

Deep Learning: Lecture 7

Alexander Schönhuth

UU

November 20, 2019

Backpropagation for CNNs

CONVOLUTIONAL BACKPROPAGATION

NOTATION

- In fully connected layers we had [w_{jk}^l is for connecting the k -th neuron in layer $l - 1$ with the j -th neuron in layer l]

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} \quad \text{where} \quad z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l \quad (1)$$

and $a_j^l = \sigma(z_j^l)$ (where σ is any activation function)

- For sake of simplicity, we assume that there is only one hidden sublayer (corresponding to one feature map, or one channel) at each level of depth, so only one convolutional filter per level of depth required.

CONVOLUTIONAL BACKPROPAGATION

NOTATION

- In fully connected layers we had [w_{jk}^l is for connecting the k -th neuron in layer $l - 1$ with the j -th neuron in layer l]

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} \quad \text{where} \quad z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l \quad (1)$$

and $a_j^l = \sigma(z_j^l)$ (where σ is any activation function)

- For sake of simplicity, we assume that there is only one hidden sublayer (corresponding to one feature map, or one channel) at each level of depth, so only one convolutional filter per level of depth required.

CONVOLUTIONAL BACKPROPAGATION

NOTATION

- ▶ We index neurons using two coordinates, so $z_{x,y}^l$ is the input for the x, y -th neuron of the hidden layer at level of depth l .
- ▶ The $M \times M$ filter that connects neurons from level l with neurons at level $l + 1$ has weights w_{ab}^{l+1} , $1 \leq a, b \leq M$
- ▶ By applying the convolution operation [and neglecting the exact indexing in the following]

$$z_{x,y}^{l+1} = \sum_a \sum_b w_{ab}^{l+1} \sigma(z_{x-a,y-b}^l) + b_{x,y}^{l+1} \quad (2)$$

CONVOLUTIONAL BACKPROPAGATION

NOTATION

- ▶ We index neurons using two coordinates, so $z_{x,y}^l$ is the input for the x, y -th neuron of the hidden layer at level of depth l .
- ▶ The $M \times M$ filter that connects neurons from level l with neurons at level $l + 1$ has weights w_{ab}^{l+1} , $1 \leq a, b \leq M$
- ▶ By applying the convolution operation [and neglecting the exact indexing in the following]

$$z_{x,y}^{l+1} = \sum_a \sum_b w_{ab}^{l+1} \sigma(z_{x-a,y-b}^l) + b_{x,y}^{l+1} \quad (2)$$

CONVOLUTIONAL BACKPROPAGATION

- We compute

$$\delta_{x,y}^l = \frac{\partial C}{\partial z_{x,y}^l} = \sum_{x'} \sum_{y'} \frac{\partial C}{\partial z_{x',y'}^{l+1}} \frac{\partial z_{x',y'}^{l+1}}{\partial z_{x,y}^l} \quad (3)$$

- Moving on, we get

$$\begin{aligned} \frac{\partial C}{\partial z_{x,y}^l} &= \sum_{x'} \sum_{y'} \frac{\partial C}{\partial z_{x',y'}^{l+1}} \frac{\partial z_{x',y'}^{l+1}}{\partial z_{x,y}^l} \\ &= \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} \frac{\partial (\sum_a \sum_b w_{ab}^{l+1} \sigma(z_{x'-a,y'-b}^l) + b_{x',y'}^{l+1})}{\partial z_{x,y}^l} \end{aligned} \quad (4)$$

CONVOLUTIONAL BACKPROPAGATION

- We compute

$$\delta_{x,y}^l = \frac{\partial C}{\partial z_{x,y}^l} = \sum_{x'} \sum_{y'} \frac{\partial C}{\partial z_{x',y'}^{l+1}} \frac{\partial z_{x',y'}^{l+1}}{\partial z_{x,y}^l} \quad (3)$$

- Moving on, we get

$$\begin{aligned} \frac{\partial C}{\partial z_{x,y}^l} &= \sum_{x'} \sum_{y'} \frac{\partial C}{\partial z_{x',y'}^{l+1}} \frac{\partial z_{x',y'}^{l+1}}{\partial z_{x,y}^l} \\ &= \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} \frac{\partial (\sum_a \sum_b w_{ab}^{l+1} \sigma(z_{x'-a,y'-b}^l) + b_{x',y'}^{l+1})}{\partial z_{x,y}^l} \end{aligned} \quad (4)$$

CONVOLUTIONAL BACKPROPAGATION

- All terms in (4) where $x \neq x' - a$ or $y \neq y' - b$ are zero, so

$$\begin{aligned} \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} \frac{\partial(\sum_a \sum_b w_{ab}^{l+1} \sigma(z_{x'-a,y'-b}^l) + b_{x',y'}^{l+1})}{\partial z_{x,y}^l} \\ = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{ab}^{l+1} \sigma'(z_{x,y}^l) \quad (5) \end{aligned}$$

- Since now $x = x' - a, y = y' - b$, we have $a = x' - x, b = y' - y$, so

$$\sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{ab}^{l+1} \sigma'(z_{x,y}^l) = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{x'-x,y'-y}^{l+1} \sigma'(z_{x,y}^l) \quad (6)$$

CONVOLUTIONAL BACKPROPAGATION

- All terms in (4) where $x \neq x' - a$ or $y \neq y' - b$ are zero, so

$$\begin{aligned} \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} \frac{\partial(\sum_a \sum_b w_{ab}^{l+1} \sigma(z_{x'-a,y'-b}^l) + b_{x',y'}^{l+1})}{\partial z_{x,y}^l} \\ = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{ab}^{l+1} \sigma'(z_{x,y}^l) \quad (5) \end{aligned}$$

- Since now $x = x' - a, y = y' - b$, we have $a = x' - x, b = y' - y$, so

$$\sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{ab}^{l+1} \sigma'(z_{x,y}^l) = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{x'-x,y'-y}^{l+1} \sigma'(z_{x,y}^l) \quad (6)$$

CONVOLUTIONAL BACKPROPAGATION

- Summary:

$$\delta_{x,y}^l = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{x'-x,y'-y}^{l+1} \sigma'(z_{x,y}^l) \quad (7)$$

- A closer look reveals this as a convolution operation in its own right, applying the filter

$$\sigma'(z_{x,y}^l) \cdot \begin{pmatrix} w_{MM} & \cdots & w_{M1} \\ \vdots & \ddots & \vdots \\ w_{1M} & \cdots & w_{11} \end{pmatrix} \quad (8)$$

to the $l+1$ -th layer of gradients $\delta_{x',y'}^{l+1}$, for computing values $\delta_{x,y}^l$ of the l -th layer of gradients.

- Note that the weights of the original filter have been rotated by 180°
- For further illustrations, see <https://medium.com/@2017csm1006/forward-and-backpropagation-in-convolutional-neural-network-4dfa9>

Note that notation differs (error E there is cost C here, X there is z here, and F are the weights w here, and O is a here, and there is no bias)

CONVOLUTIONAL BACKPROPAGATION

- Summary:

$$\delta_{x,y}^l = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{x'-x,y'-y}^{l+1} \sigma'(z_{x,y}^l) \quad (7)$$

- A closer look reveals this as a convolution operation in its own right, applying the filter

$$\sigma'(z_{x,y}^l) \cdot \begin{pmatrix} w_{MM} & \cdots & w_{M1} \\ \vdots & \ddots & \vdots \\ w_{1M} & \cdots & w_{11} \end{pmatrix} \quad (8)$$

to the $l + 1$ -th layer of gradients $\delta_{x',y'}^{l+1}$, for computing values $\delta_{x,y}^l$ of the l -th layer of gradients.

- Note that the weights of the original filter have been rotated by 180°
- For further illustrations, see <https://medium.com/@2017csm1006/forward-and-backpropagation-in-convolutional-neural-network-4dfa9>
Note that notation differs (error E there is cost C here, X there is z here, and F are the weights w here, and O is a here, and there is no bias)

CONVOLUTIONAL BACKPROPAGATION

- Summary:

$$\delta_{x,y}^l = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{x'-x,y'-y}^{l+1} \sigma'(z_{x,y}^l) \quad (7)$$

- A closer look reveals this as a convolution operation in its own right, applying the filter

$$\sigma'(z_{x,y}^l) \cdot \begin{pmatrix} w_{MM} & \cdots & w_{M1} \\ \vdots & \ddots & \vdots \\ w_{1M} & \cdots & w_{11} \end{pmatrix} \quad (8)$$

to the $l + 1$ -th layer of gradients $\delta_{x',y'}^{l+1}$, for computing values $\delta_{x,y}^l$ of the l -th layer of gradients.

- Note that the weights of the original filter have been rotated by 180°
- For further illustrations, see <https://medium.com/@2017csm1006/forward-and-backpropagation-in-convolutional-neural-network-4dfa9>
Note that notation differs (error E there is cost C here, X there is z here, and F are the weights w here, and O is a here, and there is no bias)

CONVOLUTIONAL BACKPROPAGATION

- Summary:

$$\delta_{x,y}^l = \sum_{x'} \sum_{y'} \delta_{x',y'}^{l+1} w_{x'-x,y'-y}^{l+1} \sigma'(z_{x,y}^l) \quad (7)$$

- A closer look reveals this as a convolution operation in its own right, applying the filter

$$\sigma'(z_{x,y}^l) \cdot \begin{pmatrix} w_{MM} & \cdots & w_{M1} \\ \vdots & \ddots & \vdots \\ w_{1M} & \cdots & w_{11} \end{pmatrix} \quad (8)$$

to the $l + 1$ -th layer of gradients $\delta_{x',y'}^{l+1}$, for computing values $\delta_{x,y}^l$ of the l -th layer of gradients.

- Note that the weights of the original filter have been rotated by 180°
- For further illustrations, see <https://medium.com/@2017csm1006/forward-and-backpropagation-in-convolutional-neural-network-4dfa9>

Note that notation differs (error E there is cost C here, X there is z here, and F are the weights w here, and O is a here, and there is no bias)

Training Variations

HESSIAN TECHNIQUE

- ▶ As usual, let $\mathbf{w} = (w_1, w_2, \dots)$ be NN parameters (weights in particular), and C a cost function.
- ▶ By *Taylor's theorem*, we can write

$$\begin{aligned} C(w + \delta w) = C(w) + \sum_j \frac{\partial C}{\partial w_j} \delta w_j \\ + \frac{1}{2} \sum_{jk} \delta w_j H_{jk} \delta w_k + \dots \text{ terms of higher order} \end{aligned} \quad (9)$$

where H is the *Hessian matrix*, given by

$$H_{jk} = \frac{\partial^2 C}{\partial w_j \partial w_k} \quad (10)$$

HESSIAN TECHNIQUE

- ▶ As usual, let $\mathbf{w} = (w_1, w_2, \dots)$ be NN parameters (weights in particular), and C a cost function.
- ▶ By *Taylor's theorem*, we can write

$$\begin{aligned} C(w + \delta w) = C(w) + \sum_j \frac{\partial C}{\partial w_j} \delta w_j \\ + \frac{1}{2} \sum_{jk} \delta w_j H_{jk} \delta w_k + \dots \text{ terms of higher order} \end{aligned} \quad (9)$$

where H is the *Hessian matrix*, given by

$$H_{jk} = \frac{\partial^2 C}{\partial w_j \partial w_k} \quad (10)$$

HESSIAN TECHNIQUE

- ▶ Then (12) can further be written as

$$C(w + \delta w) = C(w) + \nabla C \delta w + \frac{1}{2} \delta w^T H \delta w + \dots \quad (11)$$

- ▶ Discarding all terms of order greater than 2, we obtain

$$C(w + \delta w) \approx C(w) + \nabla C \delta w + \frac{1}{2} \delta w^T H \delta w \quad (12)$$

- ▶ The right hand side of (12) can be minimized by choosing

$$\delta w = -H^{-1} \nabla C \quad (13)$$

HESSIAN TECHNIQUE

- Then (12) can further be written as

$$C(w + \delta w) = C(w) + \nabla C \delta w + \frac{1}{2} \delta w^T H \delta w + \dots \quad (11)$$

- Discarding all terms of order greater than 2, we obtain

$$C(w + \delta w) \approx C(w) + \nabla C \delta w + \frac{1}{2} \delta w^T H \delta w \quad (12)$$

- The right hand side of (12) can be minimized by choosing

$$\delta w = -H^{-1} \nabla C \quad (13)$$

HESSIAN TECHNIQUE

- ▶ Then (12) can further be written as

$$C(w + \delta w) = C(w) + \nabla C \delta w + \frac{1}{2} \delta w^T H \delta w + \dots \quad (11)$$

- ▶ Discarding all terms of order greater than 2, we obtain

$$C(w + \delta w) \approx C(w) + \nabla C \delta w + \frac{1}{2} \delta w^T H \delta w \quad (12)$$

- ▶ The right hand side of (12) can be minimized by choosing

$$\delta w = -H^{-1} \nabla C \quad (13)$$

HESSIAN TECHNIQUE

This suggests the following algorithm for updating weights in a gradient descent like scheme:

1. Choose a starting point \mathbf{w}
2. Update \mathbf{w} to $\mathbf{w}' = \mathbf{w} - \eta H^{-1} \nabla C$, where H and ∇C are computed at \mathbf{w}
3. Iterate ... (until appropriate criteria, like convergence, are met)

HESSIAN TECHNIQUE

ADVANTAGES AND DISADVANTAGES

- ▶ The Hessian technique takes into account *how fast the gradient changes*
- ▶ Theoretical and empirical evidence says that less iterations are needed
- ▶ Issue is the size of H , which of size N^2 if N is the number of parameters
- ▶ ⓘ Note that there could be $N = 10^7$ many parameters
- ▶ *Summary:*
 - ▶ The Hessian technique can often not be applied because computations are too expensive
 - ▶ However, it provided inspiration for other techniques

HESSIAN TECHNIQUE

ADVANTAGES AND DISADVANTAGES

- ▶ The Hessian technique takes into account *how fast the gradient changes*
- ▶ Theoretical and empirical evidence says that less iterations are needed
- ▶ Issue is the size of H , which is of size N^2 if N is the number of parameters
- ▶ 📖 Note that there could be $N = 10^7$ many parameters
- ▶ *Summary:*
 - ▶ The Hessian technique can often not be applied because computations are too expensive
 - ▶ However, it provided inspiration for other techniques

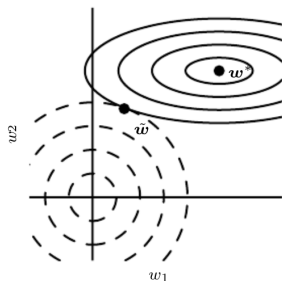
HESSIAN TECHNIQUE

ADVANTAGES AND DISADVANTAGES

- ▶ The Hessian technique takes into account *how fast the gradient changes*
- ▶ Theoretical and empirical evidence says that less iterations are needed
- ▶ Issue is the size of H , which is of size N^2 if N is the number of parameters
- ▶ 📌 Note that there could be $N = 10^7$ many parameters
- ▶ *Summary:*
 - ▶ The Hessian technique can often not be applied because computations are too expensive
 - ▶ However, it provided inspiration for other techniques

REGULARIZATION REVISITED

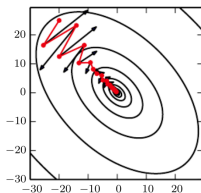
MOTIVATION



- ▶ *Reminder:* L2 regularization shrinks weights along Hessian eigenvectors
- ▶ The ball then moves as being pulled by the origin $(0,0)$ in the landscape induced by the eigenvectors of the Hessian

MOMENTUM-BASED GRADIENT DESCENT

MOTIVATION

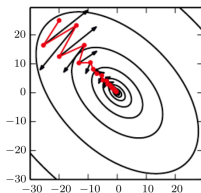


Black: gradients at each step, zig-zagging through the “valley”

- *Motivation:* Going back and forth, without making progress, during (stochastic) gradient descent
- *Reasons:*
 - Poorly conditioned Hessian matrix because of “valleys” (see Figure)
 - Variance between batches
 - High curvature in general
 - Noisy gradients

MOMENTUM-BASED GRADIENT DESCENT

MOTIVATION

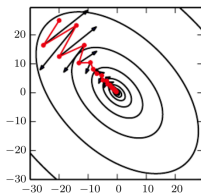


Black: gradients at each step, zig-zagging through the “valley”

- *Motivation:* Going back and forth, without making progress, during (stochastic) gradient descent
- *Reasons:*
 - Poorly conditioned Hessian matrix because of “valleys” (see Figure)
 - Variance between batches
 - High curvature in general
 - Noisy gradients

MOMENTUM-BASED GRADIENT DESCENT

MOTIVATION

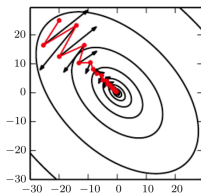


Black: gradients at each step, zig-zagging through the “valley”

- *Motivation:* Going back and forth, without making progress, during (stochastic) gradient descent
- *Reasons:*
 - Poorly conditioned Hessian matrix because of “valleys” (see Figure)
 - Variance between batches
 - High curvature in general
 - Noisy gradients

MOMENTUM-BASED GRADIENT DESCENT

MOTIVATION

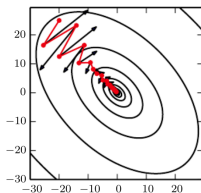


Black: gradients at each step, zig-zagging through the “valley”

- *Motivation:* Going back and forth, without making progress, during (stochastic) gradient descent
- *Reasons:*
 - Poorly conditioned Hessian matrix because of “valleys” (see Figure)
 - Variance between batches
 - High curvature in general
 - Noisy gradients

MOMENTUM-BASED GRADIENT DESCENT

MOTIVATION

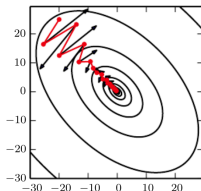


Black: gradients at each step, zig-zagging through the “valley”

- *Motivation:* Going back and forth, without making progress, during (stochastic) gradient descent
- *Reasons:*
 - Poorly conditioned Hessian matrix because of “valleys” (see Figure)
 - Variance between batches
 - High curvature in general
 - Noisy gradients

MOMENTUM-BASED GRADIENT DESCENT

SOLUTION

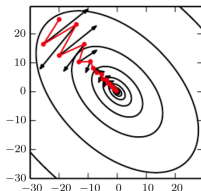


Red: averaged gradients, less zig-zagging

- Keep track of earlier gradients and take the average

MOMENTUM-BASED GRADIENT DESCENT

SOLUTION



Red: averaged gradients, less zig-zagging

- Formally, maintain *velocity* \mathbf{v} in addition to parameters \mathbf{w} themselves
- In each iteration, update
 - velocity, where α is the *momentum hyperparameter*

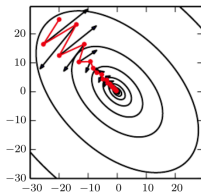
$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\mathbf{w}} C \quad (14)$$

- parameters:

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v} \quad (15)$$

MOMENTUM-BASED GRADIENT DESCENT

SOLUTION

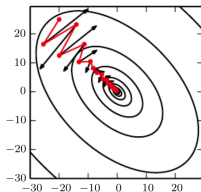


Red: averaged gradients, less zig-zagging

- ▶ In momentum based gradient descent, \mathbf{w} is an exponentially decaying average over gradients
- ▶ Variant: *Nesterov's accelerated gradient technique*
- ▶ See also <http://neuralnetworksanddeeplearning.com/chap3.html> (Nielsen, chapter 3) for more details

MOMENTUM-BASED GRADIENT DESCENT

SOLUTION



Red: averaged gradients, less zig-zagging

- ▶ In momentum based gradient descent, \mathbf{w} is an exponentially decaying average over gradients
- ▶ Variant: *Nesterov's accelerated gradient technique*
- ▶ See also <http://neuralnetworksanddeeplearning.com/chap3.html> (Nielsen, chapter 3) for more details

ALTERNATIVE OPTIMIZATION

FURTHER READING

- ▶ *Alternative methods:*
 - ▶ *Conjugate gradient descent*
 - ▶ *BFGS (Broyden-Fletcher-Goldfarb-Shanno) method*
 - ▶ *L-BFGS (Limited-memory-BFGS) method*
- ▶ *Illustrations / Literature:*
 - ▶ Bengio's deep learning book:
<http://www.deeplearningbook.org/contents/optimization.html>
 - ▶ "Efficient BackProp", Y. LeCun, L. Bottou, G. Orr, K.-R. Müller, 1998, see <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>
 - ▶ "On the importance of initialization and momentum in deep learning", I. Sutskever, J. Martens, G. Dahl and G. Hinton, 2012, <http://www.cs.toronto.edu/~hinton/absps/momentum.pdf>

ALTERNATIVE OPTIMIZATION

FURTHER READING

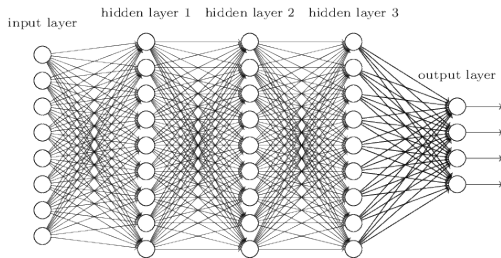
- ▶ *Alternative methods:*
 - ▶ *Conjugate gradient descent*
 - ▶ *BFGS (Broyden-Fletcher-Goldfarb-Shanno) method*
 - ▶ *L-BFGS (Limited-memory-BFGS) method*
- ▶ *Illustrations / Literature:*
 - ▶ Bengio's deep learning book:
<http://www.deeplearningbook.org/contents/optimization.html>
 - ▶ "Efficient BackProp", Y. LeCun, L. Bottou, G. Orr, K.-R. Müller, 1998, see <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>
 - ▶ "On the importance of initialization and momentum in deep learning", I. Sutskever, J. Martens, G. Dahl and G. Hinton, 2012, <http://www.cs.toronto.edu/~hinton/absps/momentum.pdf>

The Vanishing Gradient

WHY IS DEEP LEARNING TOUGH?

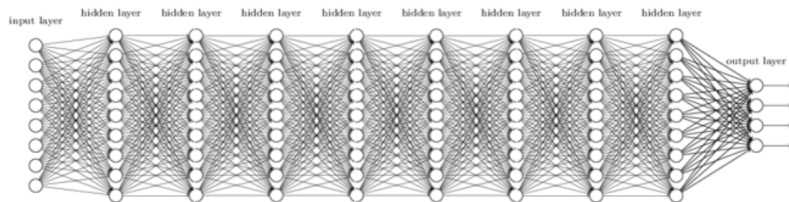
- ▶ Deep is supposed to be better than shallow
 - ▶ Less hidden nodes necessary to approximate the true functional relationship
 - ▶ See the “Universal Approximation Theorem” by Montufar, 2014
 - ▶ See further “Learning Deep Architectures”, Bengio, 2009, http://www.iro.umontreal.ca/~bengioy/papers/ftml_book.pdf for a more informal discussion
- ▶ *However*: On increasing depth in a naive way, performance usually drops
- ▶ What is going wrong?

WHY IS DEEP LEARNING TOUGH?

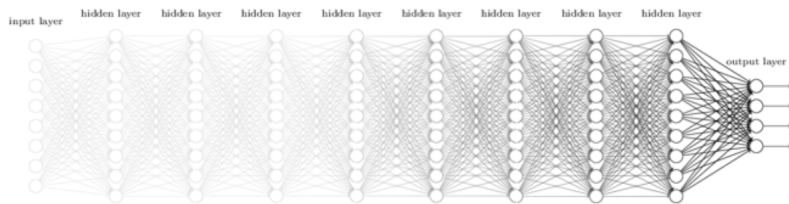


Training Deeper NN's: either the earlier layers (more common; here hidden layer 1) or the later layers (here: hidden layer 3) do not train well

THE VANISHING GRADIENT PROBLEM



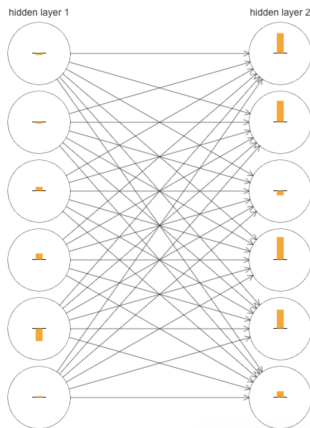
Deep Neural Network



Vanishing Gradient

Most commonly: gradients converge to zero in earlier layers

THE VANISHING GRADIENT PROBLEM



- ▶ Changes larger in later hidden layer
- ▶ Learning works better in later layers
- ▶ Are neurons likely to learn at different rates in different layers in general?

Yellow bars: $\frac{\partial C}{\partial b}$ for each hidden neuron

THE VANISHING GRADIENT PROBLEM

- ▶ Let b_j^l be the j -th bias in layer l , and $\frac{\partial C}{\partial b_j^l}$ be the respective partial derivative of the cost C .
- ▶ Let

$$\nabla_{\mathbf{b}^l} C := \left(\frac{\partial C}{\partial b_1^l}, \dots, \frac{\partial C}{\partial b_{d(l)}^l} \right) \quad (16)$$

- ▶ Then, in the example from the slide before:

$$\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.07 \text{ and } \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.31 \quad (17)$$

THE VANISHING GRADIENT PROBLEM

- ▶ Let b_j^l be the j -th bias in layer l , and $\frac{\partial C}{\partial b_j^l}$ be the respective partial derivative of the cost C .

- ▶ Let

$$\nabla_{\mathbf{b}^l} C := \left(\frac{\partial C}{\partial b_1^l}, \dots, \frac{\partial C}{\partial b_{d(l)}^l} \right) \quad (16)$$

- ▶ Then, in the example from the slide before:

$$\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.07 \quad \text{and} \quad \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.31 \quad (17)$$

THE VANISHING GRADIENT PROBLEM

- ▶ Then, in this example,

$$\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.07 \text{ and } \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.31 \quad (18)$$

- ▶ Formal quantification shows: learning faster in hidden layer 2.
- ▶ When running the identical training task (MNIST), we obtain
 - ▶ $\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.012, \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.06, \|\nabla_{\mathbf{b}}^{(3)} C\| = 0.283$
for three hidden layers
 - ▶ $\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.003, \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.017, \|\nabla_{\mathbf{b}}^{(3)} C\| = 0.07, \|\nabla_{\mathbf{b}}^{(4)} C\| = 0.285$ for four hidden layers
 - ▶ and so on...

THE VANISHING GRADIENT PROBLEM

- ▶ Then, in this example,

$$\|\nabla_{\mathbf{b}}^{(1)}C\| = 0.07 \text{ and } \|\nabla_{\mathbf{b}}^{(2)}C\| = 0.31 \quad (18)$$

- ▶ Formal quantification shows: learning faster in hidden layer 2.
- ▶ When running the identical training task (MNIST), we obtain
 - ▶ $\|\nabla_{\mathbf{b}}^{(1)}C\| = 0.012, \|\nabla_{\mathbf{b}}^{(2)}C\| = 0.06, \|\nabla_{\mathbf{b}}^{(3)}C\| = 0.283$
for three hidden layers
 - ▶ $\|\nabla_{\mathbf{b}}^{(1)}C\| = 0.003, \|\nabla_{\mathbf{b}}^{(2)}C\| = 0.017, \|\nabla_{\mathbf{b}}^{(3)}C\| = 0.07, \|\nabla_{\mathbf{b}}^{(4)}C\| = 0.285$ for four hidden layers
 - ▶ and so on...

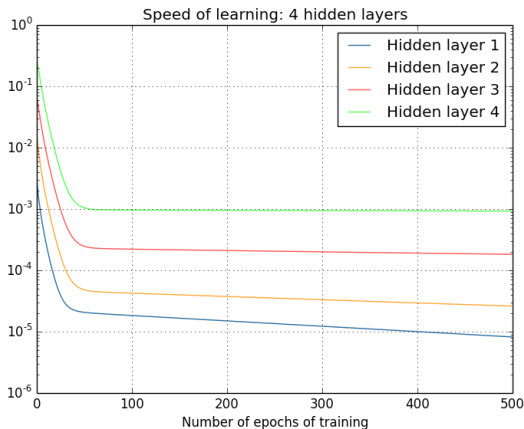
THE VANISHING GRADIENT PROBLEM

- ▶ Then, in this example,

$$\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.07 \text{ and } \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.31 \quad (18)$$

- ▶ Formal quantification shows: learning faster in hidden layer 2.
- ▶ When running the identical training task (MNIST), we obtain
 - ▶ $\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.012, \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.06, \|\nabla_{\mathbf{b}}^{(3)} C\| = 0.283$
for three hidden layers
 - ▶ $\|\nabla_{\mathbf{b}}^{(1)} C\| = 0.003, \|\nabla_{\mathbf{b}}^{(2)} C\| = 0.017, \|\nabla_{\mathbf{b}}^{(3)} C\| = 0.07, \|\nabla_{\mathbf{b}}^{(4)} C\| = 0.285$
for four hidden layers
 - ▶ and so on...

THE VANISHING GRADIENT PROBLEM



Training speed in [784,30,30,30,30,10]-NN on MNIST

THE VANISHING GRADIENT PROBLEM

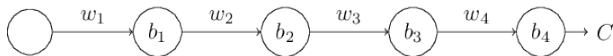
- ▶ *Vanishing gradient problem*: Neurons in earlier layers learn more slowly
- ▶ *Exploding gradient problem*: Neurons in earlier layers learn faster
- ▶ In general, gradients in NN's are unstable across layers
- ▶ And: vanishing gradients do not mean that there is nothing left to be learnt
- ▶ 📖 Fundamental problem for gradient-based learning in NN's

THE VANISHING GRADIENT PROBLEM

- ▶ *Vanishing gradient problem*: Neurons in earlier layers learn more slowly
- ▶ *Exploding gradient problem*: Neurons in earlier layers learn faster
- ▶ In general, gradients in NN's are unstable across layers
- ▶ And: vanishing gradients do not mean that there is nothing left to be learnt
- ▶ 🖱 Fundamental problem for gradient-based learning in NN's

THE VANISHING GRADIENT PROBLEM

EXPLANATION

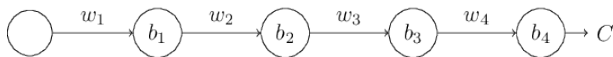


Simple NN with 3 hidden layers of one neuron each

Let w_1, w_2, w_3, w_4 be the weights, b_1, b_2, b_3, b_4 be the biases and C the cost. Let all neurons be sigmoid, so the output a_j from the j -th neuron is $\sigma(z_j)$ where $z_j = w_j a_{j-1} + b_j$ is the input of the j -th neuron (notation as usual earlier).

THE VANISHING GRADIENT PROBLEM

EXPLANATION



Simple NN with 3 hidden layers of one neuron each

Let w_1, w_2, w_3, w_4 be the weights, b_1, b_2, b_3, b_4 be the biases and C the cost. Let all neurons be sigmoid, so the output a_j from the j -th neuron is $\sigma(z_j)$ where $z_j = w_j a_{j-1} + b_j$ is the input of the j -th neuron (notation as usual earlier).

For understanding the Vanishing Gradient Problem, consider $\frac{\partial C}{\partial b_1}$. By repeated application of the backpropagation rules, we see that

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4} \quad (19)$$

THE VANISHING GRADIENT PROBLEM

EXPLANATION

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Computing $\frac{\partial C}{\partial b_1}$

There is an alternative explanation for (19). Let Δ indicate small changes. We know that

$$\frac{\partial C}{\partial b_1} \approx \frac{\Delta C}{\Delta b_1} \quad (20)$$

THE VANISHING GRADIENT PROBLEM

EXPLANATION

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Computing $\frac{\partial C}{\partial b_1}$

There is an alternative explanation for (19). Let Δ indicate small changes. We know that

$$\frac{\partial C}{\partial b_1} \approx \frac{\Delta C}{\Delta b_1} \quad (20)$$

From $a_1 = \sigma(z_1) = \sigma(w_1 a_0 + b_1)$ we further obtain

$$\Delta a_1 \approx \frac{\partial \sigma(w_1 a_0 + b_1)}{\partial b_1} \Delta b_1 = \sigma'(z_1) \Delta b_1 \quad (21)$$

THE VANISHING GRADIENT PROBLEM

EXPLANATION

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Computing $\frac{\partial C}{\partial b_1}$

There is an alternative explanation for (19). Let Δ indicate small changes. We know that

$$\frac{\partial C}{\partial b_1} \approx \frac{\Delta C}{\Delta b_1} \quad (20)$$

From $a_1 = \sigma(z_1) = \sigma(w_1 a_0 + b_1)$ we further obtain

$$\Delta a_1 \approx \frac{\partial \sigma(w_1 a_0 + b_1)}{\partial b_1} \Delta b_1 = \sigma'(z_1) \Delta b_1 \quad (21)$$

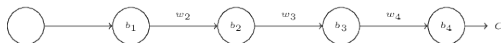
further leading to

$$\Delta z_2 \approx \frac{\partial z_2}{\partial a_1} \Delta a_1 = w_2 \Delta a_1 \quad \text{implying} \quad \Delta z_2 \approx \sigma'(z_1) w_2 \Delta b_1 \quad (22)$$

THE VANISHING GRADIENT PROBLEM

EXPLANATION

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Computing $\frac{\partial C}{\partial b_1}$

Repeated application of the computations from the slide before eventually yield

$$\Delta C \approx \sigma'(z_1)w_2\sigma'(z_2) \dots \sigma'(z_4) \frac{\partial C}{\partial a_4} \Delta b_1 \quad (23)$$

Dividing by b_1 results in the desired expression (19):

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4} \quad (24)$$

THE VANISHING GRADIENT PROBLEM

EXPLANATION

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4} \quad (25)$$

Except from the last term, this is a product of terms of the form

$$w_j \sigma'(z_j) \quad (26)$$

It holds that $0 \leq \sigma'(z_j) \leq 1/4$, while, in practice, when employing standard initialization of weights, typically $|w_j| < 1$, so

$$|w_j \sigma'(z_j)| \leq \frac{1}{4} \quad (27)$$

so in combination

$$\frac{\partial C}{\partial b_1} \leq \sigma'(z_1) \left(\frac{1}{4}\right)^3 \frac{\partial C}{\partial a_4} \quad (28)$$

THE VANISHING GRADIENT PROBLEM

EXPLANATION

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \sigma'(z_1) \overbrace{w_2 \sigma'(z_2)}^{< \frac{1}{4}} \overbrace{w_3 \sigma'(z_3)}^{< \frac{1}{4}} \underbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}_{\text{common terms}} \\ &\quad \updownarrow \\ \frac{\partial C}{\partial b_3} &= \sigma'(z_3) \overbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}^{\text{common terms}} \end{aligned}$$

Comparing $\frac{\partial C}{\partial b_1}$ with $\frac{\partial C}{\partial b_3}$

So, $\frac{\partial C}{\partial b_1}$ is about a factor of 16 (or more) smaller than $\frac{\partial C}{\partial b_3}$. Similar conclusions are drawn for $\frac{\partial C}{\partial w_j}$.

THE EXPLODING GRADIENT PROBLEM

The “Exploding Gradient Problem” occurs when

- ▶ Weights are too large (say on the order of 100 each)
- ▶ Biases b_j are such that $\sigma'(z_j)$ never is small
- ▶ *Example:* $b_j = -100 \times a_{j-1}$, so $z_j = 100 \times a_{j-1} - 100 \times a_{j-1} = 0$, implying $\sigma'(z_j) = 1/4$, yielding $w_j \sigma'(z_j) > 20$ as a gradient
- ▶ In such situations gradients iteratively explode

THE EXPLODING GRADIENT PROBLEM

The “Exploding Gradient Problem” occurs when

- ▶ Weights are too large (say on the order of 100 each)
- ▶ Biases b_j are such that $\sigma'(z_j)$ never is small
- ▶ *Example:* $b_j = -100 \times a_{j-1}$, so $z_j = 100 \times a_{j-1} - 100 \times a_{j-1} = 0$, implying $\sigma'(z_j) = 1/4$, yielding $w_j \sigma'(z_j) > 20$ as a gradient
- ▶ In such situations gradients iteratively explode

THE EXPLODING GRADIENT PROBLEM

The “Exploding Gradient Problem” occurs when

- ▶ Weights are too large (say on the order of 100 each)
- ▶ Biases b_j are such that $\sigma'(z_j)$ never is small
- ▶ *Example:* $b_j = -100 \times a_{j-1}$, so $z_j = 100 \times a_{j-1} - 100 \times a_{j-1} = 0$, implying $\sigma'(z_j) = 1/4$, yielding $w_j \sigma'(z_j) > 20$ as a gradient
- ▶ In such situations gradients iteratively explode

GRADIENTS ARE UNSTABLE

- ▶ The fundamental problem is that gradients in earlier layers are products of gradients from (all the) later layers.
- ▶ If there are many layers, the situation is unstable, unless the gradients are *balanced out*.
- ▶ Balancing is very unlikely to happen by chance, so one needs to fix this explicitly.
- ▶ Fixing this seems daunting at first glance: when making weights w_j large,

$$\sigma'(z_j) = \sigma'(w_j a_{j-1} + b_j)$$

will get small.

- ▶ *Solutions:*

- ▶ *Rectified Linear Units* instead of sigmoid activation
- ▶ *Batch Normalization* (discussed later in the lecture)

GRADIENTS ARE UNSTABLE

- ▶ The fundamental problem is that gradients in earlier layers are products of gradients from (all the) later layers.
- ▶ If there are many layers, the situation is unstable, unless the gradients are *balanced out*.
- ▶ Balancing is very unlikely to happen by chance, so one needs to fix this explicitly.
- ▶ Fixing this seems daunting at first glance: when making weights w_j large,

$$\sigma'(z_j) = \sigma'(w_j a_{j-1} + b_j)$$

will get small.

- ▶ *Solutions:*

- ▶ *Rectified Linear Units* instead of sigmoid activation
- ▶ *Batch Normalization* (discussed later in the lecture)

GRADIENTS ARE UNSTABLE

- ▶ The fundamental problem is that gradients in earlier layers are products of gradients from (all the) later layers.
- ▶ If there are many layers, the situation is unstable, unless the gradients are *balanced out*.
- ▶ Balancing is very unlikely to happen by chance, so one needs to fix this explicitly.
- ▶ Fixing this seems daunting at first glance: when making weights w_j large,

$$\sigma'(z_j) = \sigma'(w_j a_{j-1} + b_j)$$

will get small.

- ▶ *Solutions:*
 - ▶ *Rectified Linear Units* instead of sigmoid activation
 - ▶ *Batch Normalization* (discussed later in the lecture)

WHY IS DEEP LEARNING TOUGH?

LITERATURE

There are other issues that prevent easy training of neural networks with deep architectures. For further reading, see for example

- ▶ “Understanding the difficulty of training deep feedforward neural networks”, X. Glorot, Y. Bengio, 2010,
<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

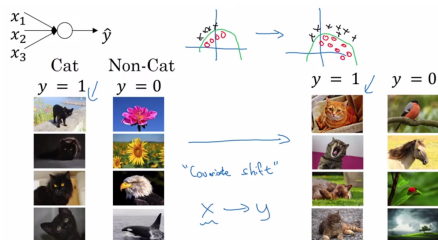
or the earlier

- ▶ “Efficient BackProp”, Y. LeCun, L. Bottou, G. Orr, K.-R. Müller, 1998, <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>
- ▶ “On the importance of initialization and momentum in deep learning”, I. Sutskever, J. Martens, G. Dahl, G. Hinton, 2013, <http://www.cs.toronto.edu/~hinton/absps/momentum.pdf>

Batch Normalization

BATCH NORMALIZATION

MOTIVATION

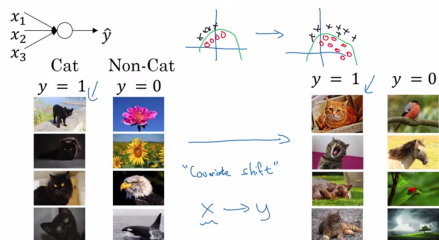


Learning black cats might not help to recognize cats of other colors

- ▶ The network might not be able to predict well if presented with examples not present in the training data (batch)
- ▶ The function learned can only be guaranteed to predict well in certain areas of feature space

BATCH NORMALIZATION

MOTIVATION

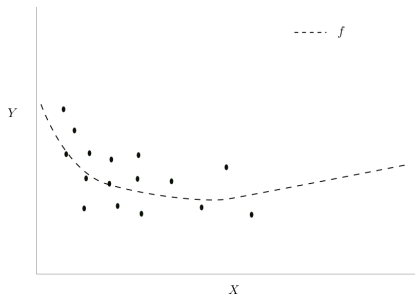


Learning black cats might not help to recognize cats of other colors

- ▶ The network might not be able to predict well if presented with examples not present in the training data (batch)
- ▶ The function learned can only be guaranteed to predict well in certain areas of feature space

BATCH NORMALIZATION

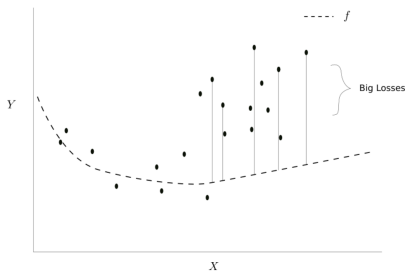
MOTIVATION



- ▶ Let f be the function that maps values from layer $l - 1$ to layer l
- ▶ After updating gradients, values in layer $l - 1$ may have been shifted to a region where f does not approximate the true function well
- ▶ The effect is referred to as *internal covariate shift* (although this is not necessarily the correct term to describe the effect)
- ▶ This slows down training

BATCH NORMALIZATION

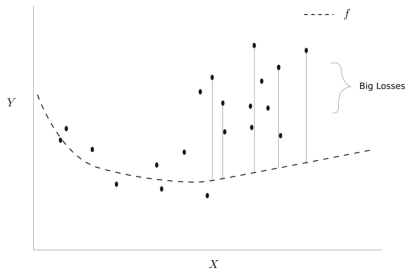
MOTIVATION



- ▶ Let f be the function that maps values from layer $l - 1$ to layer l
- ▶ After updating gradients, values in layer $l - 1$ may have been shifted to a region where f does not approximate the true function well
- ▶ The effect is referred to as *internal covariate shift* (although this is not necessarily the correct term to describe the effect)
- ▶ This slows down training

BATCH NORMALIZATION

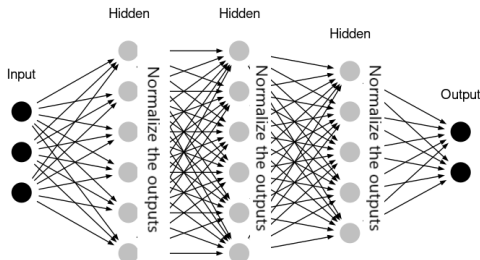
MOTIVATION



- ▶ Let f be the function that maps values from layer $l - 1$ to layer l
- ▶ After updating gradients, values in layer $l - 1$ may have been shifted to a region where f does not approximate the true function well
- ▶ The effect is referred to as *internal covariate shift* (although this is not necessarily the correct term to describe the effect)
- ▶ This slows down training

BATCH NORMALIZATION

SOLUTION



Batch Normalization: Insert normalization layers between normal-type layers

- ▶ After each layer, normalize output values
- ▶ There are parameters to be learned for normalization layers
- ▶ Parameters for normalization layers can be easily learnt with backpropagation

BATCH NORMALIZATION

DEFINITION

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

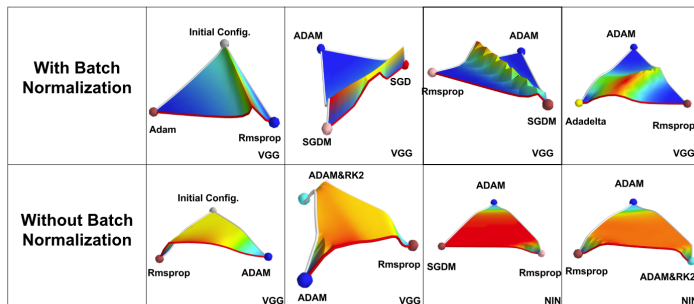
<https://arxiv.org/abs/1502.03167>

(Ioffe & Szegedy, original paper)

- Compute \hat{x}_i when forwarding training samples
- Learn γ, β during backpropagation

BATCH NORMALIZATION

EXPLANATION



[From: <http://www.aifounded.com/machine-learning/deep-loss>]

- ▶ Low error regions are larger
- ▶ Boundaries are more clearly / sharply defined
- ▶ The reshaping of the cost function surface leads to accelerated training

BATCH NORMALIZATION

SUMMARY BENEFITS

- ▶ *Gradient Vanishing*: Batch Normalization prevents gradients from vanishing
- ▶ *Internal Covariate Shift*: controversial debate whether it helps (although it is motivated by it)
- ▶ Boundaries of error regions are more clearly / sharply defined
- ▶ Reshapes cost function surface: accelerated training

BATCH NORMALIZATION

SUMMARY BENEFITS

- ▶ *Gradient Vanishing*: Batch Normalization prevents gradients from vanishing
- ▶ *Internal Covariate Shift*: controversial debate whether it helps (although it is motivated by it)
- ▶ Boundaries of error regions are more clearly / sharply defined
- ▶ Reshapes cost function surface: accelerated training

SUMMARY

- ▶ Convolutional backpropagation

- ▶ For further illustrations, see <https://medium.com/@2017csm1006/forward-and-backpropagation-in-convolutional-neural-network->
- ▶ Note that the notation differs (error E there is cost C here, X there is z here, and F are the weights w here, and O is a here, and there is no bias)

- ▶ Training variations

- ▶ <http://www.deeplearningbook.org/>, Chapter 8 (corresponding parts)
- ▶ <http://neuralnetworksanddeeplearning.com/>, Chapter 3, “Variations on stochastic gradient descent”

- ▶ The vanishing gradient problem

- ▶ <http://neuralnetworksanddeeplearning.com/>, Chapter 5

- ▶ Batch normalization

- ▶ See <http://www.deeplearningbook.org/>, 8.7.1
- ▶ See also <http://www.aifounded.com/machine-learning/deep-loss>, for example

Thanks for your attention