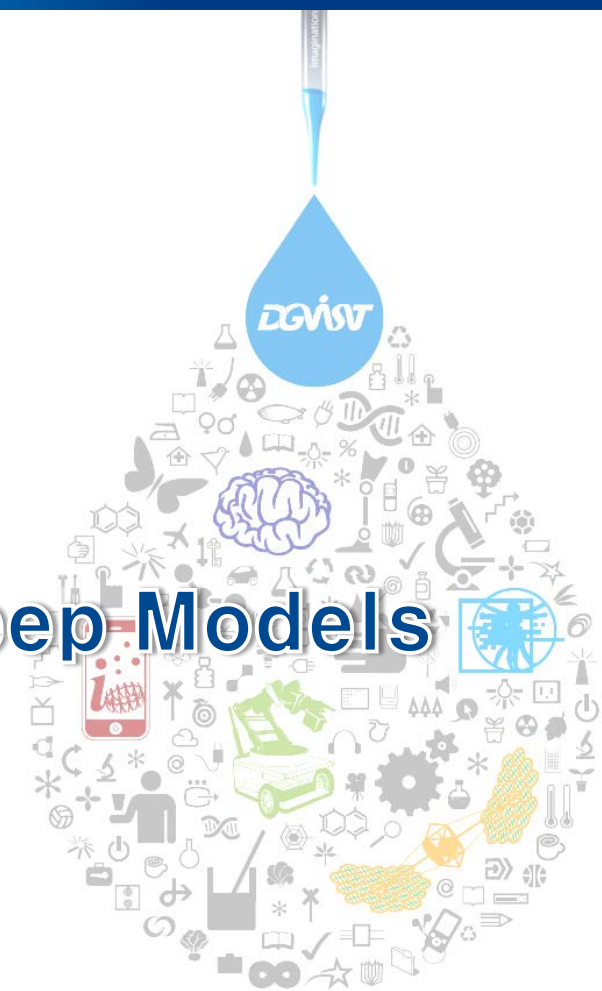


Chapter 8. Optimization for Training Deep Models

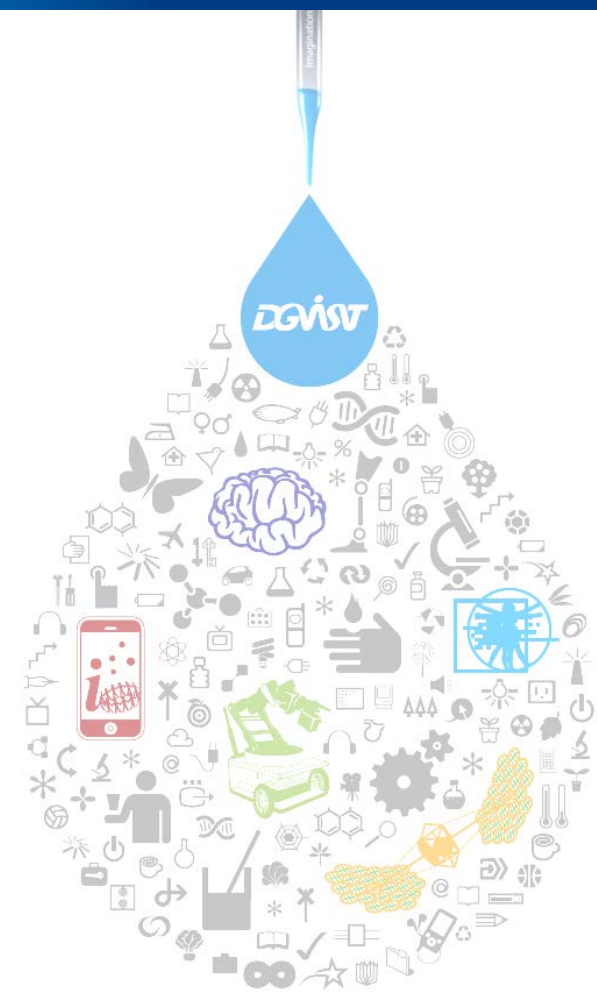
2017-08-31
Jinwook Kim



Contents

- 8.1 How Learning Differs from Pure Optimization**
- 8.2 Challenges in Neural Networks Optimization**
- 8.3 Basic Algorithms**
- 8.5 Algorithms with Adaptive Learning Rates**
- 8.4 Parameter Initialization Strategies**
- 8.6 Approximate Second-order Methods**
- 8.7 Optimization Strategies and Meta-algorithms**

How Learning Differs from Pure Optimization



Empirical Risk Minimization

- Machine learning usually acts **indirectly**
- The goal of learning is to reduce the expected generalization error (risk)
- However, learning algorithms reduce cost functions (empirical risk)
 - minimizing the expected loss on the **training dataset**
 - in the hope that the **indirect** optimization will improve the performance

Surrogate Loss Functions

■ Sometimes the loss function is difficult to be optimized

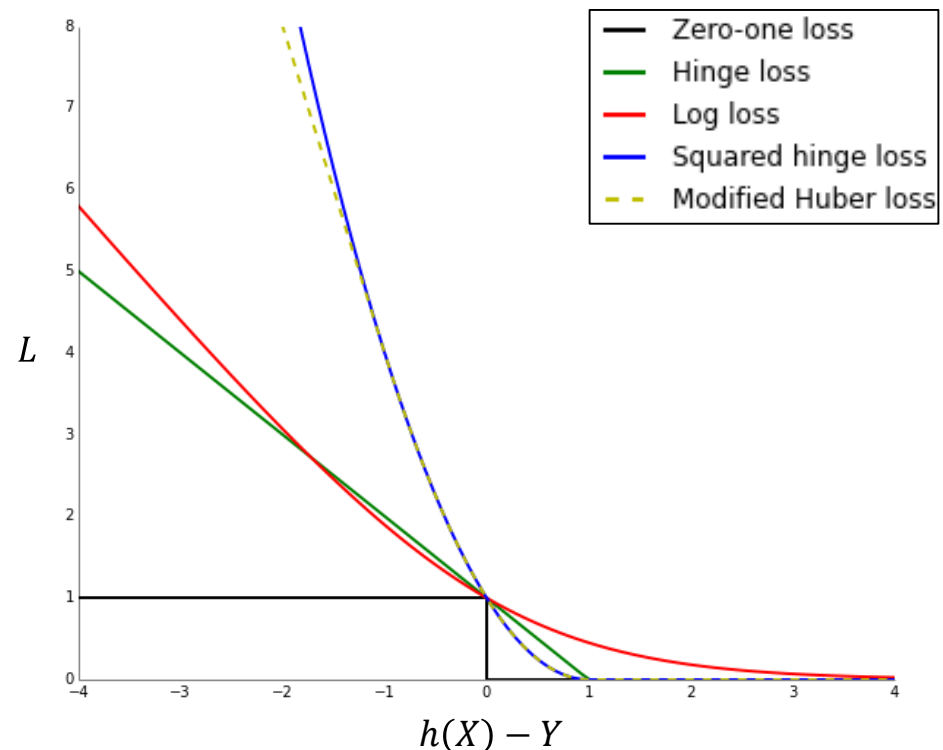
➤ e.g., 0-1 loss

$$L(h(\mathbf{x}), Y) = \begin{cases} 1 & h(\mathbf{x}) < Y \\ 0 & \text{otherwise} \end{cases}$$

■ In the case, we can use a surrogate loss for the optimization

➤ differentiable

➤ improve robustness



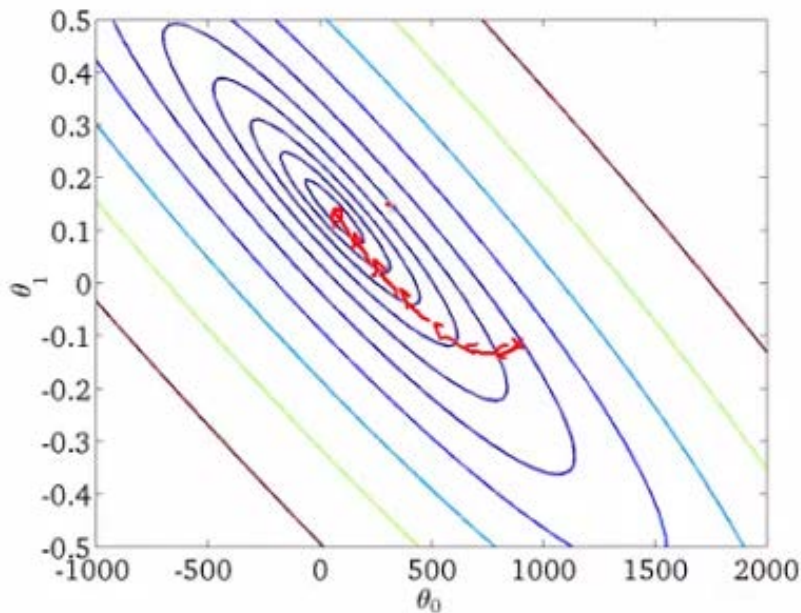
P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe, "Convexity , Classification , and Risk Bounds," J. Am. Stat. Assoc., pp. 1–36, 2003.

Minibatch Algorithms

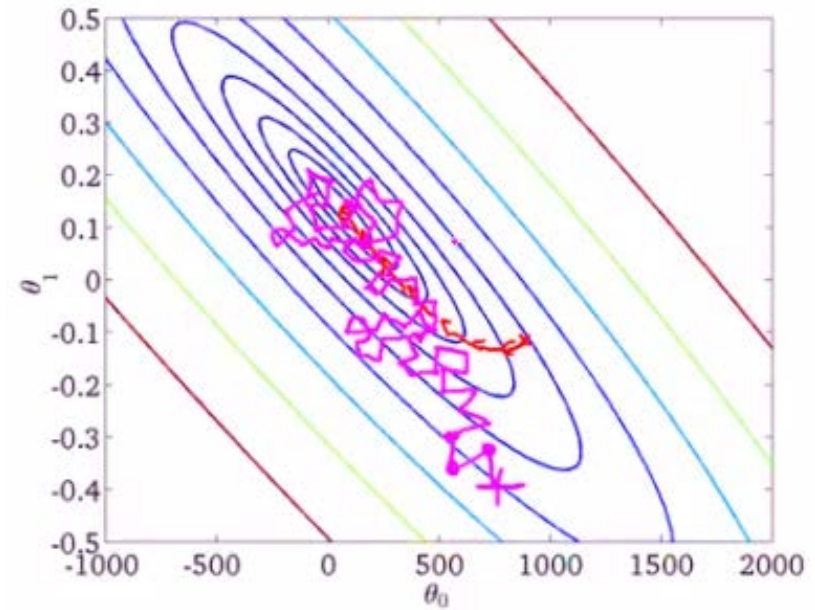
- Optimization in machine learning typically computes loss and updates parameters **iteratively**
 - e.g., gradient descent
- Using the entire training set (batch algorithm) per every iteration is computationally infeasible
- Minibatch algorithm uses only b examples in each iteration
 - m : the size of entire training examples
 - b : minibatch size ($1 \leq b < m$)

Batch vs. Minibatch

- Each iteration in minibatch algorithm may have poor optimization performance than batch algorithm
- However, after many iterations, the minibatch algorithm generally converges to optimal state



(a) Batch gradient descent optimization



(b) Minibatch gradient descent optimization

http://www.holehouse.org/mlclass/17_Large_Scale_Machine_Learning.html

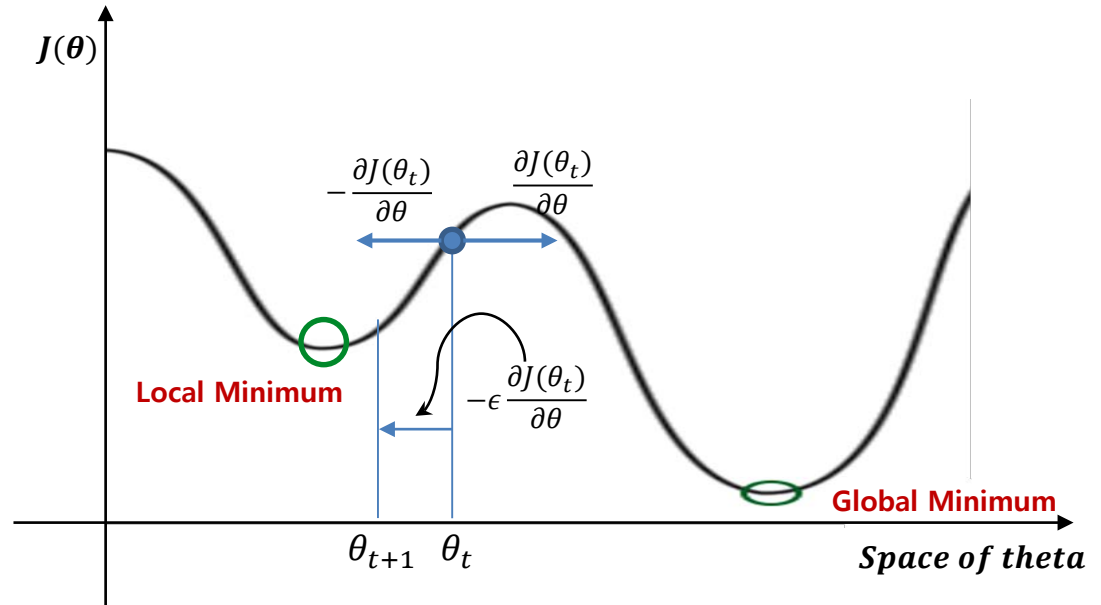
Limitation of Gradient Descent

- Issue of local minimum

Objective: $\min J(\theta)$

$$\theta_{t+1} = \theta_t - \epsilon \frac{\partial J(\theta)}{\partial \theta}$$

(ϵ : Learning rate)



- If the starting point for gradient descent was chosen inappropriately, cannot reach global minimum

Saddle Point

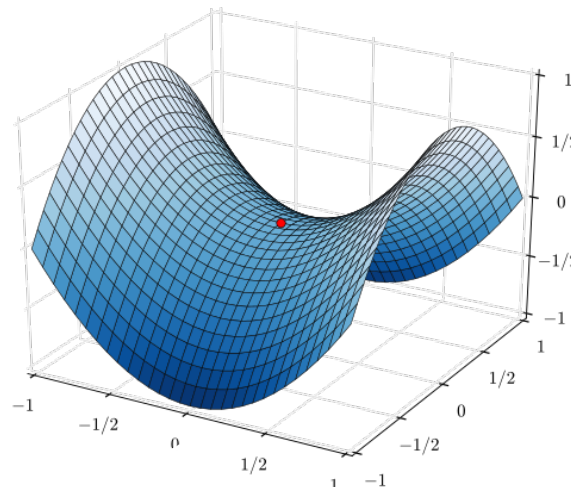
- For many high-dimensional non-convex functions, local minima are in fact rare cases

n : # of parameters in model

➤ local minima should satisfy,

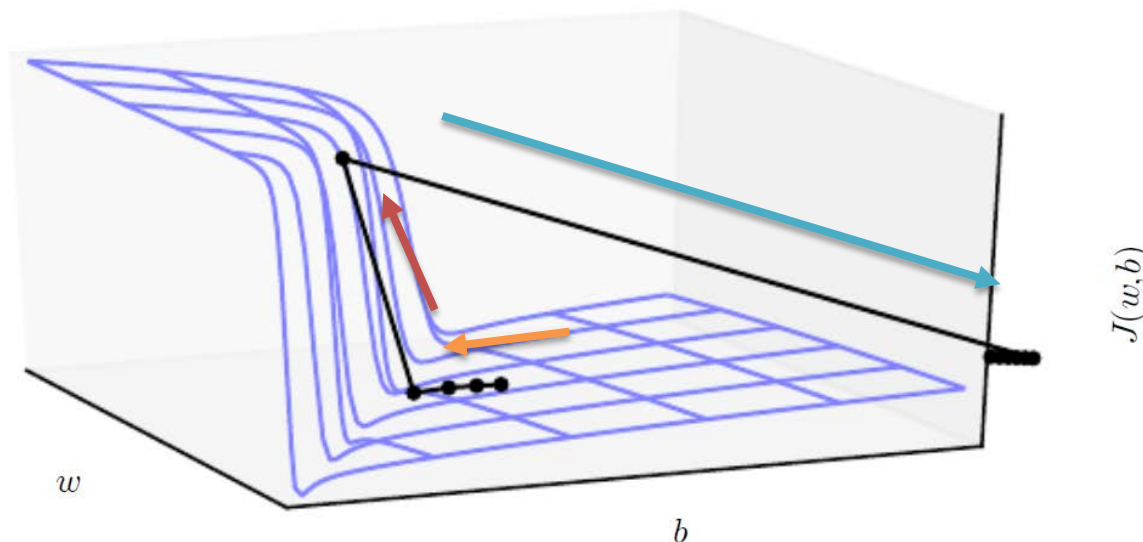
$$\frac{\partial J}{\partial \theta_i} = 0, i = 1, \dots, n$$
$$\frac{\partial^2 J}{\partial^2 \theta_i} > 0, i = 1, \dots, n$$

- Most of zero-gradient points in deep neural networks are saddle points



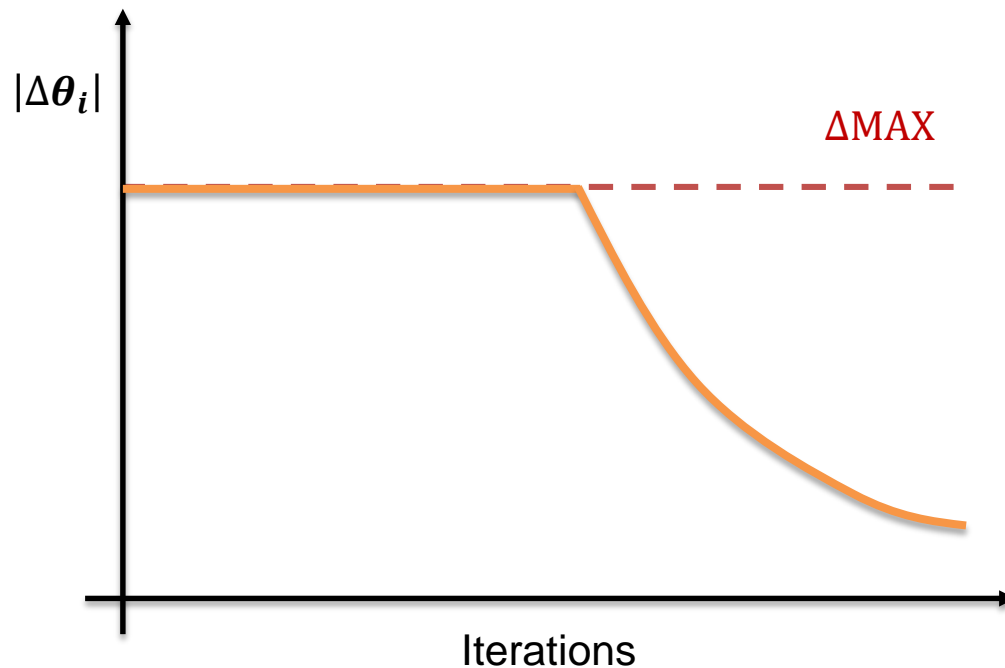
Cliffs and Exploding Gradient

- **Cliff** is one common issue in the objective functions for deep neural networks
 - resulting from the multiplication of several parameters
- **The gradient update step can move the parameter extremely far (exploding gradient)**
 - losing most of optimization works that has been done

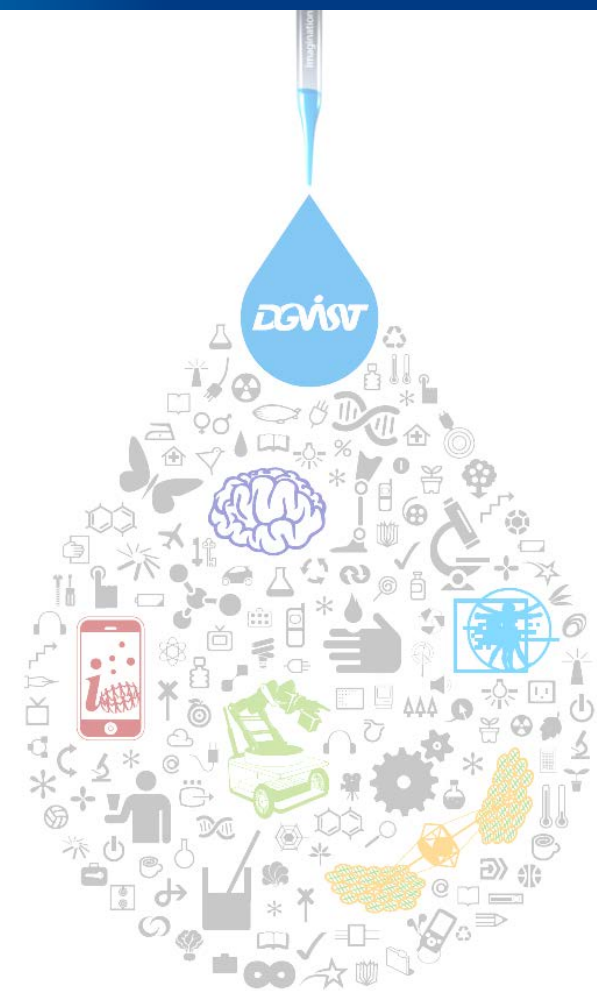


Gradient Clipping for Handling Cliffs

$$\Delta\theta_i \leftarrow \Delta\theta_i \frac{\Delta\text{MAX}}{\max(|\Delta\theta_i|, \Delta\text{MAX})}$$

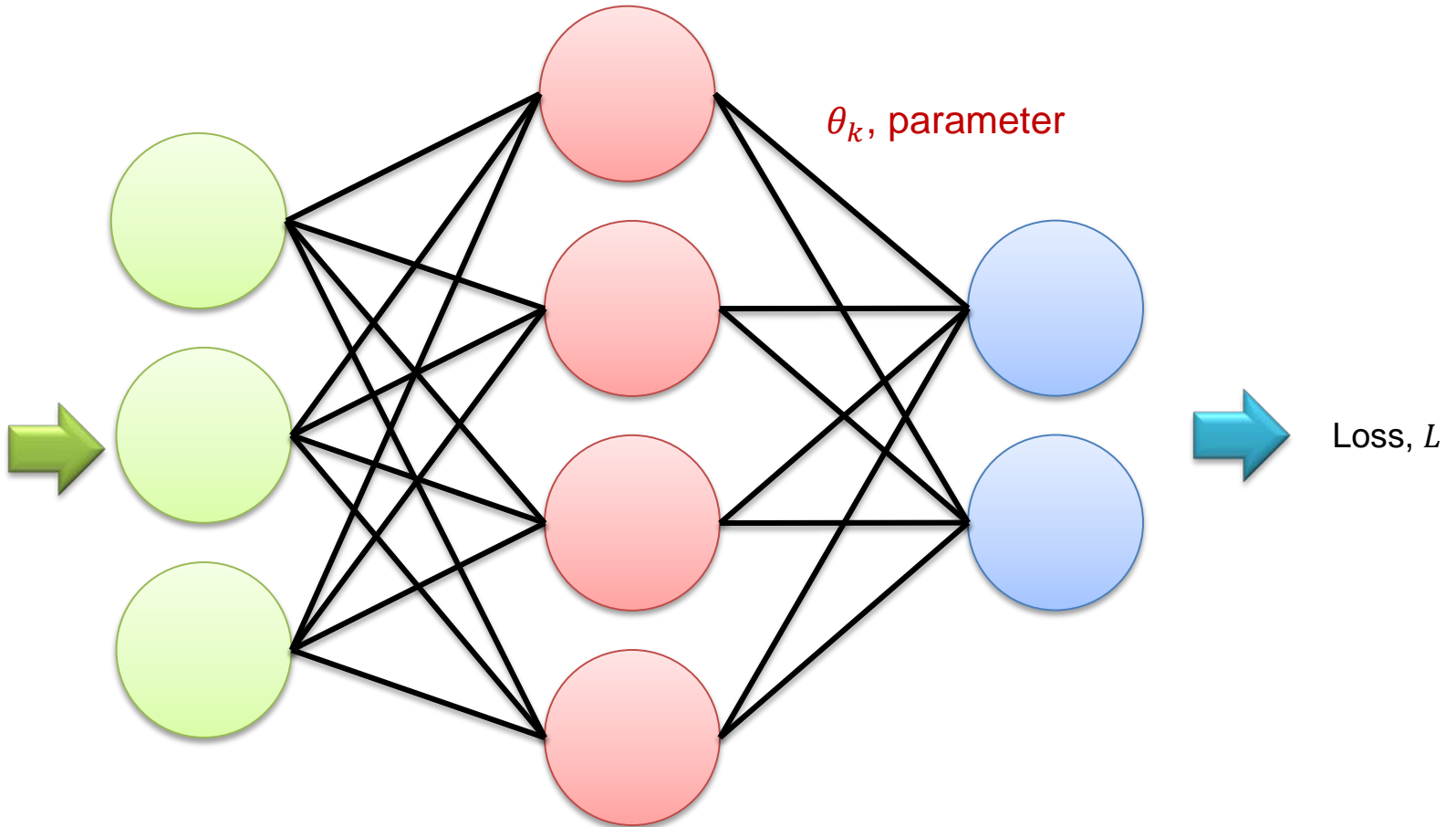


Basic Algorithms



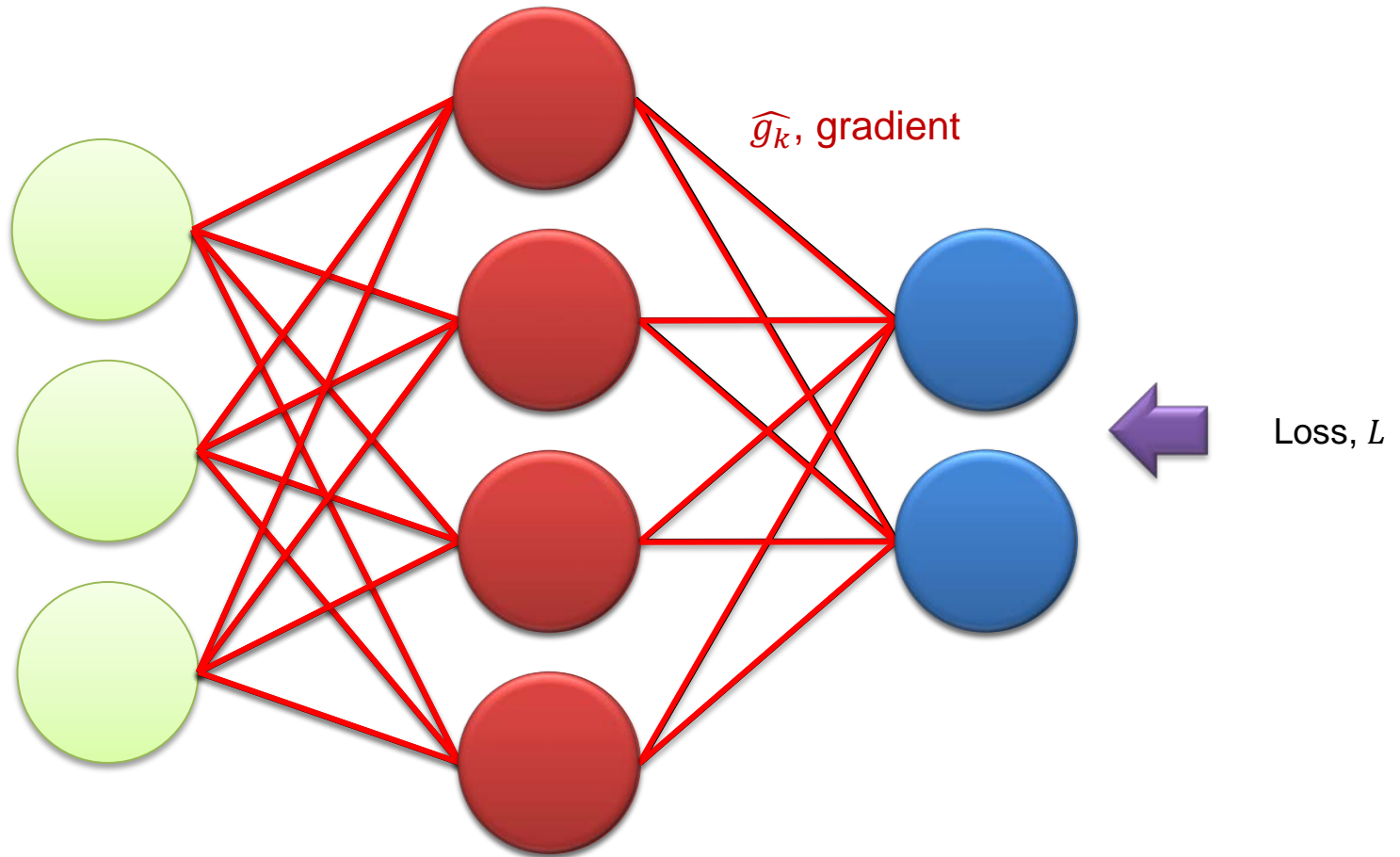
Forward Propagation (Review)

Data
(minibatch)



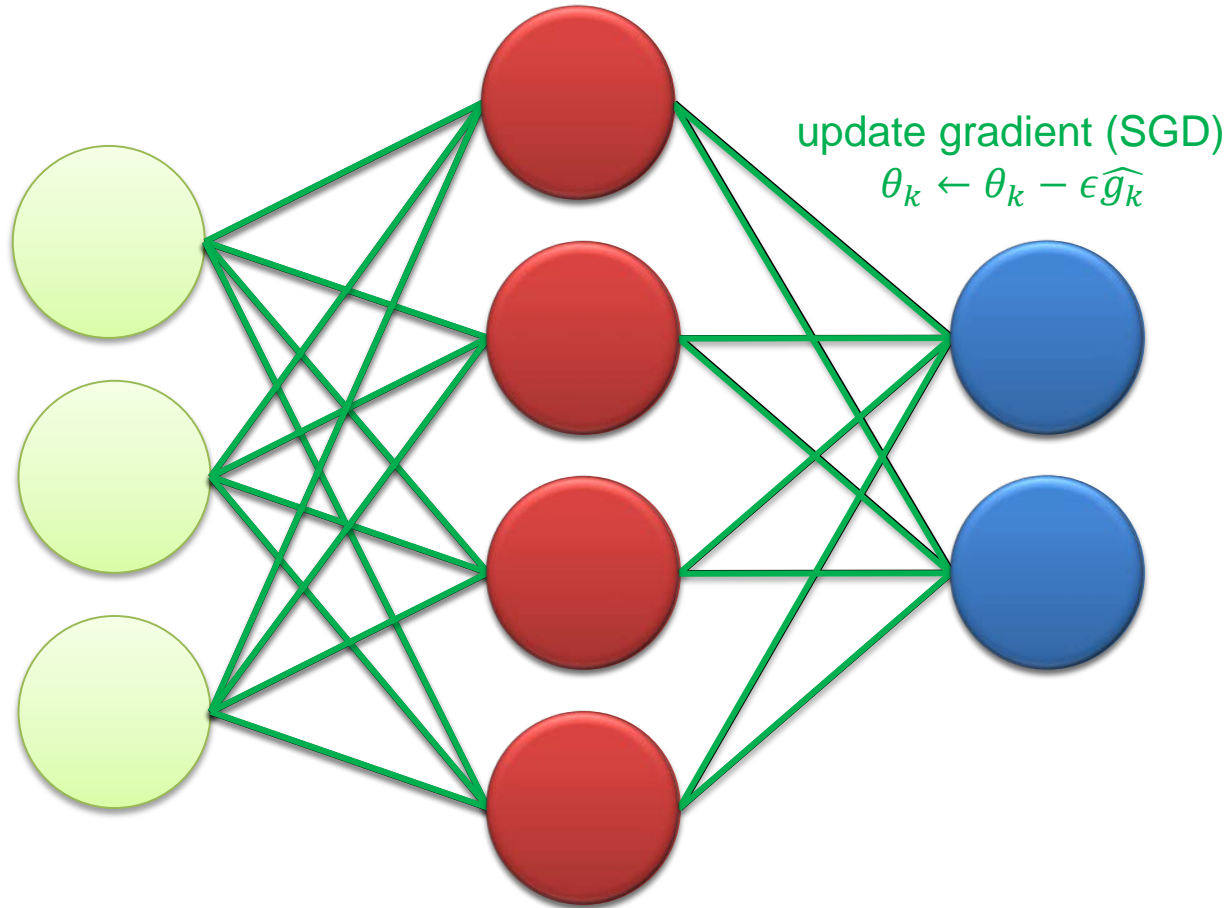
Back Propagation (Review)

Data
(minibatch)



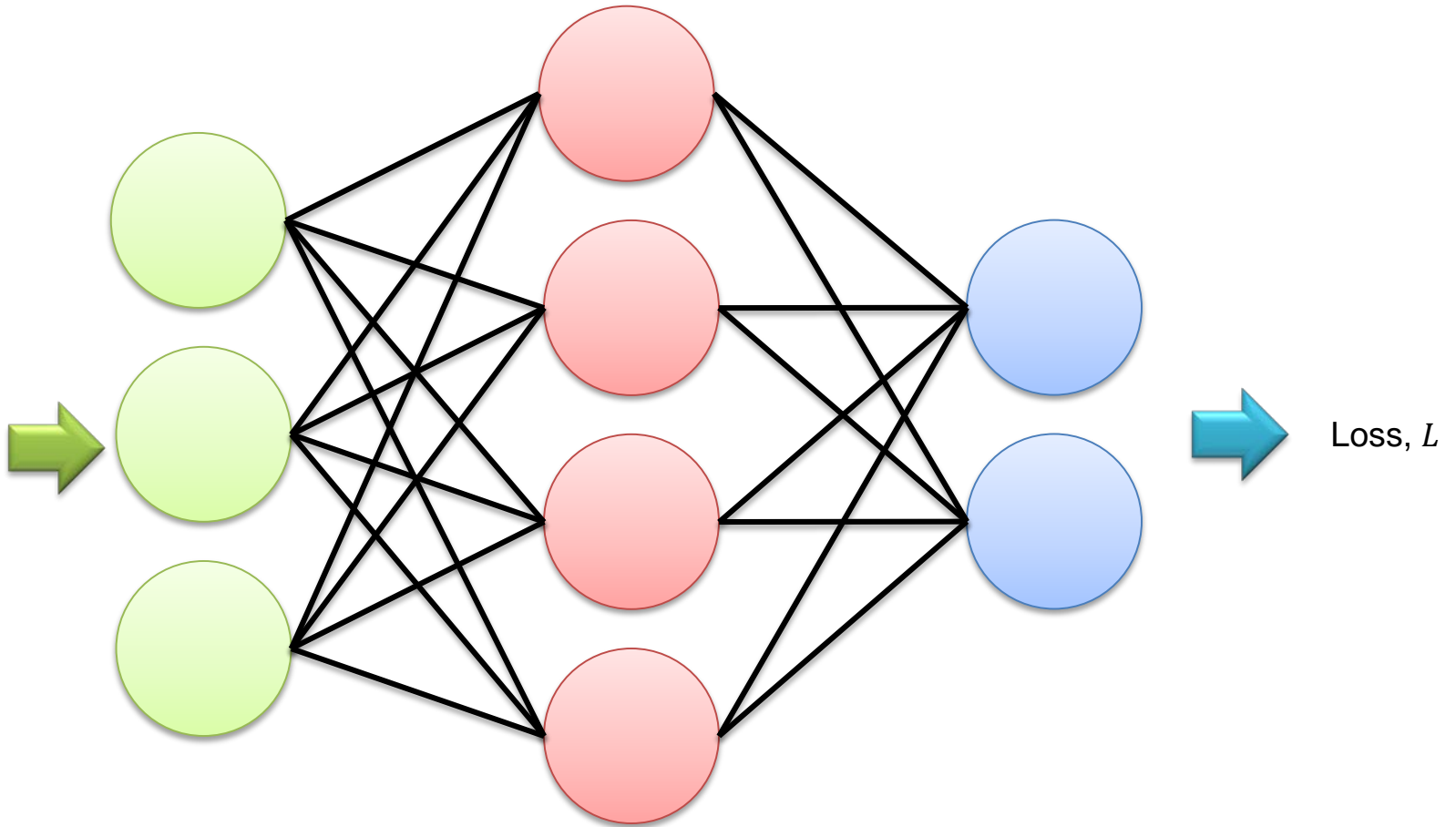
Parameter Update (Review)

Data
(minibatch)



Next Iteration

Data
(minibatch)



Stochastic Gradient Descent (SGD)

- Minibatch of the training set:

data $\{x^{(1)}, \dots, x^{(m)}\}$ with targets $\{y^{(1)}, \dots, y^{(m)}\}$

- Gradient:

$$\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(x^{(i)}; \theta), y^{(i)})$$

- Apply update:

$$\theta \leftarrow \theta - \epsilon \hat{g}$$

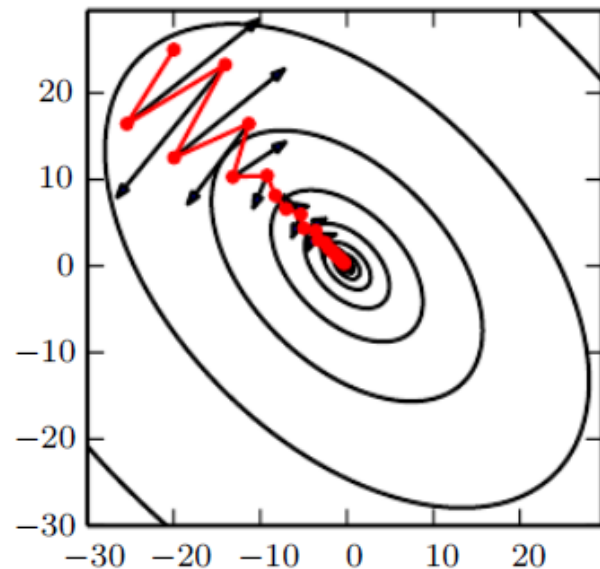
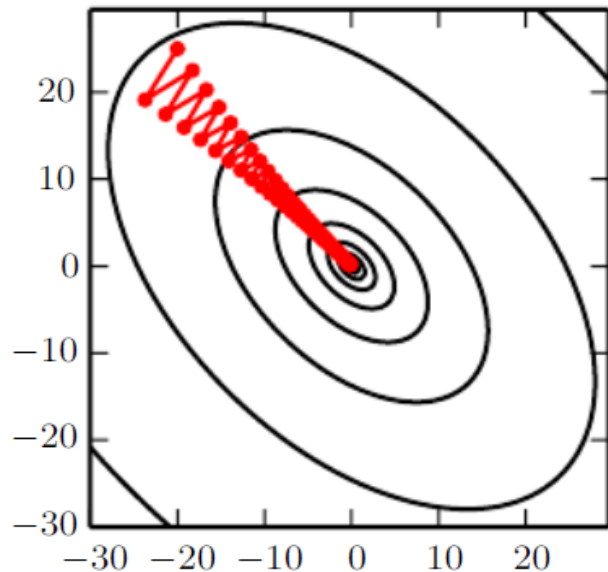
➤ learning rate ϵ is a critical hyperparameter for SGD

Momentum

■ Momentum is designed to accelerate learning

- accumulates an exponentially decaying moving-average of **past gradients**
- continues to move in their direction

■ Comparison of GD and GD with momentum:



SGD with Momentum

■ Minibatch of the training set:

data $\{x^{(1)}, \dots, x^{(m)}\}$ with targets $\{y^{(1)}, \dots, y^{(m)}\}$

■ Gradient:

$$\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(x^{(i)}; \theta), y^{(i)})$$

■ Accumulate velocity:

$$v \leftarrow \alpha v - \epsilon \hat{g}$$

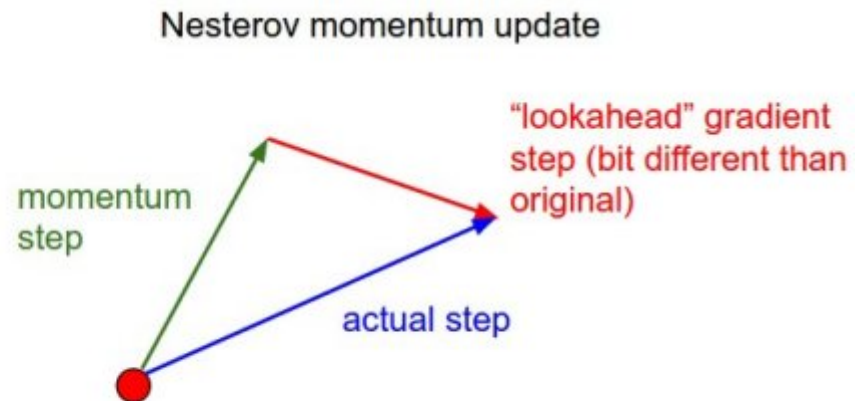
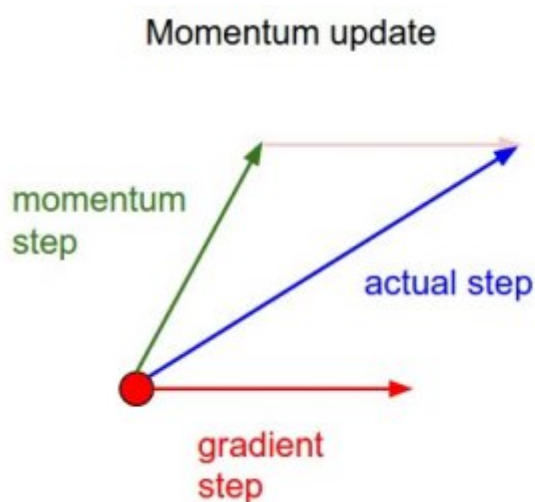
- $0 < \alpha < 1$, a hyperparameter

■ Apply update:

$$\theta \leftarrow \theta + v$$

Nesterov Momentum (Nesterov, 2004)

- With Nesterov momentum, the gradient is evaluated after current velocity is applied
- Comparison of momentum and Nesterov momentum:



From lecture material of CS231, Stanford Univ.

SGD with Nesterov Momentum

- Minibatch of the training set:

data $\{x^{(1)}, \dots, x^{(m)}\}$ with targets $\{y^{(1)}, \dots, y^{(m)}\}$

- Apply interim update:

$$\tilde{\theta} \leftarrow \theta + \alpha v$$

- Gradient:

$$\hat{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i^m L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$$

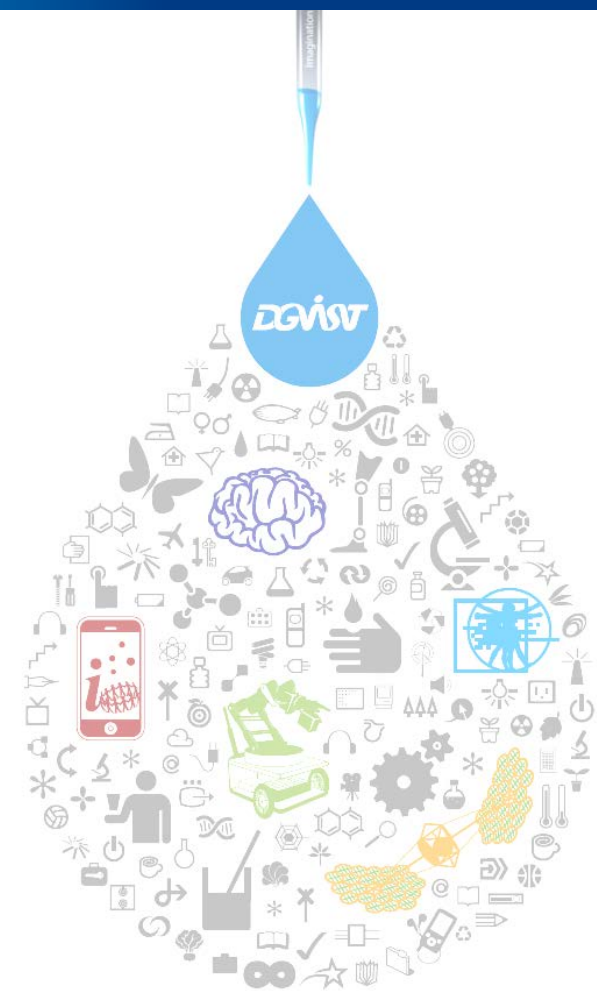
- Velocity:

$$v \leftarrow \alpha v - \epsilon \hat{g}$$

- Apply update:

$$\theta \leftarrow \theta + v$$

Algorithms with Adaptive Learning Rates



AdaGrad (Duchi et al., 2011)

- Adaptive gradient (AdaGrad) is a modified SGD with **per-parameter learning rate**

- increases the learning rate for parameters having small gradient
- decreases the learning rate for parameters having large gradient

- Gradient:

$$\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(x^{(i)}; \theta), y^{(i)})$$

- Accumulate squared gradients:

$$r \leftarrow r + g \odot g$$

\odot : element-wise multiplication

r : initialized to 0

δ : small constant value (e.g., 10^{-7})

- Apply update:

$$\theta \leftarrow \theta - \frac{\epsilon}{\delta + \sqrt{r}} \odot g$$

element-wise operation

RMSProp (Hinton, 2012)

- Root mean square propagation (RMSProp) is the modified version of AdaGrad

➤ by changing the gradient accumulation into an exponentially weighted moving average

- Gradient:

$$\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$$

- Accumulate squared gradients:

$$\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$$

$$0 < \rho < 1$$

- Apply update:

$$\theta \leftarrow \theta - \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$$

Adam (Kingma and Ba, 2014)

- Adaptive moments (Adam) is the combination of RMSProp and momentum

- Accumulate first and second moment estimates:

$$\begin{aligned}s &\leftarrow \rho_1 s + (1 - \rho_1)g \\ r &\leftarrow \rho_2 r + (1 - \rho_2)g \odot g\end{aligned}$$

- Correct bias in first and second moments:

$$\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}, \quad \hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$$

- Apply update:

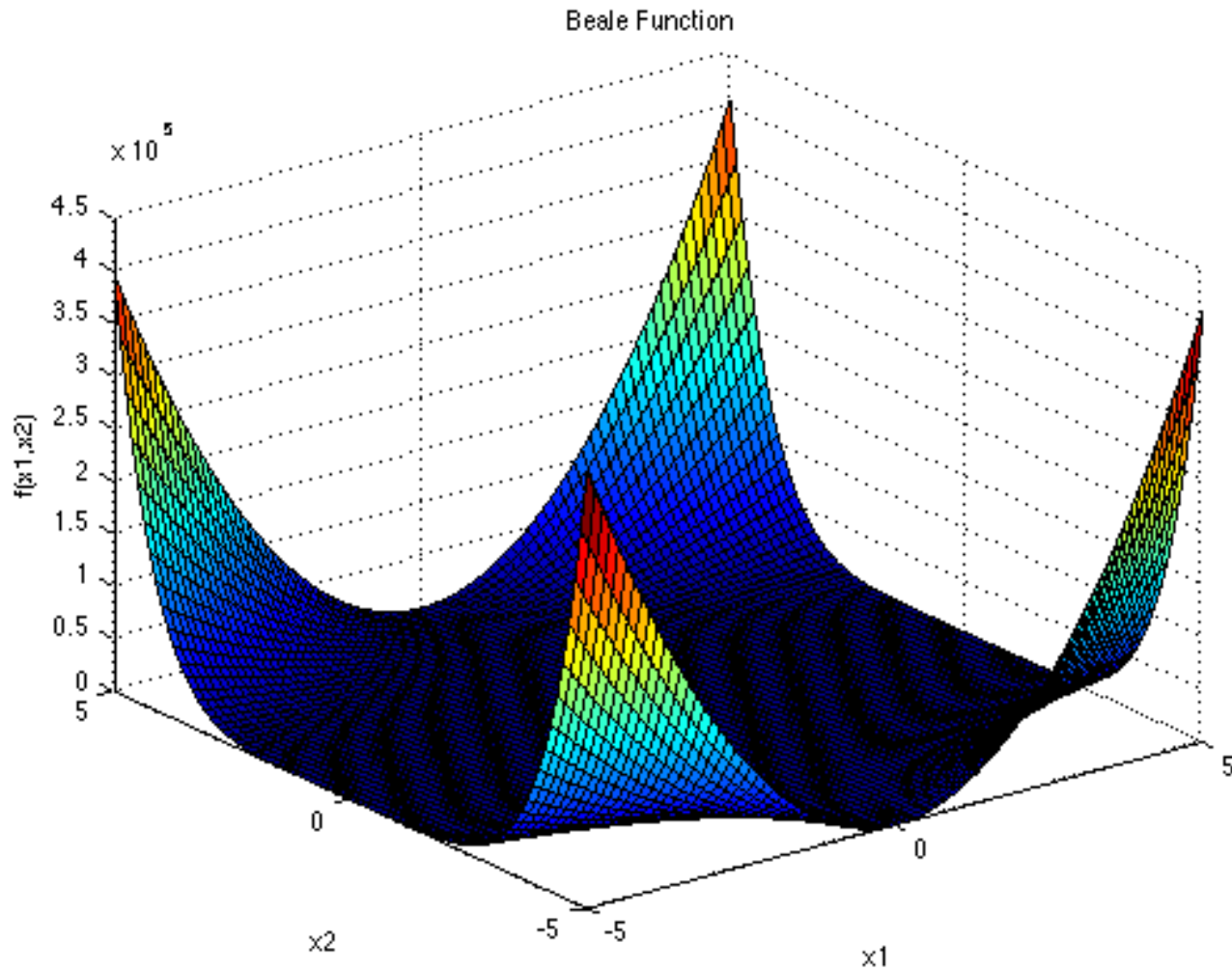
$$\theta \leftarrow \theta - \epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$$

ρ_1 : decay rate for 1st moment estimate
(usually 0.9)

ρ_2 : decay rate for 2nd moment estimate
(usually 0.999)

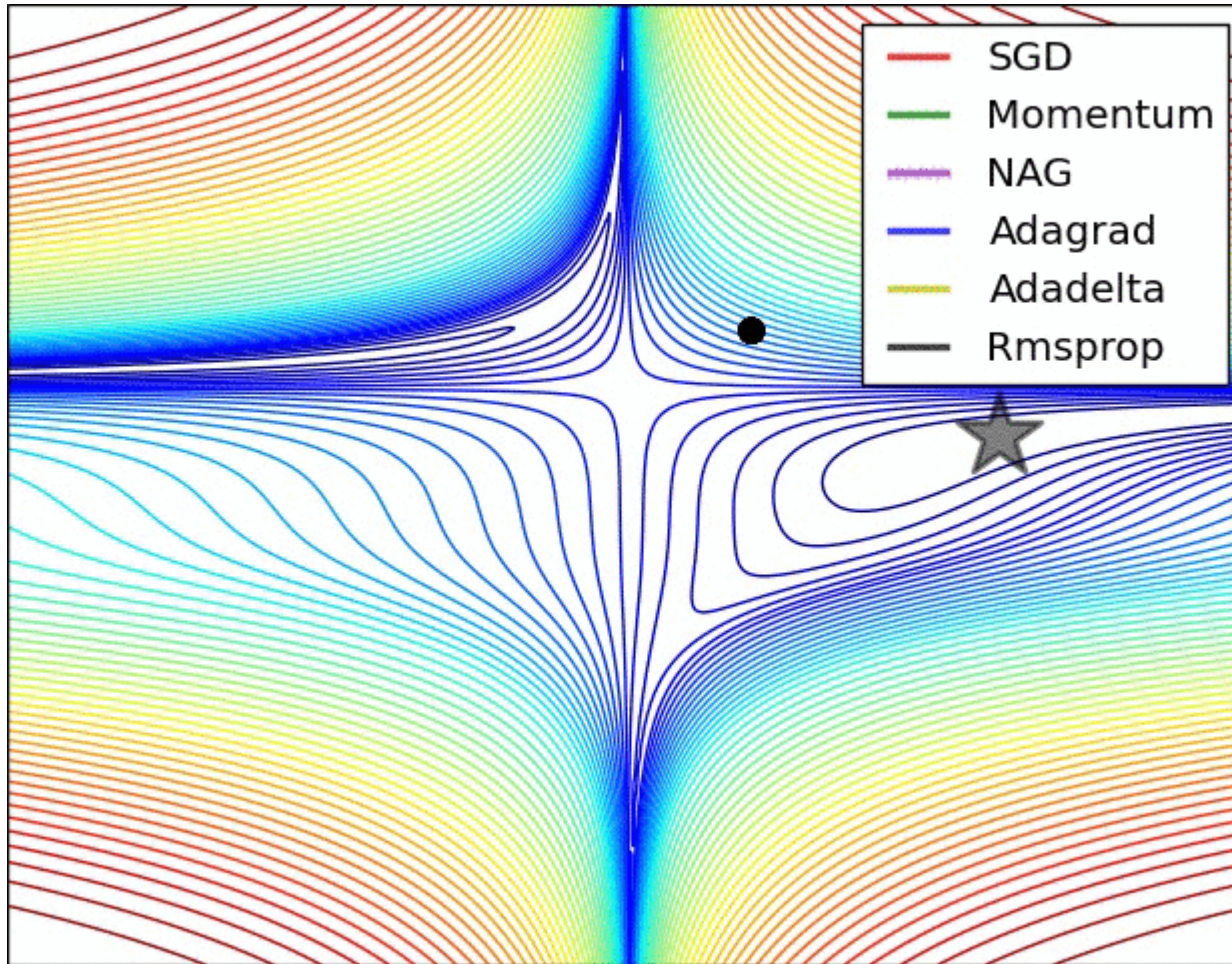
t : current number of iteration

Example: Beale's Function



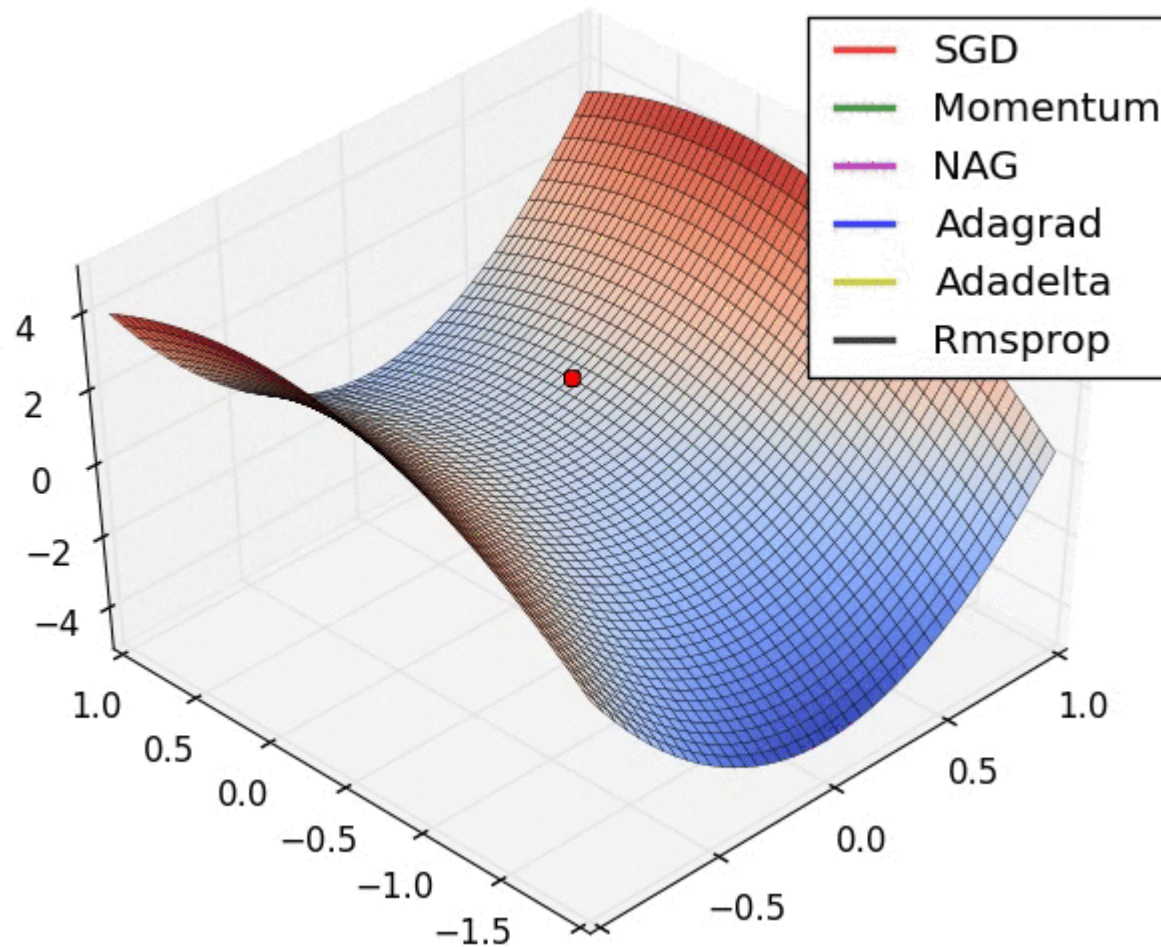
$$f(\mathbf{x}) = (1.5 - x_1 + x_1x_2)^2 + (2.25 - x_1 + x_1x_2^2)^2 + (2.625 - x_1 + x_1x_2^3)^2$$

Algorithms at Beale's Function

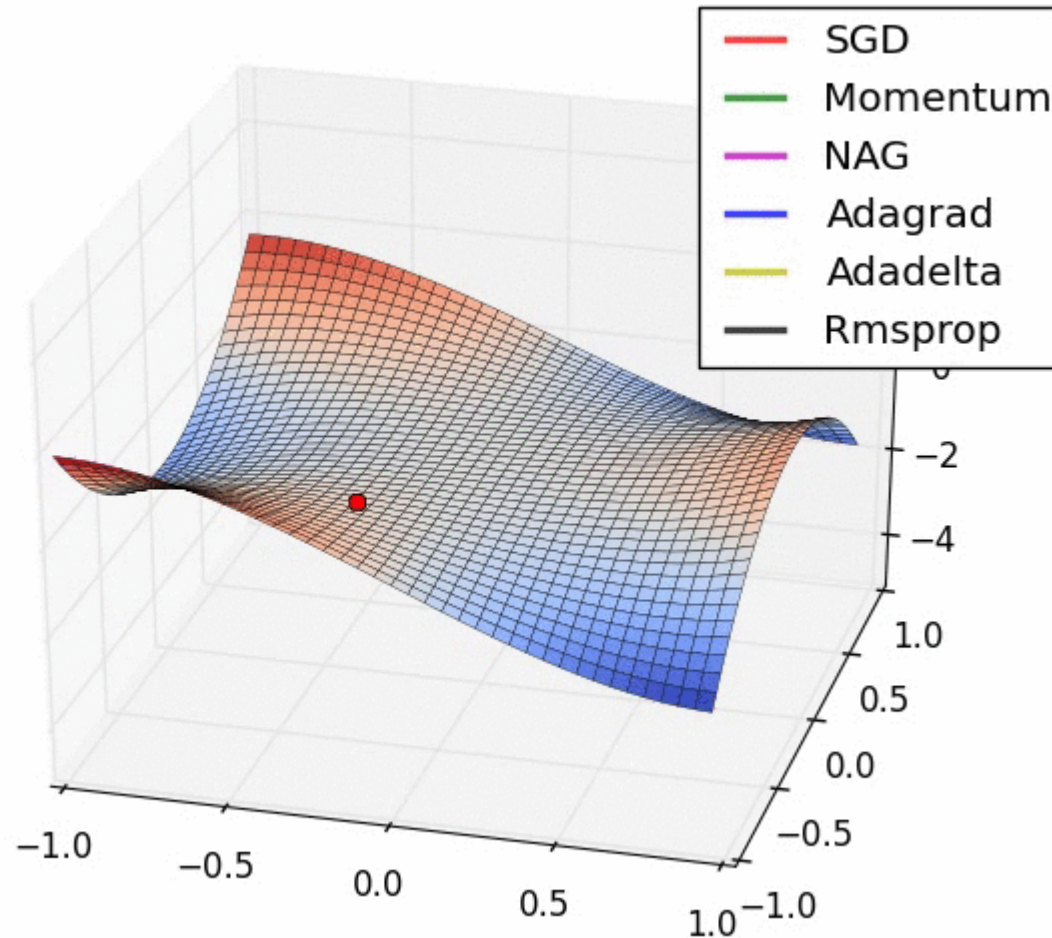


Lecture material in CS231n, "Convolutional Neural Networks for Visual Recognition", Stanford Univ. (Image credit: Alec Radford)

Algorithms at Long Valley



Algorithms at Saddle Point



Hyperparameter Setting in Caffe

```
net: "svhn/fix2-svhn-3conv-model.prototxt"

test_iter: 100
test_interval: 1000
max_iter: 100000

lr_policy: "fixed"
base_lr: 0.001
momentum: 0.9
weight_decay: 0.004
type: "SGD"

display: 10000
snapshot: 80000
snapshot_prefix: "svhn/snapshot/fix2_svhn_3conv"
solver_mode: GPU
```

- Stochastic Gradient Descent (type: "SGD"),
- AdaDelta (type: "AdaDelta"),
- Adaptive Gradient (type: "AdaGrad"),
- Adam (type: "Adam"),
- Nesterov's Accelerated Gradient (type: "Nesterov") and
- RMSProp (type: "RMSProp")

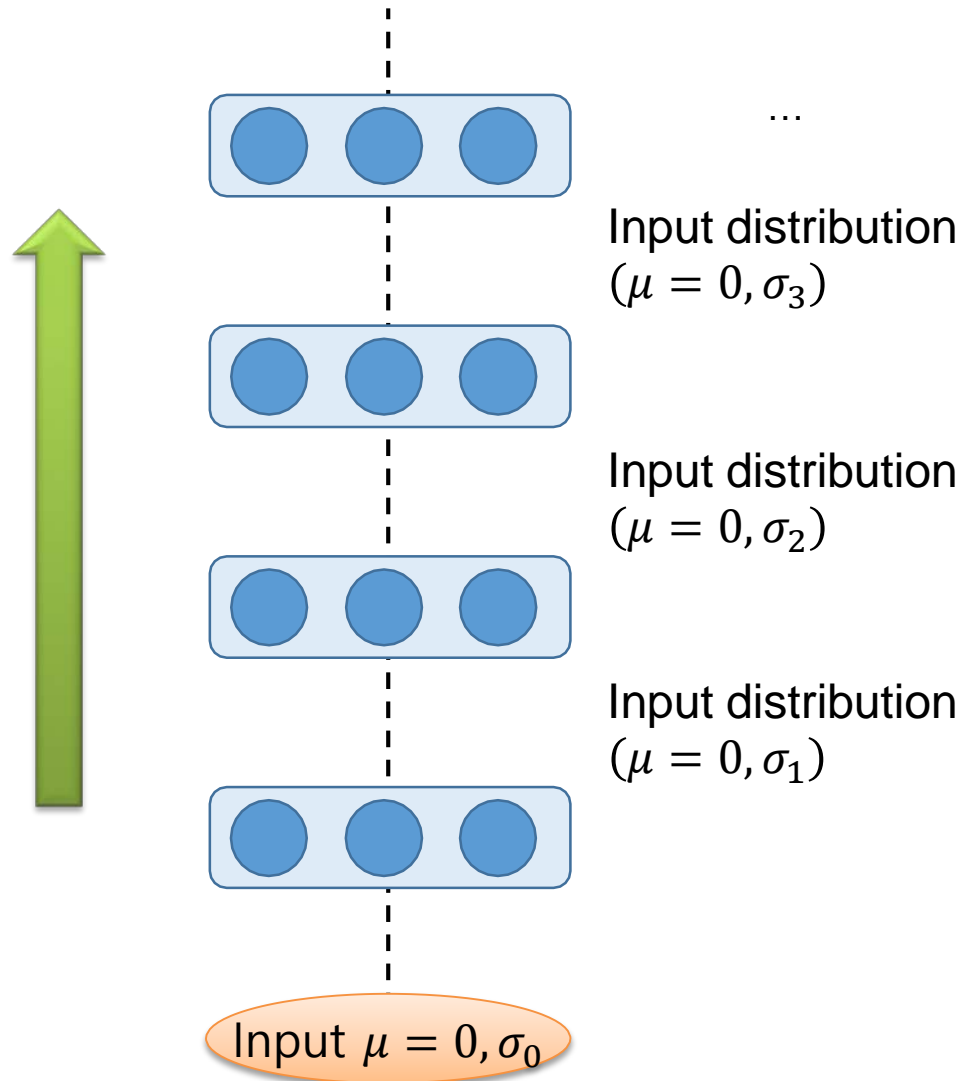
Initial Parameters in Deep Learning

- **Determining appropriate initial parameters for a deep model is one of the most important part of optimization**
- **Initial parameters can determine,**
 - whether the algorithm converges or not
 - whether the model converges to a point with a high or low cost
 - how quickly learning converges
- **However, the initialization strategies are still heuristic**
 - initial values are drawn randomly from a Gaussian or uniform distribution

Initialization with Appropriate Scale

- If the weights in a network start too small,
 - the signal shrinks as it passes through each layer
- If the weights in a network start too large,
 - the signal explodes as it passes through each layer
- **Xavier** initialization makes the signal keep in a **reasonable range** of values through many layers
- Xavier algorithm initializes W 's with drawing them from a **zero mean** and $\text{Var}(W_i) = \frac{1}{n_{\text{in}}}$

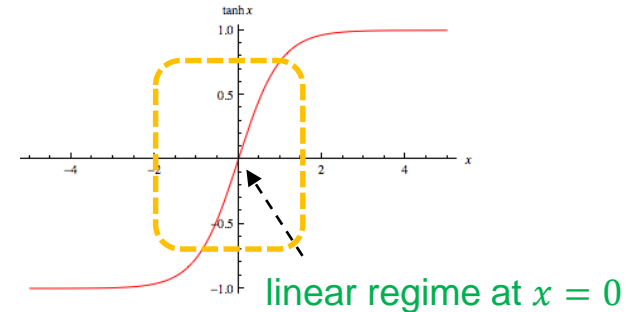
Basic Idea of Xavier Initialization



Linear Activation at Each Layer

■ $Y_i = W_1X_1 + \dots + W_nX_n$

- input X with n components
- linear neuron with random weight W



■ $\text{Var}(W_iX_i) =$
 $E[X_i]^2 \text{Var}(W_i) + E[W_i]^2 \text{Var}(X_i) + \text{Var}(W_i)\text{Var}(X_i)$

- assume, $E[X_i] = 0$ and $E[W_i] = 0$

$\therefore \text{Var}(W_iX_i) = \text{Var}(W_i)\text{Var}(X_i)$

■ $\text{Var}(Y) = \text{Var}(W_1X_1 + \dots + W_nX_n) = n\text{Var}(W_i)\text{Var}(X_i)$

- assume, X_i and W_i are i.i.d.

Xavier Initialization Algorithm

- $\text{Var}(Y) = n_{\text{in}} \text{Var}(W_i) \text{Var}(X_i)$

- the variance of the output Y is the variance of input X scaled by $n_{\text{in}} \text{Var}(W_i)$

- **Ideal case:** $\text{Var}(Y) = \text{Var}(X_i)$

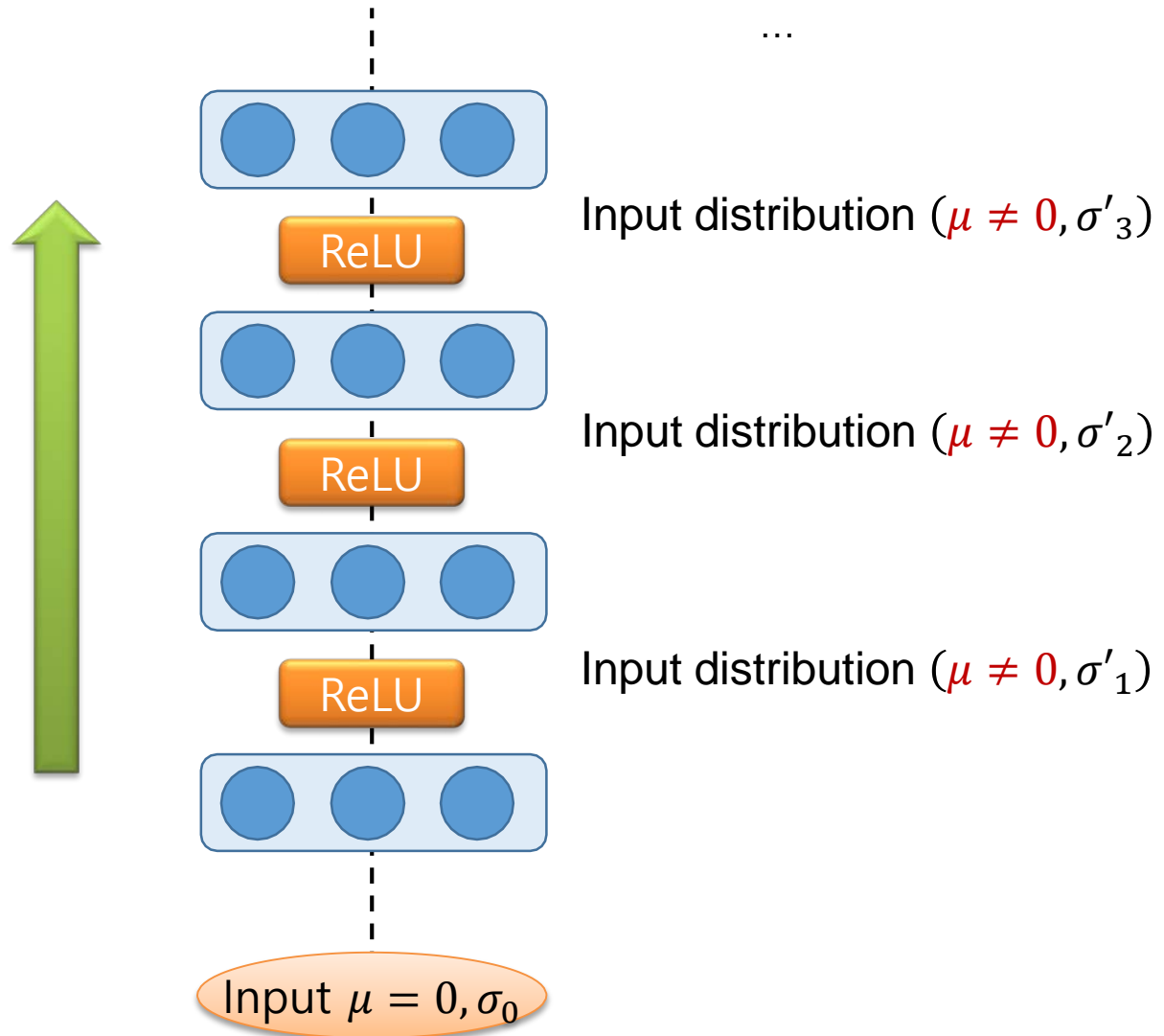
- i.e. $n_{\text{in}} \text{Var}(W_i) = 1$

- **So, for the weight initialization,**

$$\text{Var}(W_i) = \frac{1}{n_{\text{in}}}$$

$$\text{(or } \text{Var}(W_i) = \frac{2}{n_{\text{in}} + n_{\text{out}}}\text{)}$$

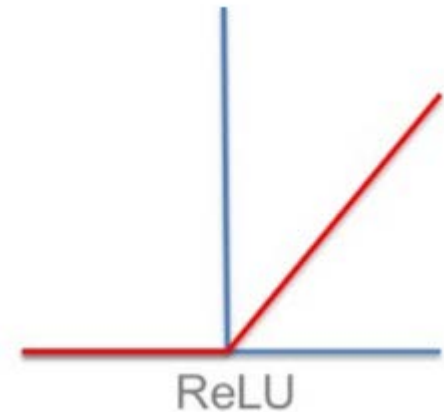
Basic Idea of MSRA Initialization



Initialization for Rectifiers (a.k.a. MSRA)

- $y_l = w_l x_l + b_l$

- n : # inputs for current layer
- l : an index of layer
- b : a bias term



- ReLU **does not have** zero mean

$$\because x_l = \max(0, y_{l-1})$$

- $\text{Var}(y_l) = n_l \text{Var}(w_l x_l) = n_l \text{Var}(w_l) E(x_l^2)$

- assume, $E[w] = 0$, but $E[x] \neq 0$
- x and w are i.i.d.

Distribution in l -th Layer

$$\begin{aligned} E[x_l^2] &= \int_{-\infty}^{\infty} \max(0, y_{l-1})^2 p(y) dy \\ &= \int_0^{\infty} y_{l-1}^2 p(y_{l-1}) dy \end{aligned}$$

- $\because y_{l-1}^2$ is symmetric around 0
- $p(y_{l-1})$ is assumed to symmetric around 0
($\because w_{l-1}$ is assumed to symmetric around 0)

$$\begin{aligned} &= \frac{1}{2} \int_{-\infty}^{\infty} y_{l-1}^2 p(y_{l-1}) dy \\ &= \frac{1}{2} \int_{-\infty}^{\infty} (y_{l-1} - E[y_{l-1}])^2 p(y_{l-1}) dy \\ &= \frac{1}{2} \text{Var}[y_{l-1}] \end{aligned}$$

Variance Propagation for All L Layers

- $E(x_l^2) = \frac{1}{2} \text{Var}(y_{l-1})$, when f is ReLU

- $\text{Var}(y_l) = \frac{1}{2} n_l \text{Var}(w_l) \text{Var}(y_{l-1})$

- With L layers put together,

$$\text{Var}(y_L) = \text{Var}(y_1) \left(\prod_{l=2}^L \frac{1}{2} n_l \text{Var}(w_l) \right)$$

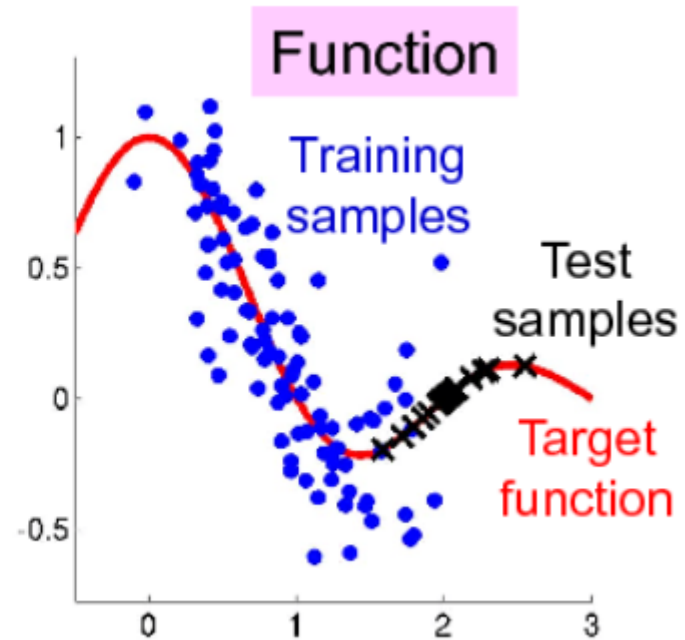
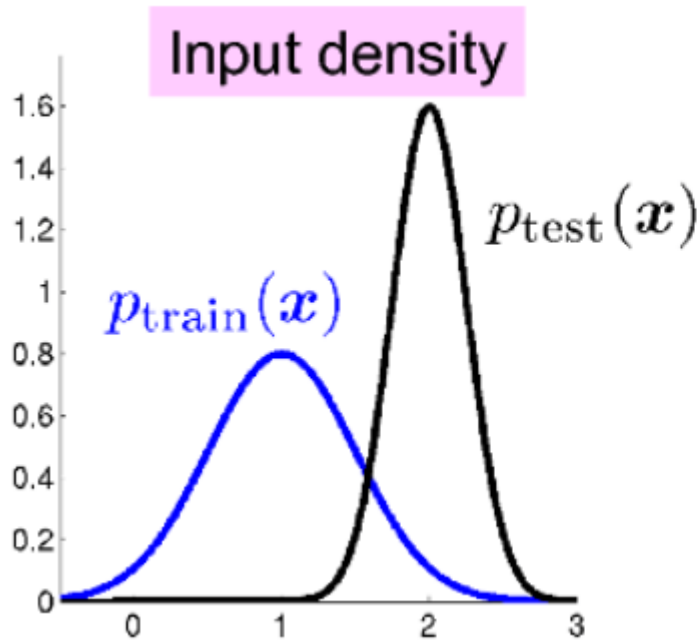
should avoid exponential growth

$$\rightarrow \frac{1}{2} n_l \text{Var}(w_l) = 1, \forall l$$

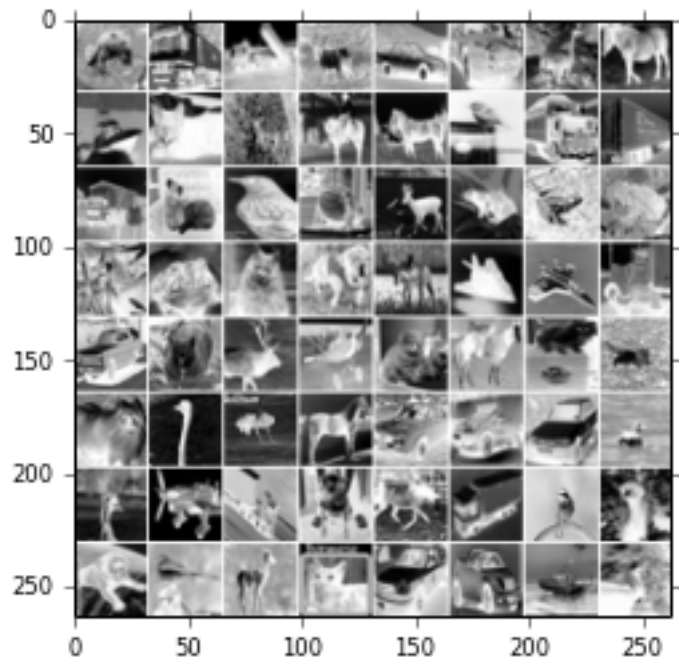
- Initialize W 's with drawing them from a zero-mean Gaussian distribution whose $\text{Var}(w_l) = \frac{2}{n_{in}}$

Covariate Shift (Shimodaira, 2000)

- Training / testing **input distributions** are different, but target function remains unchanged
 - causing weak extrapolation

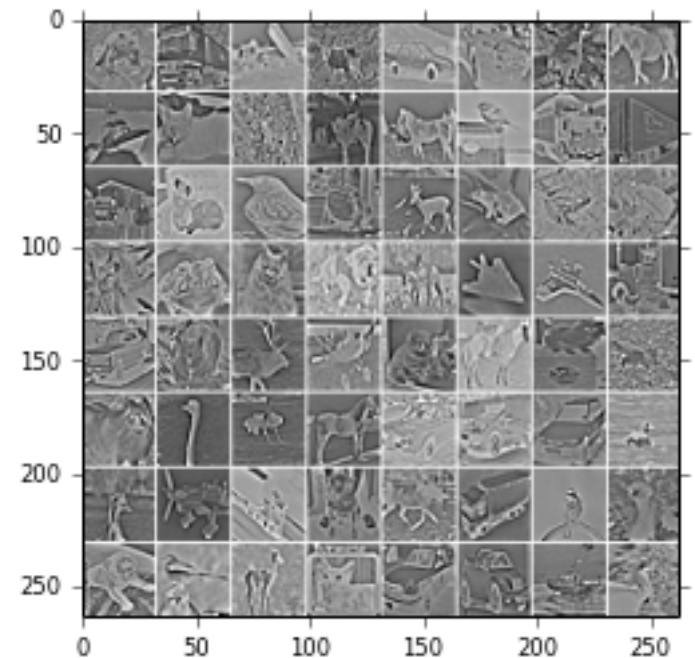


Solution for Covariate Shift



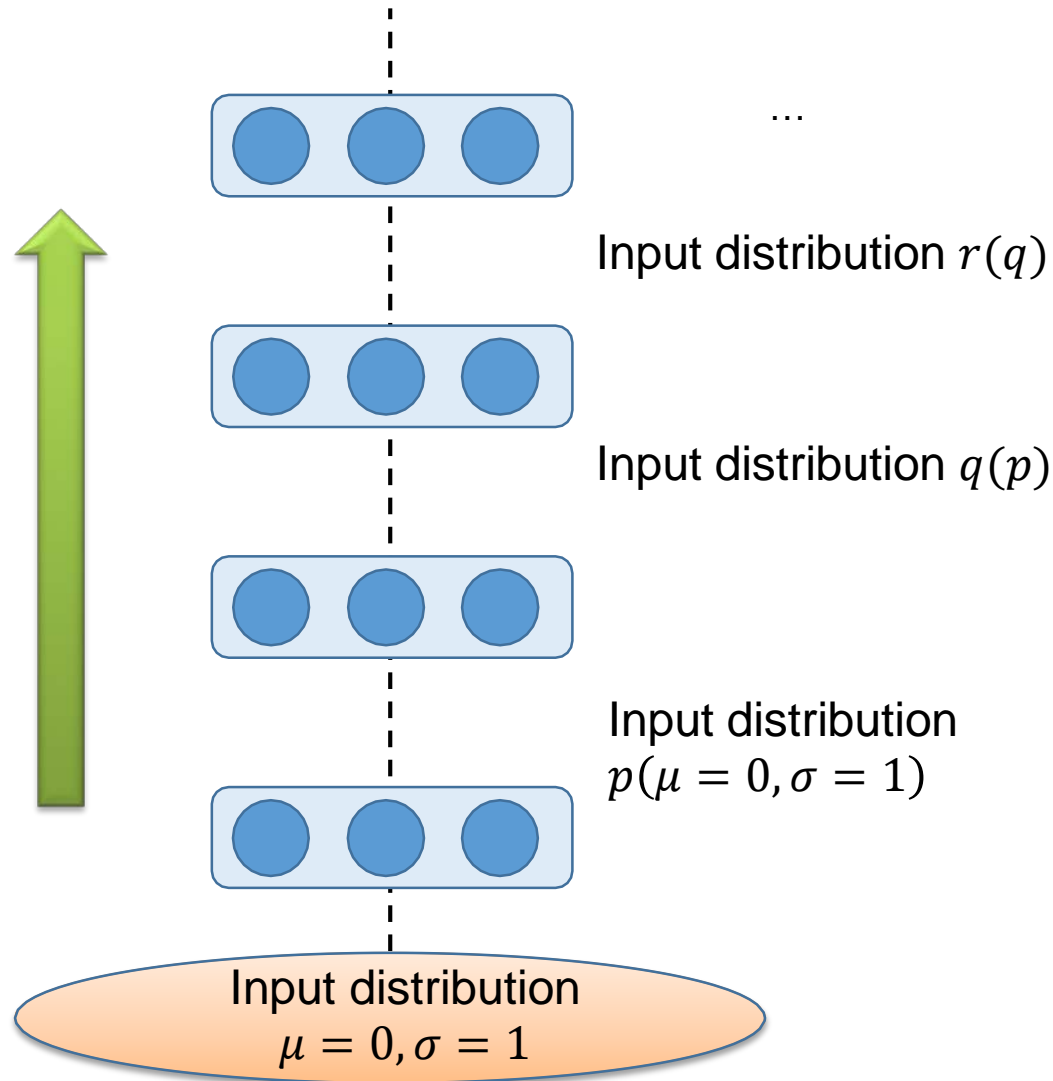
Transform:
mean = 0
variance = 1

whitening

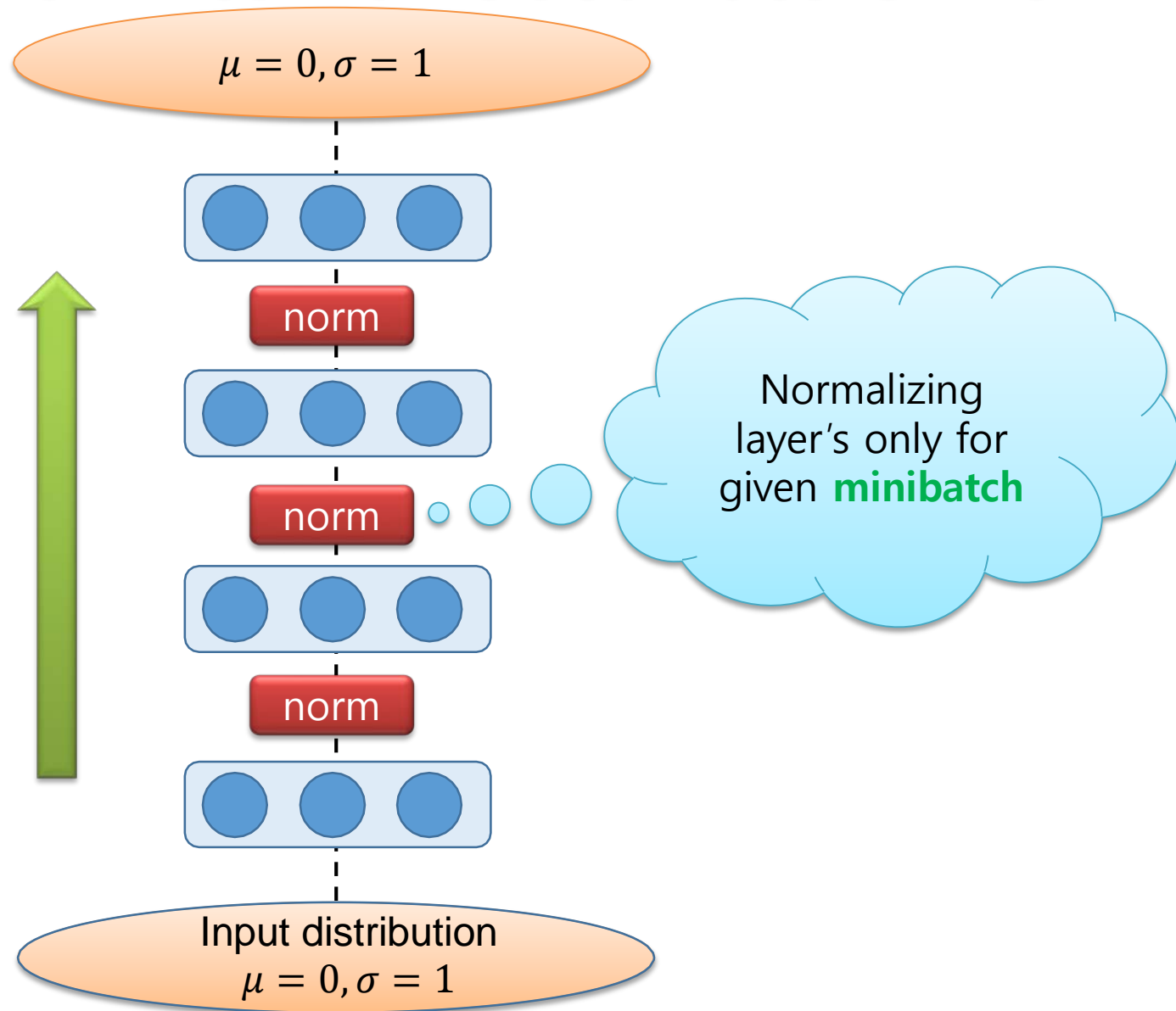


Answer by bayerj, "What is the difference between ZCA whitening and PCA whitening?", Oct 1, 2014

Internal Covariate Shift in DNN



Solution for Internal Covariate Shift



Batch Normalization (BN)

- Normalizing each layer's input of given **minibatch**

m : size of mini-batch
 $\mathcal{B} = \{x_1 \dots x_m\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad \leftarrow \text{minibatch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad \leftarrow \text{minibatch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad \leftarrow \text{normalization}$$

- Two necessary simplifications:

- normalize each feature independently, $N(0,1)$
- calculate mean and variance only for a **current training minibatch**

BN Transformation

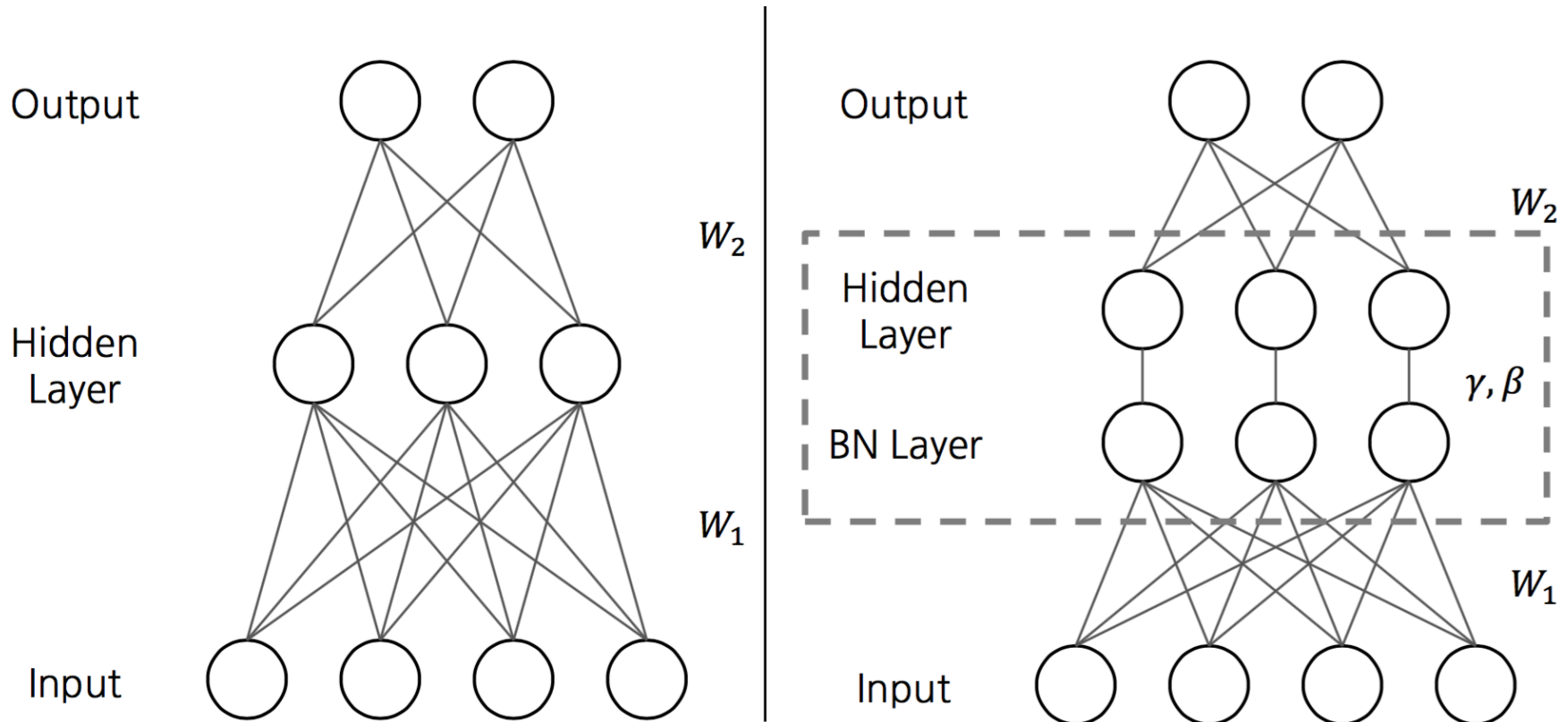
- Simply normalizing each input of a layer may change what the layer can represent
- Introduce a **linear transformation** for each activation

$$BN_{\gamma, \beta}(x_{1\dots m}) = \gamma \hat{x}_{1\dots m} + \beta$$

γ, β : parameters of
BN transformation
 \hat{x}_i : normalized input

- γ and β are learned (updated) at BP steps
 - γ : scale factor
 - β : shift factor

NN with & without BN



Comparison between neural networks with and without BN.

Next Chapter: Convolutional Networks

- The Convolution Operation
- Motivation
- Pooling
- Convolution and Pooling as an Infinitely Strong Prior
- Variants of the Basic Convolution Function
- Structured Outputs
- Data Types
- Efficient Convolution Algorithms
- Random or Unsupervised Features
- The Neuroscientific Basis for Convolutional Networks
- Convolutional Networks and the History of Deep Learning

Q&A

Thank you

