

# Batch Normalization

Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in Optimization for Deep Models

- Importance of Optimization in machine learning
- How learning differs from optimization
- Challenges in neural network optimization
- Basic Optimization Algorithms
- Parameter initialization strategies
- Algorithms with adaptive learning rates
- Approximate second-order methods
- Optimization strategies and meta-algorithms

# Topics in Optimization Strategies and Meta-Algorithms

1. Batch Normalization
2. Coordinate Descent
3. Polyak Averaging
4. Supervised Pretraining
5. Designing Models to Aid Optimization
6. Continuation Methods and Curriculum Learning

# Overview of Optimization Strategies

- Many optimization techniques are general templates that can be specialized to yield algorithms
- They can be incorporated into different algorithms

# Topics in Batch Normalization

- Batch normalization: exciting recent innovation
- Motivation is difficulty of choosing learning rate  $\epsilon$  in deep networks
- Method is to replace activations with zero-mean with unit variance activations

# Adding normalization between layers

- Motivated by difficulty of training deep models
- Method adds an additional step between layers, in which the output of the earlier layer is normalized
  - By standardizing the mean and standard deviation of each individual unit
- It is a method of adaptive re-parameterization
  - It is not an optimization algorithm at all
  - A method to reduce internal covariate shift in neural networks

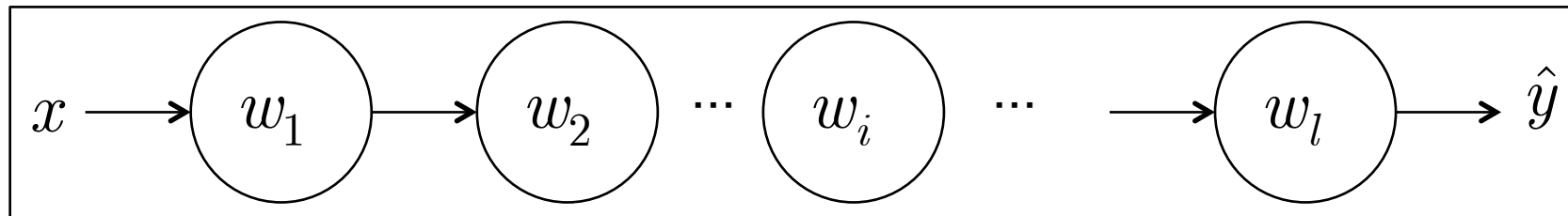
# Motivation: Difficulty of composition

- Very deep models involve compositions of several functions or layers
  - **Ex:**  $f(\mathbf{x}, \mathbf{w}) = f^{(l)}[\dots f^{(3)}[f^{(2)}[f^{(1)}(\mathbf{x})]]]$  has depth of  $l$
- The gradient tells us how to update each parameter
 
$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \varepsilon \nabla_{\mathbf{w}} E_{\mathbf{x}, y \in \hat{p}_{data}} L(f(\mathbf{x}; \mathbf{w}), y)$$
  - Under assumption that other layers do not change
  - We update all  $l$  layers simultaneously
  - When we make the change unexpected results can happen
    - Because many functions are changed simultaneously
    - Under the assumption that other functions remain constant

# Choosing learning rate $\varepsilon$ in multilayer

- Simple example:

- $l$  layers, one unit per layer, no activation function



- Network computes

$$\hat{y} = x \cdot w_1 \cdot w_2 \cdot w_3 \dots \cdot w_i \dots \cdot w_l$$

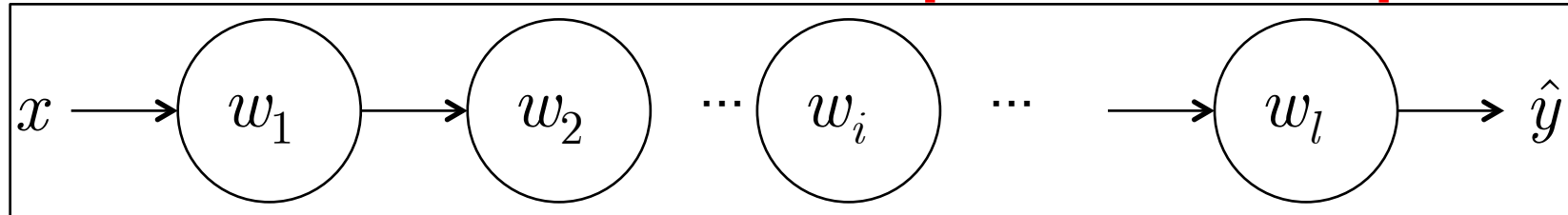
- where  $w_i$  provides weight of layer  $i$

- Output of layer  $i$  is  $h_i = h_{i-1} w_i$

- Output is a linear function of input  $x$  but a nonlinear function of the weights  $w_i$



# Gradient in Simple example



- Suppose cost has put a gradient of 1 on  $\hat{y}$ , so we wish to decrease  $\hat{y}$  slightly
  - Note: value of 1 is quite a large gradient  $\mathbf{g} = \nabla_{\mathbf{w}} E_{\mathbf{x}, y \in \hat{p}_{data}} L(f(\mathbf{x}; \mathbf{w}), y)$
- Backpropagation can compute a gradient  $\mathbf{g} = \nabla_{\mathbf{w}} \hat{y}$ 
  - which is gradient  $\mathbf{g}$  evaluated at  $y = \hat{y}$
- Consider update  $\mathbf{w} \leftarrow \mathbf{w} - \varepsilon \mathbf{g}$ 
  - First order Taylor's series approx of  $\hat{y}$  predicts that the value of  $\hat{y}$  will decrease by  $\varepsilon \mathbf{g}^T \mathbf{g}$ , because:

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^T \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^T H (\mathbf{x} - \mathbf{x}^{(0)})$$

Substituting  
 $\mathbf{x} = \mathbf{x}^{(0)} - \varepsilon \mathbf{g}$

$$f(\mathbf{x}^{(0)} - \varepsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \varepsilon \mathbf{g}^T \mathbf{g} + \frac{1}{2} \varepsilon^2 \mathbf{g}^T H \mathbf{g}$$

# Difficulty of Multilayer learning Rate

- To decrease  $\hat{y}$  by 0.1 we could set  $\varepsilon$  to  $\frac{0.1}{g^T g}$ 

$$f(\mathbf{x}^{(0)} - \varepsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \varepsilon \mathbf{g}^T \mathbf{g} + \frac{1}{2} \varepsilon^2 \mathbf{g}^T \mathbf{H} \mathbf{g}$$
- However, update includes 2<sup>nd</sup>, 3<sup>rd</sup> .. order effects
  - Since  $w \leftarrow w - \varepsilon g$  new value of  $\hat{y}$  is  $x(w_1 - \varepsilon g_1)(w_2 - \varepsilon g_2) \dots (w_l - \varepsilon g_l)$ 
    - A second order term arising from this update is  $\varepsilon^2 g_1 g_2 \prod_{i=3}^l w_i$
    - This term can be negligible if  $\prod_{i=3}^l w_i$  is small or exponentially large if weights on layers 3 to  $l$  are greater than 1
- This makes it very hard to choose  $\varepsilon$  because the effects of an update for one layer depend so strongly on all other layers
  - Second-order optimization algorithms address this issue but it seems hopeless for  $l > 2$

# The Batch Normalization Solution

- Provides an elegant way of reparameterizing almost any network
- Significantly reduces the problem of coordinating updates across many layers
- Can be applied to any input or hidden layer in a network

# Batch Normalization Equations

- **$H$** : minibatch of activations of layer to normalize

– arranged as a design matrix

- With activations for each example appearing in a row

$$H = \begin{array}{c} \text{units} \longrightarrow \\ \left[ \begin{array}{ccc} a_{1,1} & a_{1,2} & a_{1,3} \\ \cdot & \cdot & \cdot \\ a_{N,1} & a_{N,2} & a_{N,3} \end{array} \right] \text{samples} \end{array}$$

- To normalize  **$H$**  we replace it with  $H' = \frac{H - \mu}{\sigma}$

– Where  $\mu$  is a vector containing the mean of each unit and  $\sigma$  is the std deviation of each unit

- The arithmetic here is based on broadcasting the vector  $\mu$  and the vector  $\sigma$  to be applied to every row of  $H$

– Within each row, the arithmetic is element-wise

- $H_{i,j}$  is normalized by subtracting  $\mu_{i,j}$  & dividing by  $\sigma_j$

# Normalization Details

- Rest of the network operates on  $H'$  in the same way that the original network operated on  $H$
- At training time  $\mu = \frac{1}{m} \sum_i H_i$  and  $\sigma = \sqrt{\delta + \frac{1}{m} \sum_i (H_i - \mu)^2}$ 
  - where  $\delta$  is a small positive value such as  $10^{-8}$  imposed to avoid encountering the undefined gradient of  $\sqrt{z}$  at  $z=0$
  - Crucially we back propagate through these operations for computing the mean and std dev
  - And for applying them to normalize  $H$ 
    - This means that the gradient will never propose an operation that acts simply to increase std dev or mean of  $h_i$  the normalization operations remove the effect of such an action and zero out the component in the gradient

# Batch Normalization at Test time

- At test time,  $\mu$  and  $\sigma$  may be replaced by running averages that were collected during training time
- This allows the model to be evaluated on a single example without needing to use definitions of  $\mu$  and  $\sigma$  that depend on an entire minibatch

# Revisiting the simple example

- Revisiting the  $\hat{y} = x \cdot w_1 \cdot w_2 \cdot w_3 \dots \cdot w_i \dots \cdot w_l$  example we can mostly resolve the difficulties in learning the model by normalizing  $h_{l-1}$
- Suppose that  $x$  is drawn from a unit Gaussian
- Then  $h_{l-1}$  will also come from a Gaussian, because the transformation from  $x$  to  $h_l$  is linear
- However  $h_{l-1}$  will no longer have zero mean, unit variance

# Restoring zero-mean unit variance

- After applying batch normalization, we obtain the normalized  $\hat{h}_{l-1}$  that restores zero mean and unit variance
  - For almost any update to the lower layers,  $\hat{h}_{l-1}$  will remain a unit Gaussian
  - Output  $\hat{y}$  may be learned as a simple linear function  $\hat{y} = w_l \hat{h}_l - 1$
- Learning in this model is now very simple
  - Because parameters at the lower layers do not have an effect in most cases
    - Their output is always renormalized to a unit Gaussian



# Batch Normalization → learning easy

- Without normalization, updates would have an extreme effect of the statistics of  $h_{l-1}$
- Batch normalization has thus made this model easier to learn
- In this example the ease of learning came from making the lower layers useless
  - Lower layers not harmful but not beneficial either
    - Because we have normalized-out 1<sup>st</sup> and 2<sup>nd</sup> order stats
- In a deep neural network lower levels can perform useful nonlinear transformations

# Reintroducing Expressive Power

- Normalizing the mean and standard deviation can reduce the expressive power of the neural network containing that unit
- To maintain the expressive power replace the batch of hidden unit activations  $H'$  with  $\gamma H' + \beta$ 
  - $\gamma$  and  $\beta$  are learned parameters that allow the new variable to have any mean and standard deviation
- Has different learning dynamics than unnormalized approach