The Reversible Residual Network

Aidan Gomez, Mengye Ren, Raquel Urtasun and Roger Grosse





What are Deep Neural Networks (DNNs)?

- A 'deep' composition of functions.
- Generally, each function is parametrized by some weights
- These weights can then be augmented to manipulate the computation being performed.

$$f_1 \circ f_2 \circ \ldots \circ f_n$$

Single layer of an MLP $f_k(x, W_{f_k}) = \varphi(W_{f_k}^\top x)$

Convolution $f_k(x, W_{f_k}) = \varphi(x * W_{f_k})$



- 'gradient descent'
- Take steps in local direction of steepest descent in the loss

 $L = \text{Loss}(f_1 \circ \cdots \circ f_n, y)$

 Network weights are updated using 'gradient descent'

$$f_1 \circ f_2 \circ \ldots \circ f_n$$

- Take steps in local direction of steepest descent in the loss
- Chain-rule expansion can be recursively computed using algorithm known as backpropagation

 $L = \operatorname{Loss}(f_1 \circ \cdots \circ f_n, y)$

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$
$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$



- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$



- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$



- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$



- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$



- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$



- Chain-rule expansion can be recursively computed using algorithm known as backpropagation
- 'Backprop' occurs in two steps:
 - A forward pass that computes the output of each layer.
 - A backward pass the competes the gradients of each layer.
- The partial derivates w.r.t. the weights of each layer rely on the gradients of layer outputs, and the inputs to the layer.
- The consequence is the need for storing all layer activations at each

$$\nabla_{f_k} L = \frac{\partial L}{\partial f_1} \cdots \frac{\partial f_{k-2}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

$$\nabla_{f_k} L = \nabla_{f_{k-1}} L \cdot \frac{\partial f_{k-1}}{\partial f_k}$$

Data

$$\partial L$$

 ∂L
 $\partial Data$
 ∂f_2
 $\partial Data$
 ∂f_2
 $\partial Data$
 ∂f_2
 $\partial Data$
 ∂f_1
 ∂f_2
 ∂f_1
 ∂f_2
 ∂f_1
 ∂f_2
 ∂f_1
 ∂f_2

Training Gets Expensive

- The observation in deep networks is: `deep = better`
- Modern networks can have hundreds of layers each with large activation maps
- GPUs tend to have 8-16GB in memory and are easily exhausted by high dimensional data.

Residual Networks (ResNets)

- Class of layer originally developed for vision tasks
- Learns a perturbation of its input.
- Extremely successful in multiple domains: vision, text and audio.



Reversible Residual Networks

- By partitioning the graph and learning a pair of residuals we can enable exact reconstruction of inputs.
- Given a large stack of these layers, preserving activations is no longer necessary.
- Simply reconstruct activations from outputs during backprop.



Reversible Residual Networks (RevNet)





$$egin{aligned} y_1 &= x_1 + \mathcal{F}(x_2) \ y_2 &= x_2 + \mathcal{G}(y_1) \end{aligned}$$



 $egin{aligned} &x_2 = y_2 - \mathcal{G}(y_1) \ &x_1 = y_1 - \mathcal{F}(x_2) \end{aligned}$

Reversible Residual Networks (RevNet)





 $egin{aligned} y_1 &= x_1 + \mathcal{F}(x_2) \ y_2 &= x_2 + \mathcal{G}(y_1) \end{aligned}$



 $egin{aligned} &x_2 = y_2 - \mathcal{G}(y_1) \ &x_1 = y_1 - \mathcal{F}(x_2) \end{aligned}$

Reversible Residual Networks (RevNet)





$$egin{aligned} y_1 &= x_1 + \mathcal{F}(x_2) \ y_2 &= x_2 + \mathcal{G}(y_1) \end{aligned}$$



 $egin{aligned} &x_2 = y_2 - \mathcal{G}(y_1) \ &x_1 = y_1 - \mathcal{F}(x_2) \end{aligned}$

Results

Architecture	CIFAR-10 [15]		CIFAR-100 [15]	
	ResNet	RevNet	ResNet	RevNet
32 (38)	7.14%	7.24%	29.95%	28.96%
110	5.74%	5.76%	26.44%	25.40%
164	5.24%	5.17%	23.37%	23.69%
	ResNet-101 RevNet-104			

23.10%

- No noticeable loss in performance
- Massive gains in network depths (~600 vs ~100 layers on ImageNet, using a single GPU)

23.01%

- 4x increase in batch size (128 vs 32)
- only 1.5x the runtime cost of normal training (in practice, sometimes even less)

Future Work

- Recurrent Neural Networks
- Application to tasks requiring large activation maps (image segmentation, video, etc.)
- Implement your network within our codebase! github.com/renmengye/revnet-public