

Going Deep

Kate Farrahi

ECS Southampton

February 21, 2019

The Universal Approximation Theorem

Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, bounded, and continuous function.

The Universal Approximation Theorem

Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, bounded, and continuous function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$.

The Universal Approximation Theorem

Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, bounded, and continuous function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$. The space of real-valued continuous functions on I_m is denoted by $C(I_m)$.

The Universal Approximation Theorem

Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, bounded, and continuous function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$. The space of real-valued continuous functions on I_m is denoted by $C(I_m)$. Then, given any $\epsilon > 0$ and any function $f \in C(I_m)$, there exist an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$, such that we may define:

$F(x) = \sum_{i=1}^N v_i \psi(w_i^T x + b_i)$ as an approximate realization of the function f ; that is,

The Universal Approximation Theorem

Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, bounded, and continuous function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$. The space of real-valued continuous functions on I_m is denoted by $C(I_m)$. Then, given any $\epsilon > 0$ and any function $f \in C(I_m)$, there exist an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$, such that we may define:

$F(x) = \sum_{i=1}^N v_i \psi(w_i^T x + b_i)$ as an approximate realization of the function f ; that is,

$$|F(x) - f(x)| < \epsilon$$

for all $x \in I_m$.

Then Why Go Deep?

- ▶ There are functions you can compute with a deep neural network that shallow networks require exponentially more hidden units to compute.
- ▶ The following function is more efficient to implement using a deep neural network: $y = x_1 \oplus x_2 \oplus x_3 \oplus \cdots \oplus x_n$

Vanishing and Exploding Gradients

- ▶ The vanishing and exploding gradient problem is a difficulty found in training NN with gradient-based learning methods and backpropagation.

Vanishing and Exploding Gradients

- ▶ The vanishing and exploding gradient problem is a difficulty found in training NN with gradient-based learning methods and backpropagation.
- ▶ In training, the gradient may become vanishingly small (or large), effectively preventing the weight from changing its value (or exploding in value).

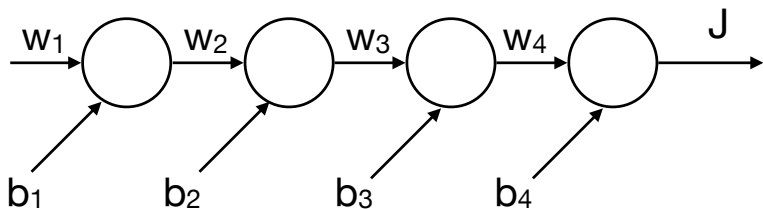
Vanishing and Exploding Gradients

- ▶ The vanishing and exploding gradient problem is a difficulty found in training NN with gradient-based learning methods and backpropagation.
- ▶ In training, the gradient may become vanishingly small (or large), effectively preventing the weight from changing its value (or exploding in value).
- ▶ This leads to the neural network not being able to train.

Vanishing and Exploding Gradients

- ▶ The vanishing and exploding gradient problem is a difficulty found in training NN with gradient-based learning methods and backpropagation.
- ▶ In training, the gradient may become vanishingly small (or large), effectively preventing the weight from changing its value (or exploding in value).
- ▶ This leads to the neural network not being able to train.
- ▶ This issue affects many-layered networks (feed-forward), as well as recurrent networks.

Issues with Going Deep



Residual Connections

- ▶ One of the most effective ways to resolve the vanishing gradient problem is with residual neural networks (ResNets)¹.

¹K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," CVPR, Las Vegas, NV, 2016, pp. 770-778.

Residual Connections

- ▶ One of the most effective ways to resolve the vanishing gradient problem is with residual neural networks (ResNets)¹.
- ▶ ResNets are artificial neural networks that use *skip connections* to jump over layers.

¹K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," CVPR, Las Vegas, NV, 2016, pp. 770-778.

Residual Connections

- ▶ One of the most effective ways to resolve the vanishing gradient problem is with residual neural networks (ResNets)¹.
- ▶ ResNets are artificial neural networks that use *skip connections* to jump over layers.
- ▶ The vanishing gradient problem is mitigated in ResNets by reusing activations from a previous layer.

¹K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," CVPR, Las Vegas, NV, 2016, pp. 770-778.

Residual Connections

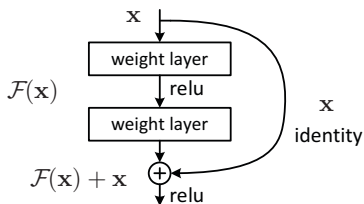


Figure 2. Residual learning: a building block.²

²K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," CVPR, Las Vegas, NV, 2016, pp. 770-778.

Residual Connections

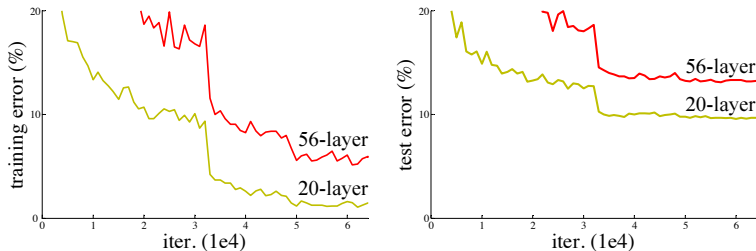


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

3

³K. He, X. Zhang, S. Ren and J. Sun, “Deep Residual Learning for Image Recognition,” CVPR, Las Vegas, NV, 2016, pp. 770-778.

Residual Connections

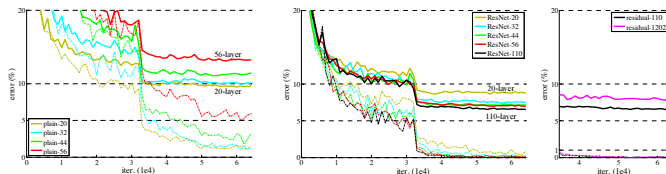


Figure 6. Training on **CIFAR-10**. Dashed lines denote training error, and bold lines denote testing error. **Left:** plain networks. The error of plain-110 is higher than 60% and not displayed. **Middle:** ResNets. **Right:** ResNets with 110 and 1202 layers.

4

⁴K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," CVPR, Las Vegas, NV, 2016, pp. 770-778.

Dropout

- ▶ Neural networks with a large number of parameters (and hidden layers) are powerful, however, overfitting is a serious problem in such systems.

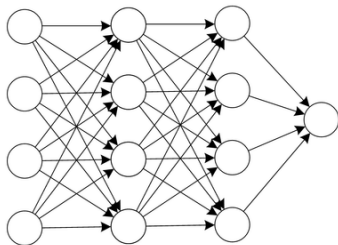
Dropout

- ▶ Neural networks with a large number of parameters (and hidden layers) are powerful, however, overfitting is a serious problem in such systems.
- ▶ Dropout is a form of regularization

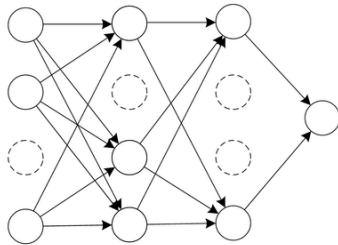
Dropout

- ▶ Neural networks with a large number of parameters (and hidden layers) are powerful, however, overfitting is a serious problem in such systems.
- ▶ Dropout is a form of regularization
- ▶ The key idea in dropout is to randomly drop neurons, including all of the connections, from the neural network during training.

Dropout



(a) Standard Neural Network



(b) Network after Dropout

5

⁵Image from: https://www.researchgate.net/figure/Dropout-neural-network-model-a-is-a-standard-neural-network-b-is-the-same-model-with-dropout-fig3_309206911

How Does Dropout Work in Practice?

- ▶ In the learning phase, we stochastically remove hidden units by setting a dropout probability for each layer in the network. We then randomly decide whether or not a neuron in a given layer is removed stochastically.

How Does Dropout Work in Practice?

- ▶ We define a random binary mask $m^{(l)}$ which is used to remove neurons, and note, $m^{(l)}$ changes for each iteration of the backpropagation algorithm.

How Does Dropout Work in Practice?

- ▶ We define a random binary mask $m^{(l)}$ which is used to remove neurons, and note, $m^{(l)}$ changes for each iteration of the backpropagation algorithm.
- ▶ For layers, $l = 1$ to $L - 1$, for the forward pass of backpropagation, we then compute

$$a^{(l)} = \sigma(w^{(l)} a^{(l-1)} + b^{(l)}) \odot m^{(l)} \quad (1)$$

How Does Dropout Work in Practice?

- ▶ We define a random binary mask $m^{(l)}$ which is used to remove neurons, and note, $m^{(l)}$ changes for each iteration of the backpropagation algorithm.
- ▶ For layers, $l = 1$ to $L - 1$, for the forward pass of backpropagation, we then compute

$$a^{(l)} = \sigma(w^{(l)} a^{(l-1)} + b^{(l)}) \odot m^{(l)} \quad (1)$$

- ▶ For layer L , $a^{(L)} = \sigma(w^{(L)} a^{(L-1)} + b^{(L)})$

How Does Dropout Work in Practice?

- ▶ We define a random binary mask $m^{(l)}$ which is used to remove neurons, and note, $m^{(l)}$ changes for each iteration of the backpropagation algorithm.
- ▶ For layers, $l = 1$ to $L - 1$, for the forward pass of backpropagation, we then compute

$$a^{(l)} = \sigma(w^{(l)} a^{(l-1)} + b^{(l)}) \odot m^{(l)} \quad (1)$$

- ▶ For layer L , $a^{(L)} = \sigma(w^{(L)} a^{(L-1)} + b^{(L)})$
- ▶ For the backward pass of the backpropagation algorithm,

$$\delta^L = \Delta_a J \odot \sigma'(z^L) \odot m^{(l)} \quad (2)$$

Why Does Dropout Work?

- ▶ Neurons cannot co-adapt to other units (they cannot assume that all of the other units will be present)

Why Does Dropout Work?

- ▶ Neurons cannot co-adapt to other units (they cannot assume that all of the other units will be present)
- ▶ By breaking co-adaptation, each unit will ultimately find more general features