

Chain Rule and Backprop

Sargur N. Srihari
srihari@buffalo.edu

Topics (Deep Feedforward Networks)

- Overview
 1. Example: Learning XOR
 2. Gradient-Based Learning
 3. Hidden Units
 4. Architecture Design
 5. Backpropagation and Other Differentiation Algorithms
 6. Historical Notes

Topics in Backpropagation

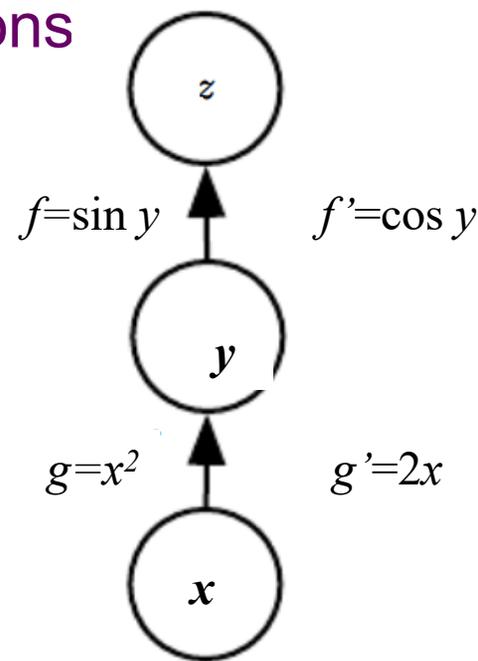
- Overview
 1. Computational Graphs
 2. Chain Rule of Calculus
 3. Recursively applying the chain rule to obtain backprop
 4. Backpropagation computation in fully-connected MLP
 5. Symbol-to-symbol derivatives
 6. General backpropagation
 7. Ex: backpropagation for MLP training
 8. Complications
 9. Differentiation outside the deep learning community
 10. Higher-order derivatives

Chain Rule of Calculus

- The chain rule states that derivative of $f(g(x))$ is $f'(g(x)) \cdot g'(x)$
 - It helps us differentiate *composite functions*
 - Note that $\sin(x^2)$ is composite, but $\sin(x) \cdot x^2$ is not
 - $\sin(x^2)$ is a composite function because it can be constructed as $f(g(x))$ for $f(x)=\sin(x)$ and $g(x)=x^2$
 - Using the chain rule and the derivatives of $\sin(x)$ and x^2 , we can then find the derivative of $\sin(x^2)$
 - Since $f(g(x))=\sin(x^2)$ it can be written as $f(y)=\sin y$ and $y=g(x)=x^2$
 - Since $f'(y)=df(y)/dy=\cos y$ and $g'(x)=dg(x)/dx=2x$,
 - we have the derivative of $\sin(x^2)$ as $\cos(y) \cdot 2x$ which is equivalent to $\cos(x^2) \cdot 2x$

Calculus' Chain Rule for Scalars

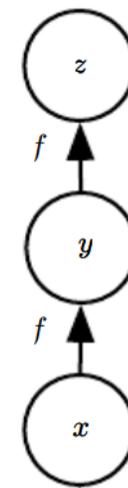
- For computing derivatives of functions formed by composing other functions whose derivatives are known
 - Back-prop computes chain rule, with a specific order of operations that is highly efficient
 - Let x be a real number, f, g be functions
 - mapping from a real no. to a real no.
 - If $y=g(x)$ and $z=f(g(x))=f(y)$
- Then chain rule $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$
- If $z=f(y)=\sin(y)$ and $y=g(x)=x^2$
 $dz/dy = \cos y$ and $dy/dx = 2x$



For vectors we need a Jacobian Matrix

- For a vector output $\mathbf{y} = \{y_1, \dots, y_m\}$ with vector input $\mathbf{x} = \{x_1, \dots, x_n\}$,
- Jacobian matrix organizes all the partial derivatives into an $m \times n$ matrix

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \quad J_{ki} = \frac{\partial y_k}{\partial x_i}$$



Determinant of Jacobian Matrix is referred to simply as the Jacobian

Generalizing Chain Rule to Vectors

- Suppose $\mathbf{x} \in R^m, \mathbf{y} \in R^n$ g maps from R^m to R^n and f from R^n to R

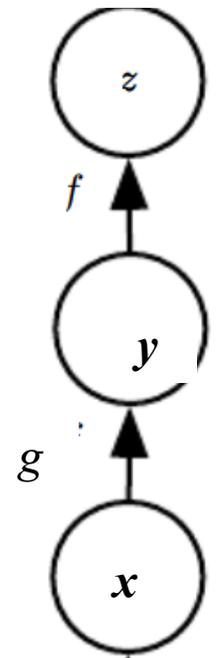
- If $\mathbf{y}=g(\mathbf{x})$ and $z=f(\mathbf{y})$ then

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

– In vector notation this is

- where $\left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right)$ is the $n \times m$ Jacobian matrix of g

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right)^T \nabla_{\mathbf{y}} z$$



- Thus gradient of z wrt \mathbf{x} is product of:

– Jacobian matrix $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ and gradient vector $\nabla_{\mathbf{y}} z$

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

$$\nabla_{\mathbf{y}} z = \begin{bmatrix} \frac{\partial z}{\partial y_1} \\ \vdots \\ \frac{\partial z}{\partial y_n} \end{bmatrix}$$

- Backprop algorithm consists of performing *Jacobian-gradient* product for each step of graph

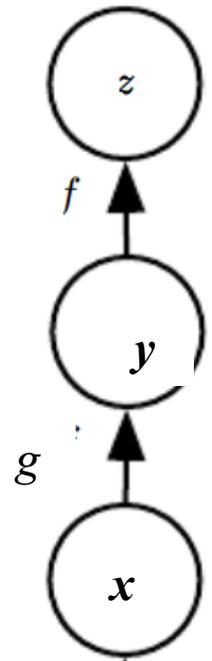
Generalizing Chain Rule to Tensors

- Backpropagation is usually applied to tensors with arbitrary dimensionality
- This is exactly the same as with vectors
 - Only difference is how numbers are arranged in a grid to form a tensor
 - We could flatten each tensor into a vector, compute a vector-valued gradient and reshape it back to a tensor
- In this view backpropagation is still multiplying Jacobians by gradients

$$\nabla_x z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_y z$$

Chain Rule for tensors

- To denote gradient of value z wrt a tensor we write $\boxed{\nabla_X z}$ as if X were a vector
- For 3-D tensor, X has three coordinates
 - We can abstract this away by using a single variable i to represent complete tuple of indices
- For all possible tuples i $\boxed{(\nabla_X z)_i}$ gives $\frac{\partial z}{\partial X_i}$
 - Exactly same as how for all possible indices i into a vector, $\boxed{(\nabla_X z)_i}$ gives $\frac{\partial z}{\partial X_i}$
- Chain rule for tensors
 - If $Y = g(X)$ and $z = f(Y)$ then



$$\boxed{\nabla_X(z) = \sum_j (\nabla_X Y_j) \frac{\partial z}{\partial Y_j}}$$

3. Recursively applying the chain rule to obtain backprop

Backprop is Recursive Chain Rule

- Backprop is obtained by recursively applying the chain rule
- Using chain rule it is straightforward to write expression for gradient of a scalar wrt any node in graph for producing that scalar
- However, evaluating that expression on a computer has some extra considerations
 - E.g., many sub-expressions may be repeated several times within overall expression
 - Whether to store sub-expressions or recompute them

Example of repeated subexpressions

- Let w be the input to the graph
 - We use same function $f: R \rightarrow R$ at every step:

$$x=f(w), y=f(x), z=f(y)$$

- To compute $\frac{\partial z}{\partial w}$ apply

$\begin{aligned}\frac{\partial z}{\partial w} &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \\ &= f'(y)f'(x)f'(w) \quad (1) \\ &= f'(f(f(w)))f'(f(w))f'(w) \quad (2)\end{aligned}$

Eq. (1): compute $f(w)$ once and store it in x

- This the approach taken by backprop

Eq. (2): expression $f(w)$ appears more than once

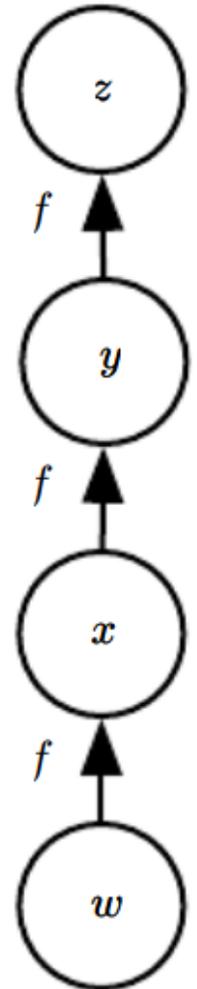
$f(w)$ is recomputed each time it is needed

For low memory, (1) preferable: reduced runtime

(2) is also valid chain rule, useful for limited memory

For complicated graphs,

exponentially wasted computations making naive implementation of chain rule infeasible



Simplified Backprop Algorithm

- Version that directly computes actual gradient
 - In the order it will actually be done according to recursive application of chain rule
 - Algorithm *Simplified Backprop* along with associated
 - *Forward Propagation*
- Could either directly perform these operations
 - or view algorithm as symbolic specification of computational graph for computing the back-prop
- This formulation does not make specific
 - Manipulation and construction of symbolic graph that performs gradient computation

Computing a single scalar

- Consider computational graph of how to compute a single scalar $u^{(n)}$
 - say loss on a training example
- | | | |
|--------|---------|---------|
| | n | |
| | n_1+1 | \dots |
| n_2 | | |
| inputs | 1 | \dots |
| | n_1 | |
- We want gradient of $u^{(n)}$ wrt n_i inputs $u^{(1)}, \dots, u^{(n_i)}$
 - i.e., we wish to compute $\frac{\partial u^{(n)}}{\partial u_i}$ for all $i = 1, \dots, n_i$
 - In application of backprop to computing gradients for gradient descent over parameters
 - $u^{(n)}$ will be cost associated with an example or a minibatch, while
 - $u^{(1)}, \dots, u^{(n_i)}$ correspond to model parameters

Nodes of Computational Graph

- Assume that nodes of the graph have been ordered such that
 - We can compute their output one after another
 - Starting at $u^{(n_i+1)}$ and going to $u^{(n)}$
- As defined in Algorithm shown next
 - Each node $u^{(i)}$ is associated with operation $f^{(i)}$ and is computed by evaluating the function
$$u^{(i)} = f(A^{(i)})$$
where $A^{(i)} = \text{Pa}(u^{(i)})$ is set of nodes that are parents of $u^{(i)}$
- Algorithm specifies a computational graph \mathcal{G}
 - Computation in reverse order gives back-propagation computational graph \mathcal{B}

Forward Propagation Algorithm

Algorithm 1: Performs computations mapping n_i inputs $u^{(1)}, \dots, u^{(n_i)}$ to an output $u^{(n)}$. This defines computational graph \mathcal{G} where each node computes numerical value $u^{(i)}$ by applying a function $f^{(i)}$ to a set of arguments $A^{(i)}$ that comprises values of previous nodes $u^{(j)}$, $j < i$ with $j \in \text{Pa}(u^{(i)})$. Input to \mathcal{G} is \mathbf{x} set to the first n_i nodes $u^{(1)}, \dots, u^{(n_i)}$. Output of \mathcal{G} is read off the last (output) node for $u^{(n)}$

- **for** $i = 1, \dots, n_i$ **do**

$$u^{(i)} \leftarrow x_i$$

- **end for**

- **for** $i = n_i + 1, \dots, n$ **do**

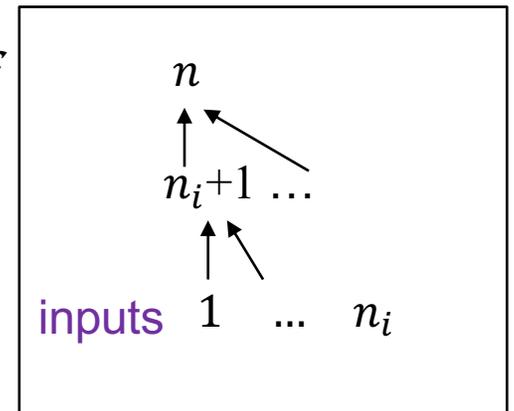
$$A^{(i)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(i)})\} \text{ parent nodes}$$

$$u^{(i)} \leftarrow f^{(i)}(A^{(i)}) \text{ activation based on parents}$$

- **end for**

- **return** $u^{(n)}$

Network \mathcal{G}



Computation in \mathcal{B}

- Proceeds exactly in reverse order of computation in \mathcal{G}
- Each node in \mathcal{B} computes the derivative $\frac{\partial u^{(n)}}{\partial u_i}$ associated with the forward graph node $u^{(i)}$
- This is done using the chain rule wrt the scalar output $u^{(n)}$

$$\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{i:j \in \text{IPa}(u^{(i)})} \frac{\partial u^{(n)}}{\partial u^{(i)}} \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

Preamble to Simplified Backprop

- Objective is to compute derivatives of $u^{(n)}$ with respect to variables in the graph
 - Here all variables are scalars and we wish to compute derivatives wrt
 - We wish to compute the derivatives wrt $u^{(1)}, \dots, u^{(ni)}$
- Algorithm computes the derivatives of all nodes in the graph

Simplified Backprop Algorithm

- Algorithm 2:** For computing derivatives of $u^{(n)}$ wrt variables in \mathcal{G}
 All variables are scalars and we wish to compute derivatives wrt $u^{(1)}, \dots, u^{(n_i)}$. We compute derivatives of all nodes in \mathcal{G}
- Run forward propagation to obtain network activations
- Initialize **grad-table**, a data structure that will store derivatives that have been computed, The entry **grad-table** $[u^{(i)}]$ will store the computed value of $\frac{\partial u^{(n)}}{\partial u_i}$
 - grad-table** $[u^{(n)}] \leftarrow 1$
 - for $j=n-1$ down to 1 do
 grad-table $[u^{(j)}] \leftarrow \sum_{i:j \in Pa(u^{(i)})} \text{grad-table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$
 - endfor
 - return $\{[u^{(i)}] \mid i=1, \dots, n_i\}$

Step 2 computes

$$\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{i:j \in Pa(u^{(i)})} \frac{\partial u^{(n)}}{\partial u^{(i)}} \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

$$\sum_{i:j \in Pa(u^{(i)})} \text{grad-table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

Computational Complexity

- Computational cost is proportional to no. of edges in graph (same as for forward prop)
 - Each $\frac{\partial u^{(i)}}{\partial u^{(j)}}$ is a function of the parents of $u^{(j)}$ and $u^{(i)}$ thus linking nodes of the forward graph to those added for \mathcal{B}
- Backpropagation thus avoids exponential explosion in repeated sub-expressions
 - By simplifications on the computational graph

Generalization to Tensors

- Backprop is designed to reduce the no. of common sub-expressions without regard to memory
- It performs on the order of one Jacobian product per node in the graph

5. Backprop in fully connected MLP

Backprop in fully connected MLP

- Consider specific graph associated with fully-connected multilayer perceptron
- Algorithm discussed next shows forward propagation
 - Maps parameters to supervised loss $L(\hat{y}, y)$ associated with a single training example (x, y) with \hat{y} the output when x is the input

Forward Prop: deep nn & cost computation

- **Algorithm 3:** The loss $L(\hat{y}, y)$ depends on output \hat{y} and on the target y . To obtain total cost J the loss may be added to a regularizer $\Omega(\theta)$ where θ contains all the parameters (weights and biases). *Algorithm 4* computes gradients of J wrt parameters W and b . This demonstration uses only single input example x
- **Require:** Net depth l ; Weight matrices $W^{(i)}, i \in \{1, \dots, l\}$; bias parameters $b^{(i)}, i \in \{1, \dots, l\}$; input x ; target output y
 1. $h^{(0)} = x$
 2. **for** $k = 1$ to l **do**
 - $a^{(k)} = b^{(k)} + W^{(k)}h^{(k-1)}$
 - $h^{(k)} = f(a^{(k)})$
 3. **end for**
 4. $\hat{y} = h^{(l)}$
 - $J = L(\hat{y}, y) + \lambda\Omega(\theta)$

Backward computation for deep NN of Algorithm 3

Algorithm 4: uses in addition to input \mathbf{x} a target \mathbf{y} . Yields gradients on activations $\mathbf{a}^{(k)}$ for each layer starting from output layer to first hidden layer. From these gradients one can obtain gradient on parameters of each layer, Gradients can be used as part of SGD

After forward computation compute gradient on output layer

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l-1, \dots, 1$ **do**

Convert gradient on layer's output into a gradient into the pre-nonlinearity activation (elementwise multiply if f is elementwise)

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights biases (incl. regularizn term)

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta), \quad \nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)T} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients wrt the next lower-level hidden layer's activations: $\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)T} \mathbf{g}$

end for