}^{d_{\text{MH}}} \land d_{\text{M}}$

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

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

- Designed for ReLU activation
- 30-layer neural network



Kaiming Initialization (He et al. '15)

Fuch layer:
$$2^h = W^h \times h$$
 | $ve-activation$
 $x^{h+1} = 6(2^h)$
 $x^h = \frac{1}{2^h} W_{ij}^h \cdot x_j^h$
 $x^h = \frac{1}{2^h} W_{ij}^h \cdot x_j^h$

Kaiming Initialization (He et al. '15)

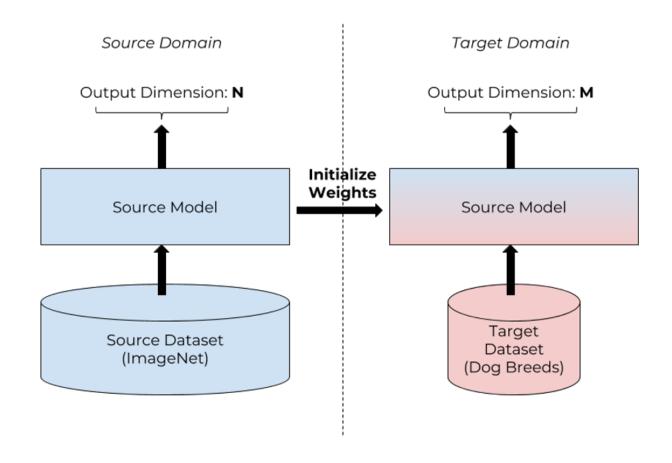
$$\begin{aligned}
& \text{If } \left((x_{i}^{h})^{2} \right) = \int_{-\infty}^{\infty} (x_{i}^{h})^{2} P(x_{i}^{h}) dx_{i}^{h} \\
& = \int_{-\infty}^{\infty} (x$$

Kaiming Initialization (He et al. '15)
$$Vav (2i) = du \cdot Vav (W_{ij}) \cdot \frac{1}{2} \quad Vav (2j)$$

$$= \int_{2}^{4} du \cdot Vav (W_{ij}) = \left[Vav (2i) \cdot \frac{1}{2} \cdot \frac{1}$$

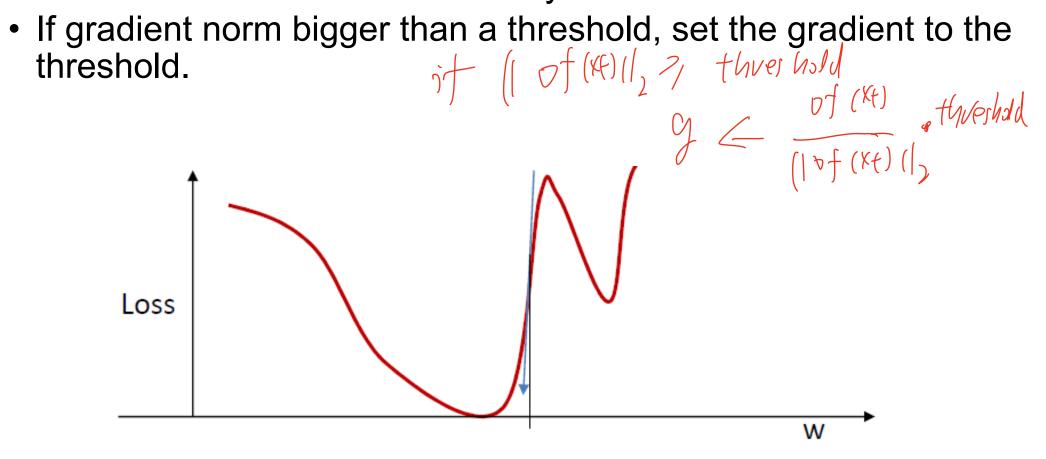
Initialization by Pre-training

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



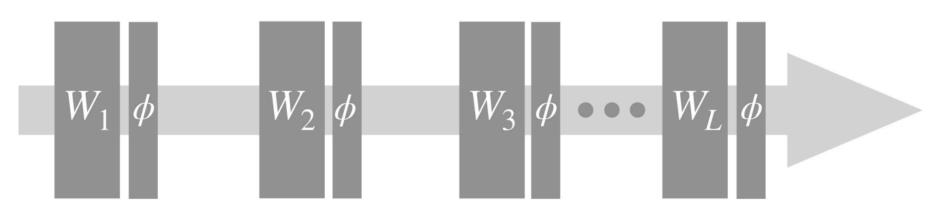
Gradient Clipping

- The loss can occasionally lead to a steep descent
- This result in immediate instability

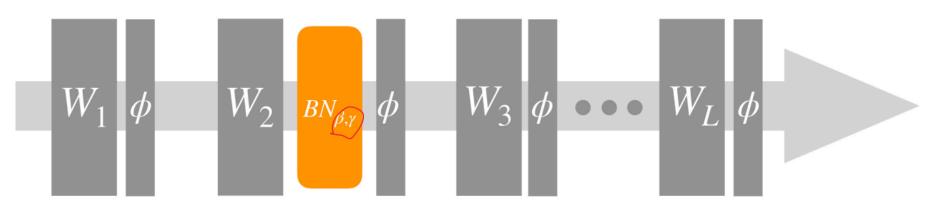


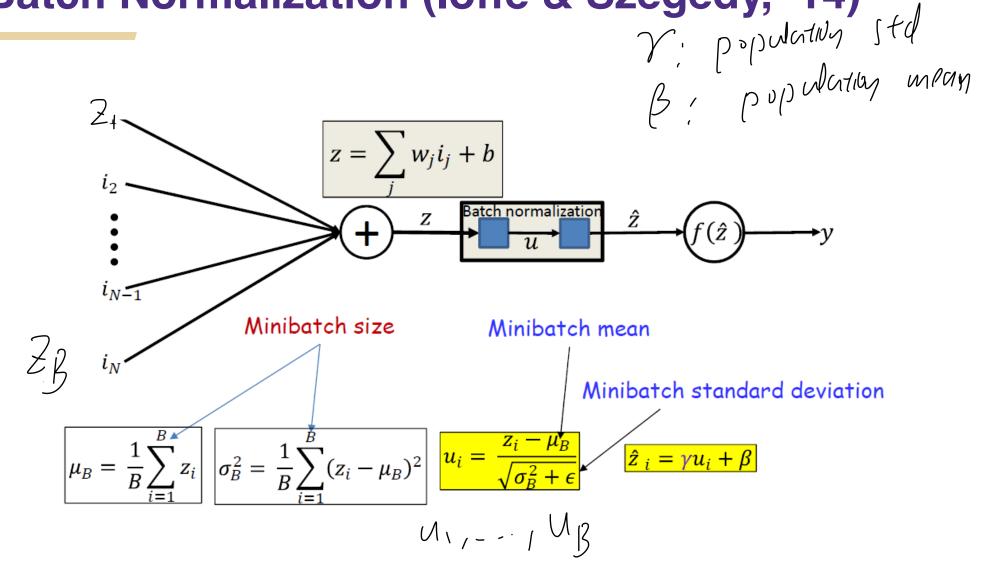
- Normalizing/whitening (mean = 0, variance = 1) the inputs is generally useful in machine learning.
 - Could normalization be useful at the level of hidden layers?
 - Internal covariate shift: the calculations of the neural networks change the distribution in hidden layers even if the inputs are normalized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

Standard Network

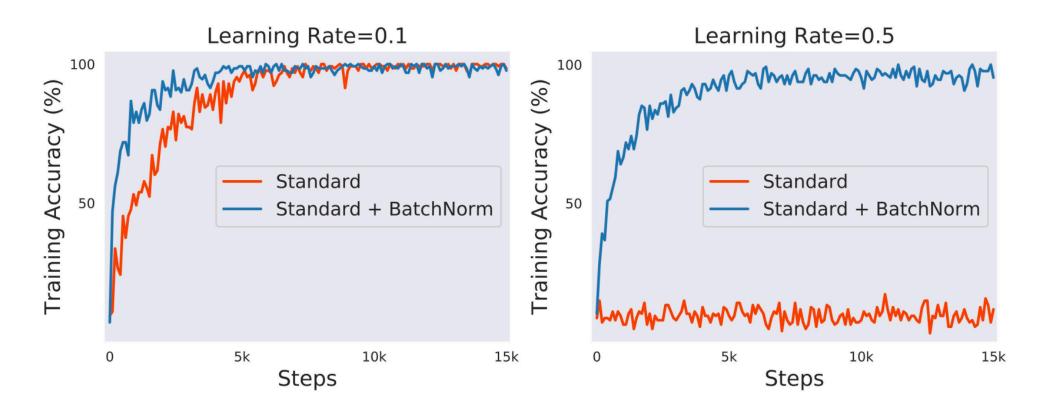


Adding a BatchNorm layer (between weights and activation function)





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



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 nontrivial performance (Frankle et al. '20)
- BN can use exponentially increasing learning rate! (Li & Arora '19)

More normalizations

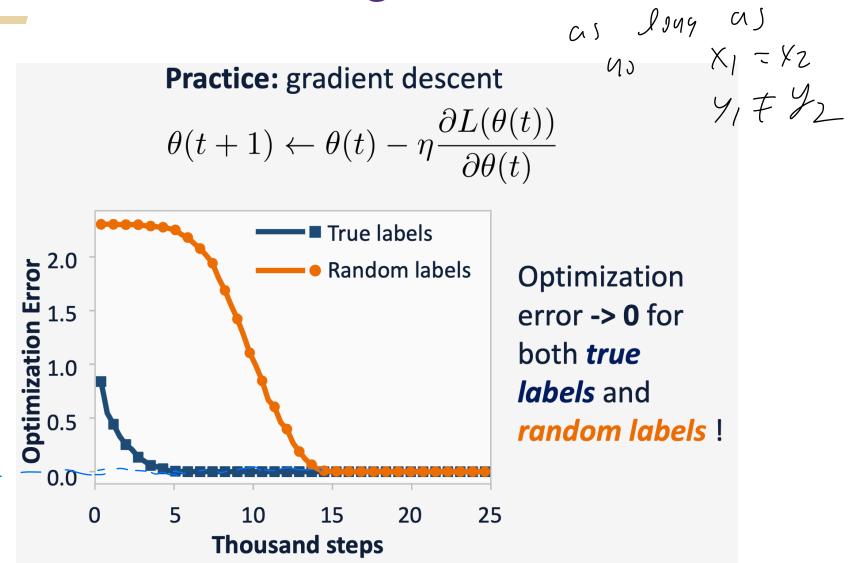
- Layer normalization (Ba, Kiros, Hinton, '16)
 - Batch-independent
 - Suitable for RNN, MLP
- Weight normalization (Salimans, Kingma, '16)
 - Suitable for meta-learning (higher order gradients are needed)

•

Non-convex Optimization Landscape



Gradient descent finds global minima



Zhang Bengio Hardt Recht Vinyals 2017 Understanding DL Requires Rethinking Generalization

Types of stationary points

- Stationary points: $x : \nabla f(x) = 0$
- Global minimum:

$$x: f(x) \le f(x') \, \forall x' \in \mathbb{R}^d$$

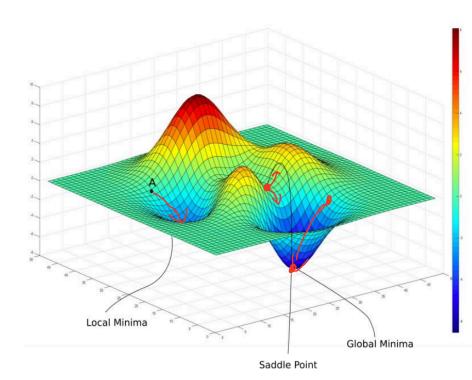
• Local minimum:

$$x: f(x) \le f(x') \forall x': ||x - x'|| \le \epsilon$$

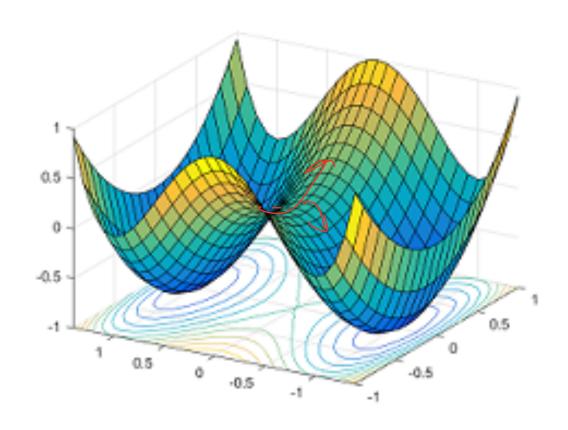
Local maximum:

$$x: f(x) \ge f(x') \forall x': ||x - x'|| \le \epsilon$$

 Saddle points: stationary points that are not a local min/max

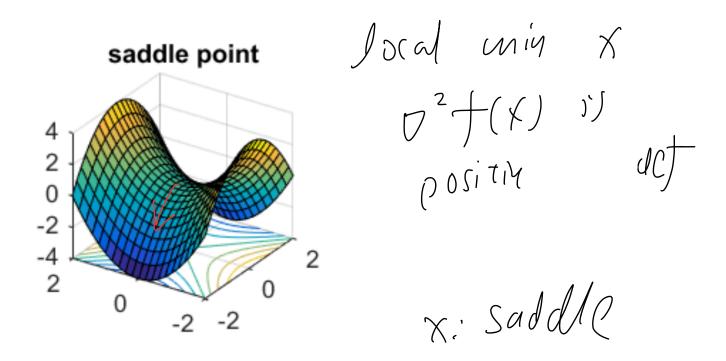


Landscape Analysis



- 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 $\lambda_{\min}(\nabla^2 f(x)) < 0$

$$m_{1} = \sqrt{X^{7}AX}$$

$$\sqrt{X^{7}X^{7}} = \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}}$$

$$= \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}}$$

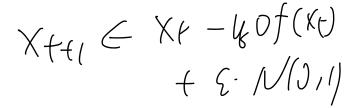
$$= \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}}$$

$$= \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^{7}}$$

$$= \sqrt{X^{7}X^{7}} + \sqrt{X^{7}X^$$

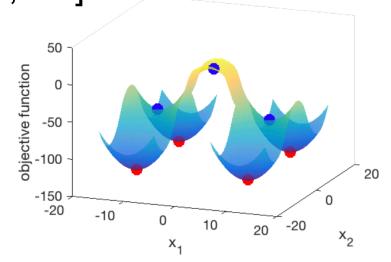
A D J: Pigen UP(t) (VIVO) PONDS +> Smalle)t Pigen Udap

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].