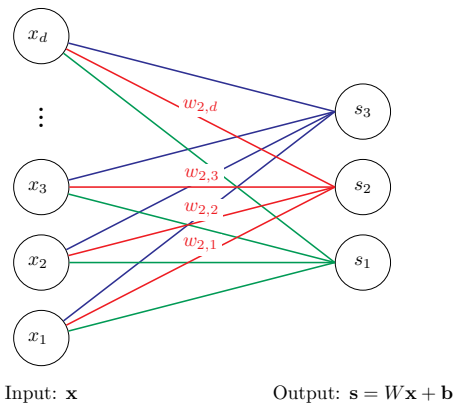


Gradient computations

Computation Graphs and the Back-propagation algorithm

Visualization of our multi-class classification function

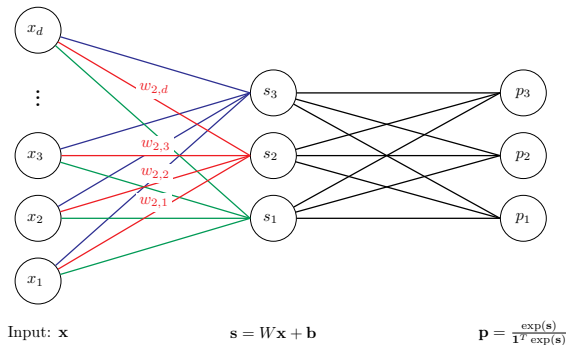
Linear with multiple outputs



Final decision: $g(\mathbf{x}) = \arg \max_j s_j$

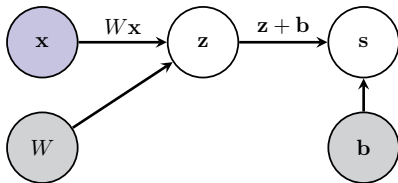
Visualization of our multi-class classification function

Linear with multiple probabilistic outputs



Final decision: $g(\mathbf{x}) = \arg \max_j p_j$

Computational graph of the multiple linear function



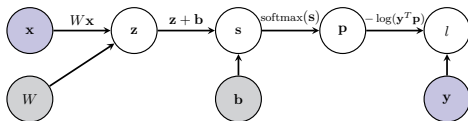
The computational graph:

- Represents order of computations.
- Displays the dependencies between the computed quantities.
- User input, parameters that have to be learnt.

Computational Graph helps automate gradient computations.

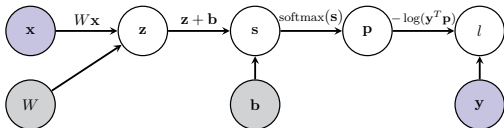
Computational graph of the complete loss function

- Linear scoring function + SOFTMAX + cross-entropy loss



where y is the 1-hot response vector induced by the label y .

To learn W , \mathbf{b} need to compute gradients



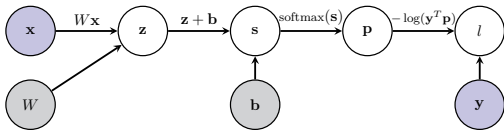
- Let l be the loss function defined by the computational graph.
- Find W , \mathbf{b} by optimizing

$$\arg \max_{W, \mathbf{b}} \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} l(\mathbf{x}, y, W, \mathbf{b})$$

- Solve using a variant of **mini-batch gradient descent**
 \implies need to efficiently compute the gradient vectors

$$\nabla_W l(\mathbf{x}, y, W, \mathbf{b})|_{(\mathbf{x}, y) \in \mathcal{D}} \quad \text{and} \quad \nabla_{\mathbf{b}} l(\mathbf{x}, y, W, \mathbf{b})|_{(\mathbf{x}, y) \in \mathcal{D}}$$

How do we compute these gradients?



- Let l be the complete loss function defined by the computational graph.
- How do we efficiently compute the gradient vectors

$$\nabla_W l(\mathbf{x}, y, W, \mathbf{b})|_{(\mathbf{x}, y) \in \mathcal{D}} \quad \text{and} \quad \nabla_{\mathbf{b}} l(\mathbf{x}, y, W, \mathbf{b})|_{(\mathbf{x}, y) \in \mathcal{D}}?$$

- Answer: **Back Propagation**
- BackProp relies on the **chain rule** applied to the **composition of functions**.
- Example: the composition of functions

$$l(\mathbf{x}, y, W, \mathbf{b}) = -\log(\mathbf{y}^T \text{softmax}(W\mathbf{x} + \mathbf{b}))$$

Chain Rule for functions with a scalar input and a scalar output

Differentiation of the composition of functions

- Have two functions $g : \mathbb{R} \rightarrow \mathbb{R}$ and $f : \mathbb{R} \rightarrow \mathbb{R}$.
- Define $h : \mathbb{R} \rightarrow \mathbb{R}$ as the composition of f and g :

$$h(x) = (f \circ g)(x) = f(g(x))$$

- How do we compute

$$\frac{dh(x)}{dx} ?$$

- Use the chain rule.

- Have functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$ and define $h : \mathbb{R} \rightarrow \mathbb{R}$ as

$$h(x) = (f \circ g)(x) = f(g(x))$$

- Derivative of h w.r.t. x is given by the Chain Rule.

- **Chain Rule**

$$\frac{dh(x)}{dx} = \frac{df(y)}{dy} \frac{dg(x)}{dx} \quad \text{where } y = g(x)$$

Example of the Chain Rule in action

- Have

$$g(x) = x^2, \quad f(x) = \sin(x)$$

- One composition of these two functions is

$$h(x) = f(g(x)) = \sin(x^2)$$

- According to the **chain rule**

$$\begin{aligned} \frac{dh(x)}{dx} &= \frac{df(y)}{dy} \frac{dg(x)}{dx} \quad \leftarrow \text{where } y = x^2 \\ &= \frac{d \sin(y)}{dy} \frac{dx^2}{dx} \\ &= \cos(y) 2x \\ &= 2x \cos(x^2) \quad \leftarrow \text{plug in } y = x^2 \end{aligned}$$

The composition of n functions

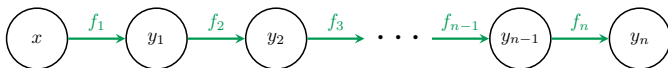
- Have functions $f_1, \dots, f_n : \mathbb{R} \rightarrow \mathbb{R}$
- Define function $h : \mathbb{R} \rightarrow \mathbb{R}$ as the composition of f_j 's

$$h(x) = (f_n \circ f_{n-1} \circ \dots \circ f_1)(x) = f_n(f_{n-1}(\dots(f_1(x))\dots))$$

- Can we compute the derivative

$$\frac{dh(x)}{dx} \quad ?$$

- Yes recursively apply the CHAIN RULE



- Define

$$g_j = f_n \circ f_{n-1} \circ \dots \circ f_j$$

- Therefore $g_1 = h$, $g_n = f_n$ and

$$g_j = g_{j+1} \circ f_j \quad \text{for } j = 1, \dots, n-1$$

- Let $y_j = f_j(y_{j-1})$ and $y_0 = x$ then

$$y_n = g_j(y_{j-1}) \quad \text{for } j = 1, \dots, n$$

- Apply the **Chain Rule**:

- For $j = 1, 2, 3, \dots, n-1$

$$\begin{aligned} \frac{dy_n}{dy_{j-1}} &= \frac{dg_j(y_{j-1})}{dy_{j-1}} = \frac{d(g_{j+1} \circ f_j)(y_{j-1})}{dy_{j-1}} = \frac{dg_{j+1}(y_j)}{dy_j} \frac{dg f_j(y_{j-1})}{dy_{j-1}} \\ &= \frac{dy_n}{dy_j} \frac{dy_j}{dy_{j-1}} \end{aligned}$$

The Chain Rule for the composition of n functions

Recursively applying this fact gives:

$$\begin{aligned}\frac{dh(x)}{dx} &= \frac{dg_1(x)}{dx} && \leftarrow \text{Apply } h = g_1 \\ &= \frac{d(g_2 \circ f_1)(x)}{dx} && \leftarrow \text{Apply } g_1 = g_2 \circ f_1 \\ &= \frac{dg_2(y_1)}{dy_1} \frac{df_1(x)}{dx} && \leftarrow \text{Apply chain rule \& } y_1 = f_1(x) \\ &= \frac{d(g_3 \circ f_2)(y_1)}{dy_1} \frac{df_1(x)}{dx} && \leftarrow \text{Apply } g_2 = g_3 \circ f_2 \\ &= \frac{dg_3(y_2)}{dy_2} \frac{df_2(y_1)}{dy_1} \frac{df_1(x)}{dx} && \leftarrow \text{Apply chain rule \& } y_2 = f_2(y_1) \\ &\quad \vdots \\ &= \frac{dg_n(y_{n-1})}{dy_{n-1}} \frac{df_{n-1}(y_{n-2})}{dy_{n-2}} \dots \frac{df_2(y_1)}{dy_1} \frac{df_1(x)}{dx} \\ &= \frac{df_n(y_{n-1})}{dy_{n-1}} \frac{df_{n-1}(y_{n-2})}{dy_{n-2}} \dots \frac{df_2(y_1)}{dy_1} \frac{df_1(x)}{dx} && \leftarrow \text{Apply } g_n = f_n\end{aligned}$$

where $y_j = (f_j \circ f_{j-1} \circ \dots \circ f_1)(x) = f_j(y_{j-1})$.

Exploit structure to compute gradient at a point

$$\frac{dh(x)}{dx} = \frac{dy_n}{dx} = \frac{dy_n}{dy_{n-1}} \frac{dy_{n-1}}{dy_{n-2}} \dots \frac{dy_2}{dy_1} \frac{dy_1}{dx}$$

Computation of $\frac{dh(x)}{dx}$ for $x = x^*$ relies on:

- **Record keeping:** Compute and record values of $y_j^* = f_j(y_{j-1}^*)$'s.
- **Iteratively aggregate** local gradients.

For $j = n, n-1, \dots, 2$

- Compute local derivative: $\left. \frac{df_j(y_{j-1})}{dy_{j-1}} \right|_{y_{j-1}=y_{j-1}^*}$

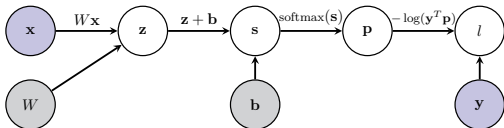
- Aggregate:

$$\left. \frac{dy_n^*}{dy_{j-1}} \right|_{y_{j-1}=y_{j-1}^*} = \left. \frac{dy_n^*}{dy_j} \right|_{y_j=y_j^*} \times \left. \frac{df_j(y_{j-1})}{dy_{j-1}} \right|_{y_{j-1}=y_{j-1}^*}$$

Remember $\frac{dy_n}{dy_{j+1}} = \frac{dy_n}{dy_{n-1}} \frac{dy_{n-1}}{dy_{n-2}} \dots \frac{dy_{j+2}}{dy_{j+1}}$

This is Backprop algorithm given a chain dependency between the y_j 's.

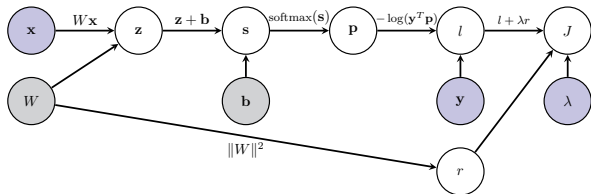
Problem 1: But what if I don't have a chain?



- This computational graph is **not a chain**.
- Some nodes have **multiple parents**.
- The function represented by graph is

$$l(\mathbf{x}, \mathbf{y}, W, \mathbf{b}) = -\log(\mathbf{y}^T \text{Softmax}(W\mathbf{x} + \mathbf{b}))$$

Problem 1a: And when a regularization term is added..

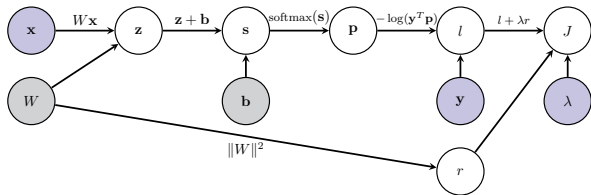


- This computational graph is **not a chain**.
- Some nodes have **multiple parents** and others **multiple children**.
- The function represented by graph is

$$J(\mathbf{x}, \mathbf{y}, W, \mathbf{b}, \lambda) = -\log(\mathbf{y}^T \text{Softmax}(W\mathbf{x} + \mathbf{b})) + \lambda \sum_{i,j} W_{i,j}^2$$

- How is the back-propagation algorithm defined in these cases?

Problem 1a: And when a regularization term is added..

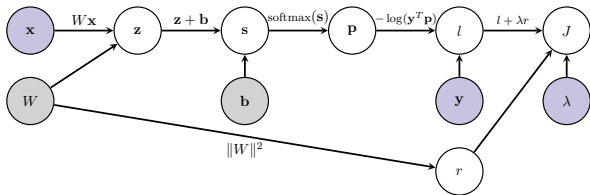


- This computational graph is **not a chain**.
- Some nodes have **multiple parents** and others **multiple children**.
- The function represented by graph is

$$J(\mathbf{x}, \mathbf{y}, W, \mathbf{b}, \lambda) = -\log(\mathbf{y}^T \text{Softmax}(W\mathbf{x} + \mathbf{b})) + \lambda \sum_{i,j} W_{i,j}^2$$

- How is the back-propagation algorithm defined in these cases?

Problem 2: Don't have scalar inputs and outputs



- The function represented by graph:

$$J(\mathbf{x}, \mathbf{y}, W, \mathbf{b}, \lambda) = -\log(\mathbf{y}^T \text{Softmax}(W\mathbf{x} + \mathbf{b})) + \lambda \sum_{i,j} W_{i,j}^2$$

- Nearly all of the inputs and intermediary outputs are **vectors** or **matrices**.
- How are the derivatives defined in this case?

Chain Rule for functions with vector inputs and vector outputs

Chain Rule for vector input and output

- Have two functions $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$ and $f : \mathbb{R}^m \rightarrow \mathbb{R}^c$.
- Define $h : \mathbb{R}^d \rightarrow \mathbb{R}^c$ as the composition of f and g :

$$h(\mathbf{x}) = (f \circ g)(\mathbf{x}) = f(g(\mathbf{x}))$$

- Consider

$$\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}}$$

- How is it defined and computed?
- What's the chain rule for vector valued functions?

Chain Rule for vector input and output

- Let $\mathbf{y} = h(\mathbf{x})$ where each $h : \mathbb{R}^d \rightarrow \mathbb{R}^c$ then

$$\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_d} \\ \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_2}{\partial x_d} \\ \vdots & \vdots & \vdots \\ \frac{\partial y_c}{\partial x_1} & \cdots & \frac{\partial y_c}{\partial x_d} \end{pmatrix} \quad \leftarrow \text{this is a Jacobian matrix}$$

and is a matrix of size $c \times d$.

- Chain Rule** says if $h = f \circ g$ ($g : \mathbb{R}^d \rightarrow \mathbb{R}^m$ and $f : \mathbb{R}^m \rightarrow \mathbb{R}^c$) then

$$\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$$

where $\mathbf{z} = g(\mathbf{x})$ and $\mathbf{y} = f(\mathbf{z})$.

- Both $\frac{\partial \mathbf{y}}{\partial \mathbf{z}}$ ($c \times m$) and $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ ($m \times d$) defined sllly to $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$.

Chain Rule for vector input and scalar output

The cost functions we will examine usually have a scalar output

- Let $\mathbf{x} \in \mathbb{R}^d$, $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^m \rightarrow \mathbb{R}$

$$\mathbf{z} = f(\mathbf{x})$$

$$s = g(\mathbf{z})$$

- The **Chain Rule** says gradient of output w.r.t. input

$$\frac{\partial s}{\partial \mathbf{x}} = \left(\frac{\partial s}{\partial x_1} \quad \cdots \quad \frac{\partial s}{\partial x_d} \right)$$

is given by a gradient times a Jacobian:

$$\frac{\partial s}{\partial \mathbf{x}} = \underbrace{\frac{\partial s}{\partial \mathbf{z}}}_{1 \times m} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{x}}}_{m \times d}$$

where

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \cdots & \frac{\partial z_1}{\partial x_d} \\ \frac{\partial z_2}{\partial x_1} & \cdots & \frac{\partial z_2}{\partial x_d} \\ \vdots & \vdots & \vdots \\ \frac{\partial z_m}{\partial x_1} & \cdots & \frac{\partial z_m}{\partial x_d} \end{pmatrix}$$

Two intermediary vector inputs and scalar output

- $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}^{m_1}$, $f_2 : \mathbb{R}^d \rightarrow \mathbb{R}^{m_2}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ ($n = m_1 + m_2$)

$$\mathbf{z}_1 = f_1(\mathbf{x}),$$

$$\mathbf{z}_2 = f_2(\mathbf{x})$$

$$s = g(\mathbf{z}_1, \mathbf{z}_2) = g(\mathbf{v}) \quad \text{where } \mathbf{v} = \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix}.$$

- **Chain Rule** says gradient of the output w.r.t. the input

$$\frac{\partial s}{\partial \mathbf{x}} = \left(\frac{\partial s}{\partial x_1} \quad \cdots \quad \frac{\partial s}{\partial x_d} \right)$$

is given by:

$$\frac{\partial s}{\partial \mathbf{x}} = \underbrace{\frac{\partial s}{\partial \mathbf{v}}}_{1 \times n} \underbrace{\frac{\partial \mathbf{v}}{\partial \mathbf{x}}}_{n \times d}$$

But

$$\frac{\partial s}{\partial \mathbf{v}} = \left(\frac{\partial s}{\partial \mathbf{z}_1} \quad \frac{\partial s}{\partial \mathbf{z}_2} \right) \quad \text{and} \quad \frac{\partial \mathbf{v}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}} \\ \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}} \end{pmatrix}$$

\implies

$$\frac{\partial s}{\partial \mathbf{x}} = \frac{\partial s}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial \mathbf{x}} = \underbrace{\frac{\partial s}{\partial \mathbf{z}_1}}_{1 \times m_1} \underbrace{\frac{\partial \mathbf{z}_1}{\partial \mathbf{x}}}_{m_1 \times d} + \underbrace{\frac{\partial s}{\partial \mathbf{z}_2}}_{1 \times m_2} \underbrace{\frac{\partial \mathbf{z}_2}{\partial \mathbf{x}}}_{m_2 \times d}$$

- $f_i : \mathbb{R}^d \rightarrow \mathbb{R}^{m_i}$ for $i = 1, \dots, t$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ ($n = m_1 + \dots + m_t$)

$$\mathbf{z}_i = f_i(\mathbf{x}), \quad \text{for } i = 1, \dots, t$$

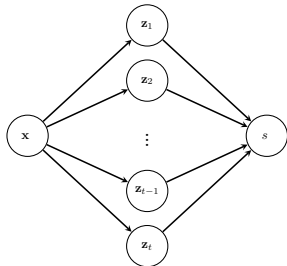
$$s = g(\mathbf{z}_1, \dots, \mathbf{z}_t)$$

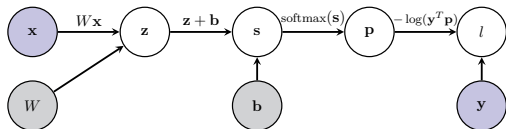
- Consequence of the **Chain Rule**

$$\frac{\partial s}{\partial \mathbf{x}} = \sum_{i=1}^t \frac{\partial s}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{x}}$$

- Computational graph interpretation. Let \mathcal{C}_x be the children nodes of x then

$$\frac{\partial s}{\partial \mathbf{x}} = \sum_{\mathbf{z} \in \mathcal{C}_x} \frac{\partial s}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$$



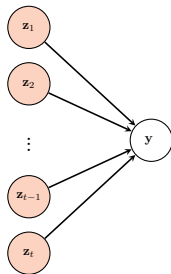


- Back-propagation when the computational graph is **not a chain**.
- Derivative computations when the inputs and outputs are not scalars. ✓
- Will now describe Back-prop for non-chains.

Back-propagation for non-chain computational graphs

- Have node y .
- Denote the set of y 's parent nodes by \mathcal{P}_y and their values by

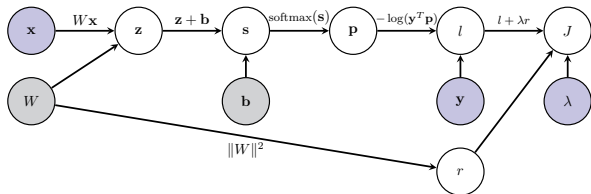
$$V_{\mathcal{P}_y} = \{\mathbf{z}.\text{value} \mid \mathbf{z} \in \mathcal{P}_y\}$$



- Given $V_{\mathcal{P}_y}$ can now apply the function $f_{\mathbf{z}}$

$$\mathbf{y}.\text{value} = f_{\mathbf{y}}(V_{\mathcal{P}_y})$$

Results that we need but already know



- Consider node W in the above graph. Its children are $\{z, r\}$. Applying the chain rule

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial r} \frac{\partial r}{\partial W} + \frac{\partial J}{\partial z} \frac{\partial z}{\partial W}$$

- In general for node c with children specified by \mathcal{C}_c :

$$\frac{\partial J}{\partial c} = \sum_{u \in \mathcal{C}_c} \frac{\partial J}{\partial u} \frac{\partial u}{\partial c}$$

Pseudo-Code for the Generic Forward Pass

procedure EVALUATEGRAPHFN(G)

$\mathcal{S} = \text{GetStartNodes}(G)$

for $s \in \mathcal{S}$ **do**

 ComputeBranch(s, G)

end for

end procedure

▷ G is the computational graph

▷ a start node has no parent and its value is already set

procedure COMPUTEBRANCH(s, G)

$\mathcal{C}_s = \text{GetChildren}(s, G)$

for each $n \in \mathcal{C}_s$ **do**

if ! n .computed **then**

$\mathcal{P}_n = \text{GetParents}(n, G)$

if CheckAllNodesComputed(\mathcal{P}_n) **then** ▷ Or not all parents of children are computed

$f_n = \text{GetNodeFn}(n)$

n .value = $f_n(\mathcal{P}_n)$

n .computed = true

 ComputeBranch(n, G)

end if

end if

end for

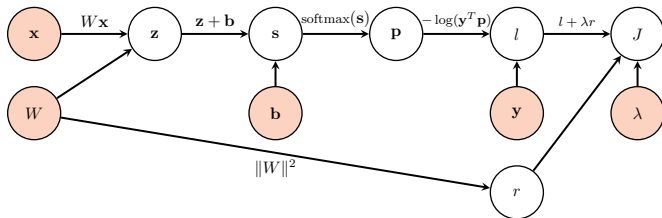
end procedure

▷ recursive fn evaluating nodes

▷ Try to evaluate each children node

▷ Unless child is already computed

Identify Start Nodes



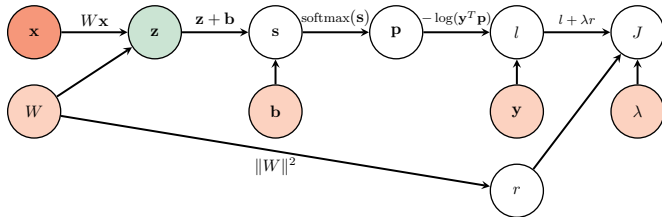
```

procedure EVALUATEGRAPHFN(G)  ▷ G is the computational graph
  S = GetStartNodes(G)
  for s ∈ S do
    ComputeBranch(s, G)
  end for
end procedure
  
```

```

procedure COMPUTEBRANCH(s, G)
  Cs = GetChildren(s, G)
  for each n ∈ Cs do
    if !n.computed then
      Pn = GetParents(n, G)
      if CheckAllNodesComputed(Pn) then
        fn = GetNodeFn(n)
        n.value = fn(Pn)
        n.computed = true
        ComputeBranch(n, G)
      end if
    end if
  end for
end procedure
  
```

Order in which nodes are evaluated



```

procedure EVALUATEGRAPHFN( $G$ )  $\triangleright G$  is the computational graph
   $S = \text{GetStartNodes}(G)$ 
  for  $s \in S$  do
     $\text{ComputeBranch}(s, G)$ 
  end for
end procedure

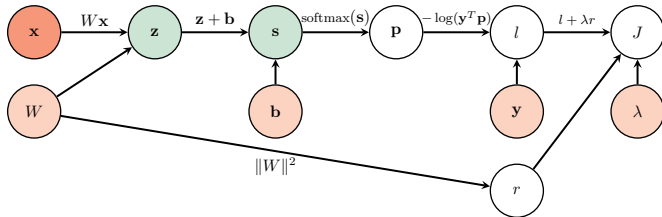
```

```

procedure COMPUTEBRANCH( $s, G$ )
   $C_s = \text{GetChildren}(s, G)$ 
  for each  $n \in C_s$  do
    if ! $n$ .computed then
       $P_n = \text{GetParents}(n, G)$ 
      if  $\text{CheckAllNodesComputed}(P_n)$  then
         $f_n = \text{GetNodeFn}(n)$ 
         $n$ .value =  $f_n(P_n)$ 
         $n$ .computed = true
         $\text{ComputeBranch}(n, G)$ 
      end if
    end if
  end for
end procedure

```


Order in which nodes are evaluated



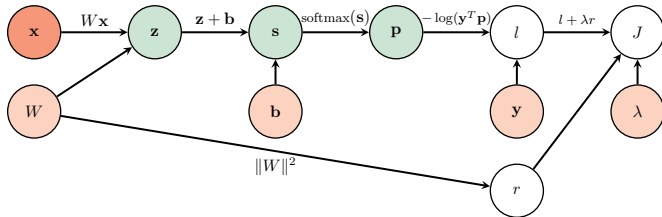
```

procedure EVALUATEGRAPHFN(G) ▷ G is the computational graph
  S = GetStartNodes(G)
  for s ∈ S do
    ComputeBranch(s, G)
  end for
end procedure
  
```

```

procedure COMPUTEBRANCH(s, G)
  Cs = GetChildren(s, G)
  for each n ∈ Cs do
    if !n.computed then
      Pn = GetParents(n, G)
      if CheckAllNodesComputed(Pn) then
        fn = GetNodeFn(n)
        n.value = fn(Pn)
        n.computed = true
        ComputeBranch(n, G)
      end if
    end if
  end for
end procedure
  
```

Order in which nodes are evaluated



```

procedure EVALUATEGRAPHFN( $G$ )  $\triangleright G$  is the computational graph
   $S = \text{GetStartNodes}(G)$ 
  for  $s \in S$  do
     $\text{ComputeBranch}(s, G)$ 
  end for
end procedure

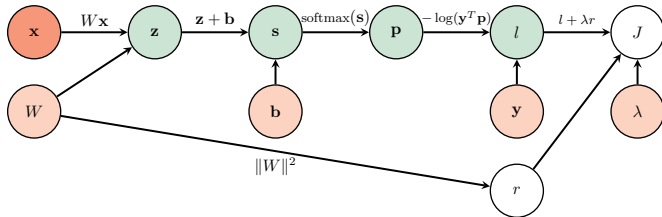
```

```

procedure COMPUTEBRANCH( $s, G$ )
   $C_s = \text{GetChildren}(s, G)$ 
  for each  $n \in C_s$  do
    if ! $n$ .computed then
       $P_n = \text{GetParents}(n, G)$ 
      if  $\text{CheckAllNodesComputed}(P_n)$  then
         $f_n = \text{GetNodeFn}(n)$ 
         $n$ .value =  $f_n(P_n)$ 
         $n$ .computed = true
         $\text{ComputeBranch}(n, G)$ 
      end if
    end if
  end for
end procedure

```

Order in which nodes are evaluated



```

procedure EVALUATEGRAPHFN(G) ▷ G is the computational graph
  S = GetStartNodes(G)
  for s ∈ S do
    ComputeBranch(s, G)
  end for
end procedure

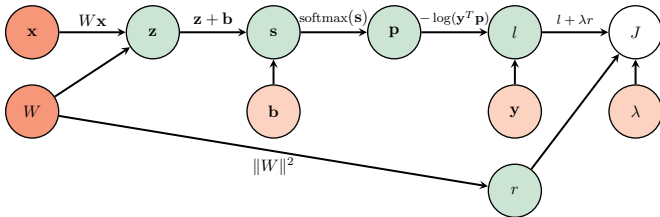
```

```

procedure COMPUTEBRANCH(s, G)
  Cs = GetChildren(s, G)
  for each n ∈ Cs do
    if !n.computed then
      Pn = GetParents(n, G)
      if CheckAllNodesComputed(Pn) then
        fn = GetNodeFn(n)
        n.value = fn(Pn)
        n.computed = true
        ComputeBranch(n, G)
      end if
    end if
  end for
end procedure

```

Order in which nodes are evaluated



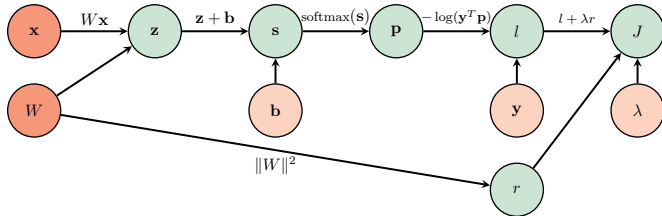
```

procedure EVALUATEGRAPHFN(G) ▷ G is the computational graph
  S = GetStartNodes(G)
  for s ∈ S do
    ComputeBranch(s, G)
  end for
end procedure
  
```

```

procedure COMPUTEBRANCH(s, G)
  Cs = GetChildren(s, G)
  for each n ∈ Cs do
    if !n.computed then
      Pn = GetParents(n, G)
      if CheckAllNodesComputed(Pn) then
        fn = GetNodeFn(n)
        n.value = fn(Pn)
        n.computed = true
        ComputeBranch(n, G)
      end if
    end if
  end for
end procedure
  
```

Order in which nodes are evaluated



```

procedure EVALUATEGRAPHFN(G) ▷ G is the computational graph
  S = GetStartNodes(G)
  for s ∈ S do
    ComputeBranch(s, G)
  end for
end procedure
  
```

```

procedure COMPUTEBRANCH(s, G)
  Cs = GetChildren(s, G)
  for each n ∈ Cs do
    if !n.computed then
      Pn = GetParents(n, G)
      if CheckAllNodesComputed(Pn) then
        fn = GetNodeFn(n)
        n.value = fn(Pn)
        n.computed = true
        ComputeBranch(n, G)
      end if
    end if
  end for
end procedure
  
```

Pseudo-Code for the Generic Backward Pass

procedure PERFORMBACKPASS(G)

$J = \text{GetResultNode}(G)$

 BackOp(J , G)

end procedure

▷ node with the value of cost function

▷ Start the Backward-pass

procedure BACKOP(s , G)

$C_s = \text{GetChildren}(s, G)$

if $C_s = \emptyset$ **then**

$s.\text{Grad} = 1$

end if

if AllGradientsComputed(C_s) **then**

$s.\text{Grad} = 0$

for each $c \in C_s$ **do**

$s.\text{Grad} += c.\text{Grad} * c.s.\text{Jacobian}$

end for

$s.\text{GradComputed} = \text{true}$

end if

for each $p \in \mathcal{P}_s$ **do**

$s.p.\text{Jacobian} = \frac{\partial f_p(\mathcal{P}_s)}{\partial p}$

 BackOp(p , G)

end for

end procedure

▷ At the result node

▷ Have computed all $\frac{\partial J}{\partial c}$ where $c \in C_s$

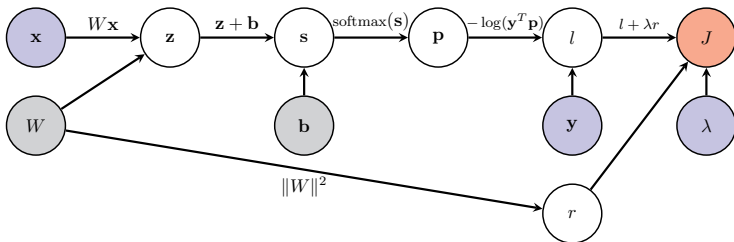
▷ $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$

▷ Compute the Jacobian of f_s w.r.t. each parent node

▷ $\frac{\partial f_s(\mathcal{P}_s)}{\partial p} = \frac{\partial s}{\partial p}$

Generic Backward Pass: Order of computations

Identify Result Node

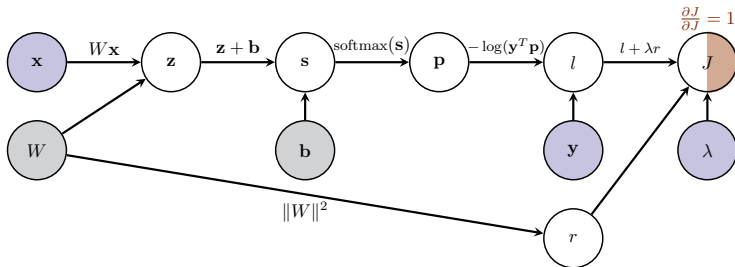


```
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
```

```
procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = ∅ then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each  $c \in C_s$  do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \sum \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each  $p \in \mathcal{P}_s$  do ▷ Compute Jacobian of  $f_s$  w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(\mathcal{P}_s)}{\partial p}$  ▷  $\frac{\partial f_s(\mathcal{P}_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
```

Generic Backward Pass: Order of computations

Compute Gradient of current node



```

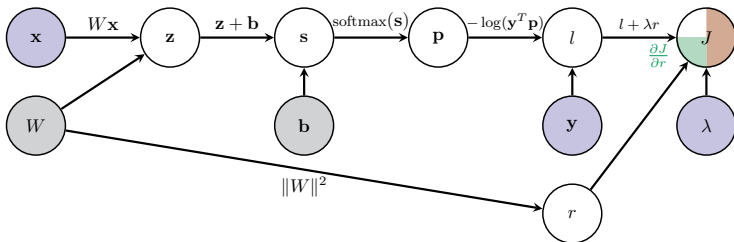
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G) ▷ At the result node
  if Cs = {} then
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```


Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

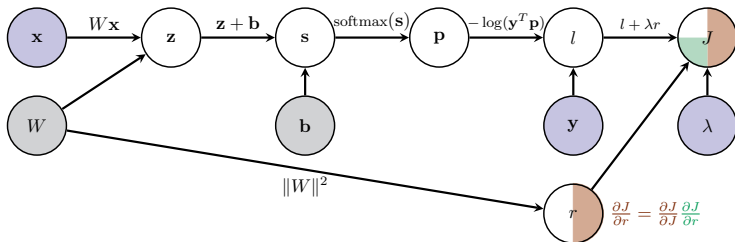
procedure PERFORMBACKPASS(G)
   $J = \text{GetResultNode}(G)$   $\triangleright$  node with the value of cost function
  BackOp( $J, G$ )  $\triangleright$  Start the Backward-pass
end procedure
    
```

```

procedure BACKOP( $s, G$ )
   $C_s = \text{GetChildren}(s, G)$ 
  if  $C_s = \emptyset$  then  $\triangleright$  At the result node
     $s.\text{Grad} = 1$ 
  else
    if AllGradientsComputed( $C_s$ ) then  $\triangleright$  All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
       $s.\text{Grad} = 0$ 
      for each  $c \in C_s$  do
         $s.\text{Grad} += c.\text{Grad} * c.s.\text{Jacobian}$   $\triangleright \frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
       $s.\text{GradComputed} = \text{true}$ 
    end if
  end if
  for each  $p \in \mathcal{P}_s$  do  $\triangleright$  Compute Jacobian of  $f_s$  w.r.t. each parent node
     $s.p.\text{Jacobian} = \frac{\partial f_s(\mathcal{P}_s)}{\partial p}$   $\triangleright \frac{\partial f_s(\mathcal{P}_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp( $p, G$ )
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Gradient of current node



```

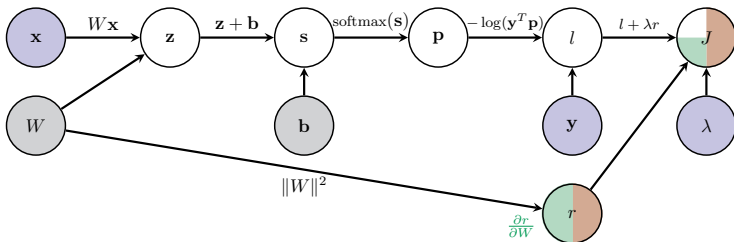
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = ∅ then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

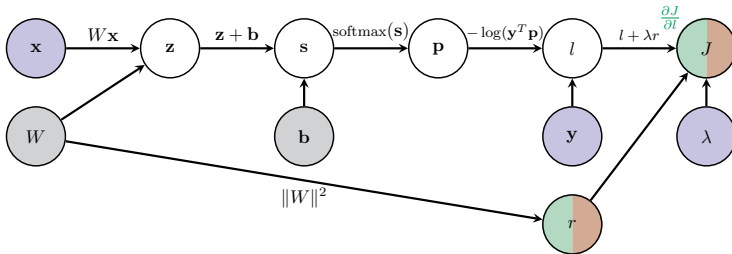
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = ∅ then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of  $f_s$  w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

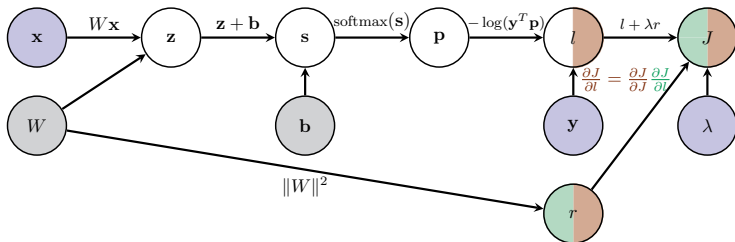
procedure PERFORMBACKPASS(G)
   $J = \text{GetResultNode}(G)$   $\triangleright$  node with the value of cost function
  BackOp( $J, G$ )  $\triangleright$  Start the Backward-pass
end procedure
  
```

```

procedure BACKOP( $s, G$ )
   $C_s = \text{GetChildren}(s, G)$ 
  if  $C_s = \emptyset$  then  $\triangleright$  At the result node
     $s.\text{Grad} = 1$ 
  else
    if AllGradientsComputed( $C_s$ ) then  $\triangleright$  All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
       $s.\text{Grad} = 0$ 
      for each  $c \in C_s$  do
         $s.\text{Grad} += c.\text{Grad} * c.s.\text{Jacobian}$   $\triangleright \frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
       $s.\text{GradComputed} = \text{true}$ 
    end if
  end if
  for each  $p \in \mathcal{P}_s$  do  $\triangleright$  Compute Jacobian of  $f_s$  w.r.t. each parent node
     $s.p.\text{Jacobian} = \frac{\partial f_s(\mathcal{P}_s)}{\partial p}$   $\triangleright \frac{\partial f_s(\mathcal{P}_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp( $p, G$ )
  end for
end procedure
  
```

Generic Backward Pass: Order of computations

Compute Gradient of current node



```

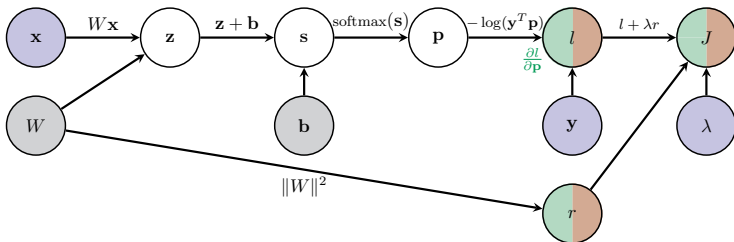
procedure PERFORMBACKPASS(G)
   $J = \text{GetResultNode}(G)$   $\triangleright$  node with the value of cost function
  BackOp( $J, G$ )  $\triangleright$  Start the Backward-pass
end procedure
  
```

```

procedure BACKOP( $s, G$ )
   $C_s = \text{GetChildren}(s, G)$   $\triangleright$  At the result node
  if  $C_s = \emptyset$  then
     $s.\text{Grad} = 1$ 
  else
    if AllGradientsComputed( $C_s$ ) then  $\triangleright$  All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
       $s.\text{Grad} = 0$ 
      for each  $c \in C_s$  do
         $s.\text{Grad} += c.\text{Grad} * c.s.\text{Jacobian}$   $\triangleright \frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
       $s.\text{GradComputed} = \text{true}$ 
    end if
  end if
  for each  $p \in \mathcal{P}_s$  do  $\triangleright$  Compute Jacobian of  $f_s$  w.r.t. each parent node
     $s.p.\text{Jacobian} = \frac{\partial f_s(\mathcal{P}_s)}{\partial p}$   $\triangleright \frac{\partial f_s(\mathcal{P}_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp( $p, G$ )
  end for
end procedure
  
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

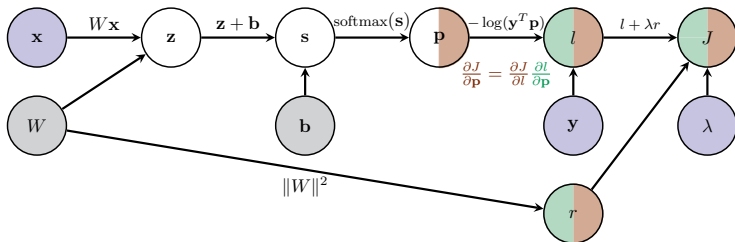
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = ∅ then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Gradient of current node



```

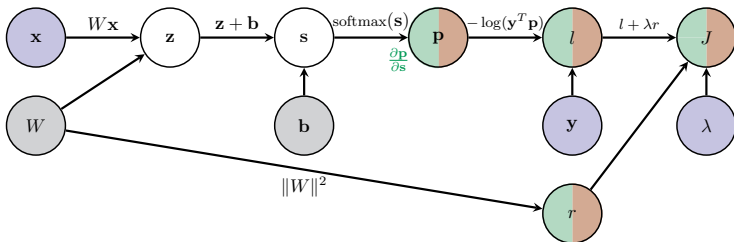
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = ∅ then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of  $f_s$  w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

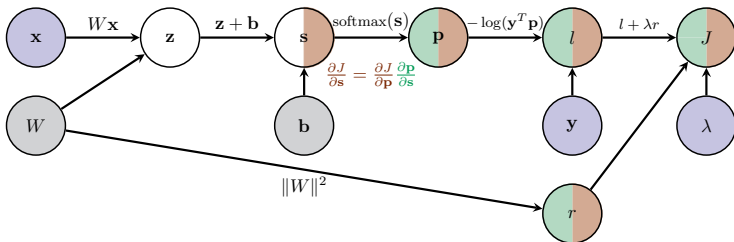
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = {} then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```


Generic Backward Pass: Order of computations

Compute Gradient of current node



```

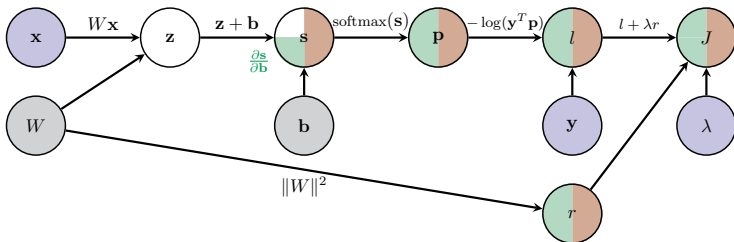
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = {} then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

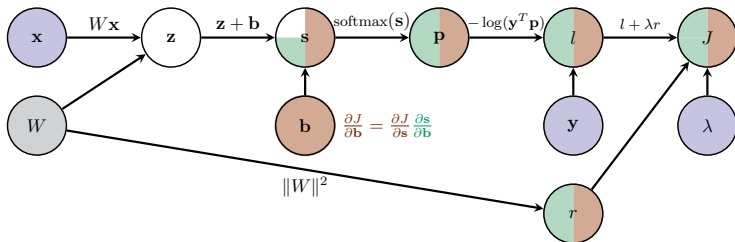
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = {} then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Gradient of current node



```

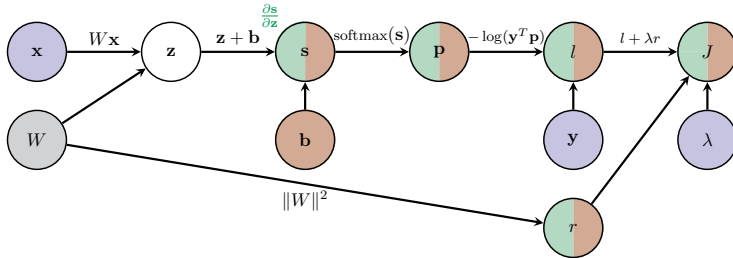
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G) ▷ At the result node
  if Cs = {} then
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of  $f_s$  w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

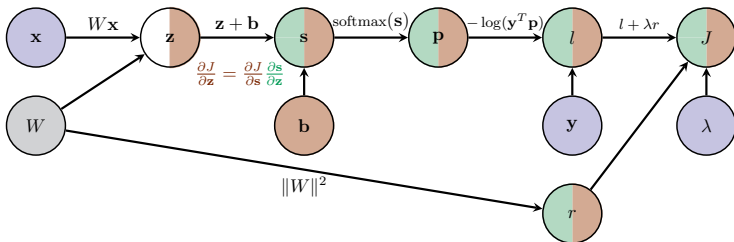
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G) ▷ At the result node
  if Cs = ∅ then
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Gradient of current node



```

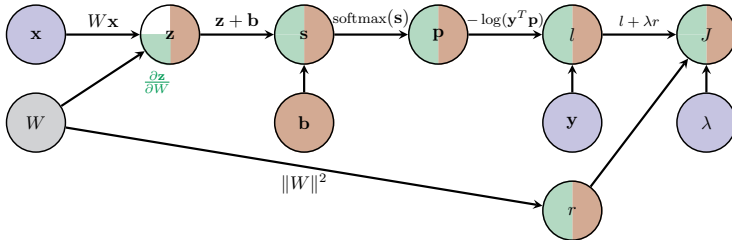
procedure PERFORMBACKPASS(G)
   $J = \text{GetResultNode}(G)$   $\triangleright$  node with the value of cost function
   $\text{BackOp}(J, G)$   $\triangleright$  Start the Backward-pass
end procedure
  
```

```

procedure BACKOP(s, G)
   $C_s = \text{GetChildren}(s, G)$ 
  if  $C_s = \emptyset$  then  $\triangleright$  At the result node
     $s.\text{Grad} = 1$ 
  else
    if  $\text{AllGradientsComputed}(C_s)$  then  $\triangleright$  All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
       $s.\text{Grad} = 0$ 
      for each  $c \in C_s$  do
         $s.\text{Grad} += c.\text{Grad} * c.s.\text{Jacobian}$   $\triangleright \frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
       $s.\text{GradComputed} = \text{true}$ 
    end if
  end if
  for each  $p \in \mathcal{P}_s$  do  $\triangleright$  Compute Jacobian of  $f_s$  w.r.t. each parent node
     $s.p.\text{Jacobian} = \frac{\partial f_s(\mathcal{P}_s)}{\partial p}$   $\triangleright \frac{\partial f_s(\mathcal{P}_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
     $\text{BackOp}(p, G)$ 
  end for
end procedure
  
```

Generic Backward Pass: Order of computations

Compute Jacobian of current node w.r.t. one parent



```

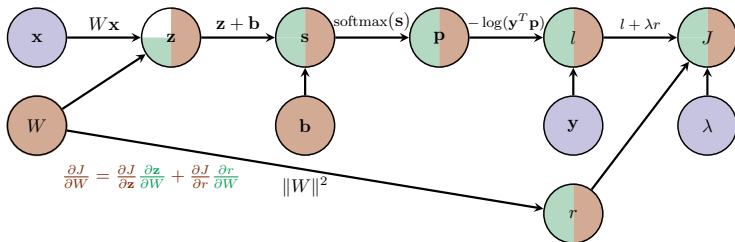
procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = {} then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```

Generic Backward Pass: Order of computations

Compute Gradient of current node

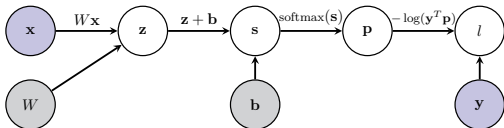


```

procedure PERFORMBACKPASS(G)
  J = GetResultNode(G) ▷ node with the value of cost function
  BackOp(J, G) ▷ Start the Backward-pass
end procedure
    
```

```

procedure BACKOP(s, G)
  Cs = GetChildren(s, G)
  if Cs = {} then ▷ At the result node
    s.Grad = 1
  else
    if AllGradientsComputed(Cs) then ▷ All  $\frac{\partial J}{\partial c}$  computed where  $c \in C_s$ 
      s.Grad = 0
      for each c ∈ Cs do
        s.Grad += c.Grad * c.s.Jacobian ▷  $\frac{\partial J}{\partial s} += \frac{\partial J}{\partial c} \frac{\partial c}{\partial s}$ 
      end for
      s.GradComputed = true
    end if
  end if
  for each p ∈ Ps do ▷ Compute Jacobian of fs w.r.t. each parent node
    s.p.Jacobian =  $\frac{\partial f_s(P_s)}{\partial p}$  ▷  $\frac{\partial f_s(P_s)}{\partial p} = \frac{\partial s}{\partial p}$ 
    BackOp(p, G)
  end for
end procedure
    
```



- Back-propagation when the computational graph is **not a chain**. ✓
- Derivative computations when the inputs and outputs are not scalars. ✓