

MATH 565: Lecture 26 (04/16/2026)

Today: * Loss function to \bar{w} derivatives
* CGs with vector variables

Loss Function to Weight Derivatives

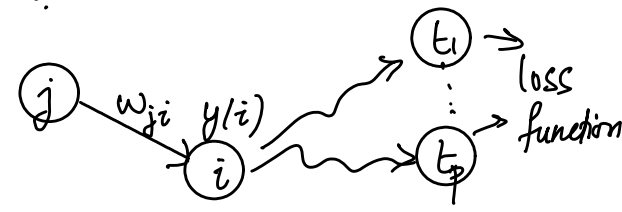
Recall: $S(s, t) = \frac{\partial y(t)}{\partial y(s)}$, node-to-node derivative and the

DP to compute $S(s, t)$ using back propagation.

The output node values ($y(t)$ or $y(t_i)$ for $i=1, \dots, p$) are used to compute loss function J . And we want to compute $\nabla_{\bar{w}} J$ (or $\frac{\partial J}{\partial \bar{w}}$) for edge weights $w_{ji}, (j, i) \in E$.

We apply the chain rule again!

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y(i)} \cdot \frac{\partial y(i)}{\partial w_{ji}}$$



If there are p output nodes t_1, \dots, t_p , then

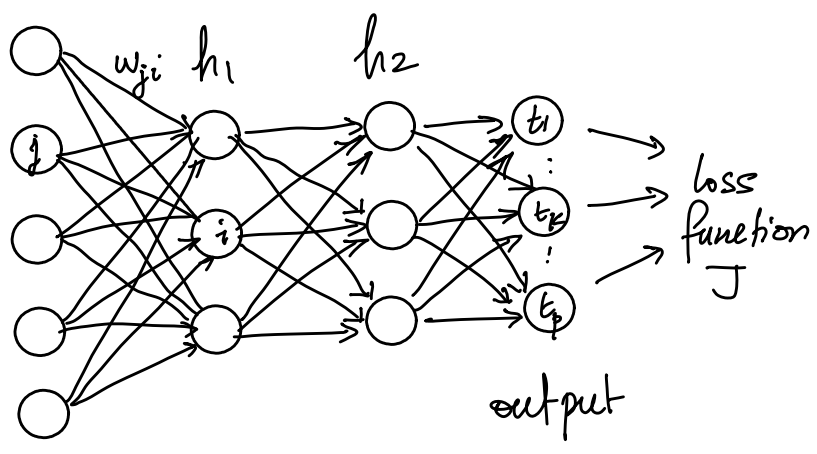
$$\frac{\partial J}{\partial w_{ji}} = \sum_{k=1}^p \left[\frac{\partial J}{\partial y(t_k)} \cdot \frac{\partial y(t_k)}{\partial y(i)} \right] \cdot \frac{\partial y(i)}{\partial w_{ji}}$$

We can modify the previous DP algorithm to compute $\frac{\partial J}{\partial w_{ji}}$ instead. Instead of $S(j, t)$, we use

$$\Delta(i) = \frac{\partial J}{\partial y(i)} \text{ and we back propagate the loss}$$

derivatives all the way to the input nodes.

We modify the original DP algorithm to the present context.



$$\Delta(t_k) = \frac{\partial J}{\partial y(t_k)}, \quad k=1, \dots, p \quad (\text{initialization})$$

do
 select unprocessed $i \in \mathcal{N}$ for which $\Delta(j)$ is known $\forall j \in A(i)$;

$$\Delta(i) \leftarrow \sum_{j \in A(i)} \Delta(j) \cdot z(i,j)$$

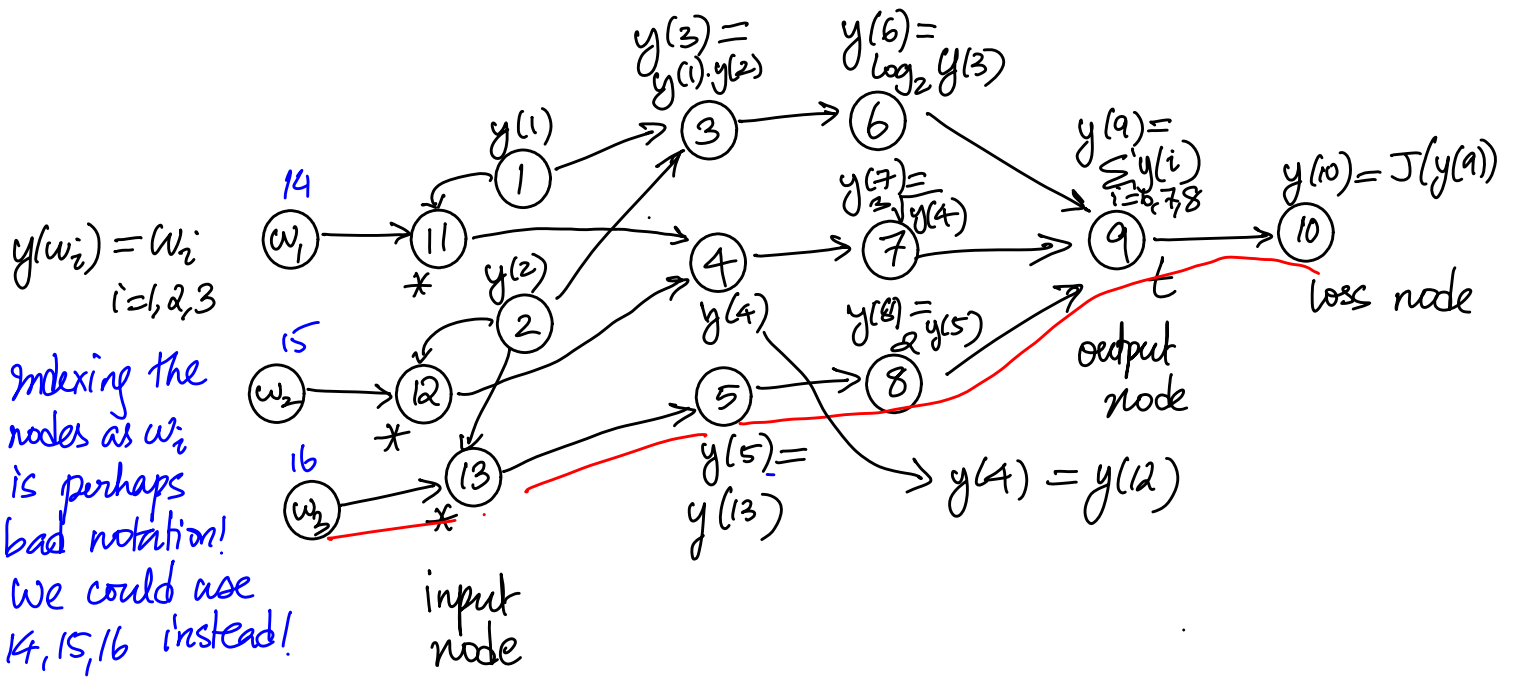
while $\exists i \in \mathcal{N}$ unprocessed
 for $(j,i) \in \mathcal{E}$ with edge weight w_{ji}

$$\frac{\partial J}{\partial w_{ji}} \leftarrow \Delta(i) \cdot \frac{\partial y(i)}{\partial w_{ji}}$$

end

Another Approach: Stick with node-to-node derivatives, and add extra nodes to capture edge weights and the loss function. This is a standard network transformation; the other transformation of replacing a node with an edge is more standard in network flows.

Illustration on Example



$$y(11) = w_1 \cdot y(1) \quad y(12) = w_2 \cdot y(2) \quad y(13) = w_3 \cdot y(2)$$

After the network is transformed in this manner, we can run the node-to-node back propagation DP.

Now, we can run standard gradient descent (GD):

$$\bar{w} \leftarrow \bar{w} - \alpha \nabla_{\bar{w}} J$$

α learning rate

Illustration

Let the loss function be $y(10) = J(y(9)) = \ln(y(9))$.

$$\Rightarrow \frac{\partial J}{\partial w_3} = \frac{\partial J}{\partial y(w_3)} = \frac{\partial y(10)}{\partial y(9)} \cdot \underbrace{S(5,9)}_{2^{y(5)} \ln 2} \cdot \frac{\partial y(5)}{\partial y(13)} \cdot \frac{\partial y(13)}{\partial y(w_3)}$$

$\frac{\partial y(10)}{\partial y(9)} = \frac{1}{y(9)}$
 $\frac{\partial y(5)}{\partial y(13)} = 1$
 $\frac{\partial y(13)}{\partial y(w_3)} = y(2)$

$$= \frac{1}{y(9)} \cdot 2^{y(5)} \ln 2 \cdot y(2), \text{ which can be computed using numerical differentiation.}$$

Computational Graphs with Vector Variables

We had seen $z(i,j) = \frac{\partial y(j)}{\partial y(i)}$ and $S(s,t) = \frac{\partial y(t)}{\partial y(s)}$. In more general CGs, $y(i) \in \mathbb{R}^d$ for some $d > 1$. The derivative is now replaced by the Jacobian.

Def For $\bar{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\bar{f}(\bar{x}) = \begin{bmatrix} f_1(\bar{x}) \\ \vdots \\ f_m(\bar{x}) \end{bmatrix}$ $\bar{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$

the Jacobian is

or \mathcal{D}_f or $\mathcal{D}\bar{f}$

$$\mathcal{D}(\bar{f}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}_{m \times n} = \begin{bmatrix} (\nabla_{\bar{x}} f_1)^T \\ \vdots \\ (\nabla_{\bar{x}} f_m)^T \end{bmatrix}$$

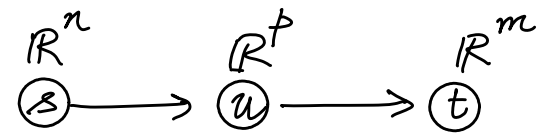
↓ gradients of f_i written as rows.

if we have $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$, then

$$\mathcal{D}(g \circ f) = \mathcal{D}(g) \cdot \mathcal{D}(f)$$

$p \times n$ $p \times m$ $m \times n$

Let $y(t) \in \mathbb{R}^m$, $y(s) \in \mathbb{R}^n$. Then $\mathcal{D}(s,t)$ is an $m \times n$ matrix of partial derivatives.



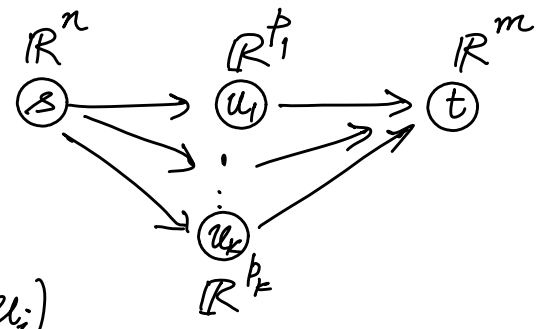
With a simple s-t path, we get

$$\mathcal{D}(s,t) = \mathcal{D}(u,t) \cdot \mathcal{D}(s,u)$$

$m \times n$ $m \times p$ $p \times n$

*! The order is important here - it should honor matrix multiplication!

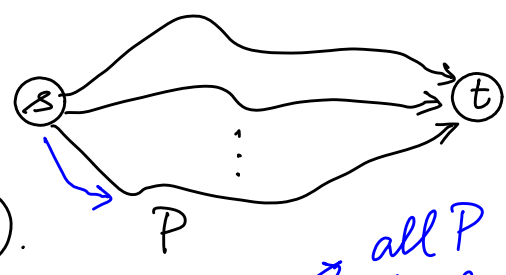
If there are k s - t paths in this fashion, we get



$$D(s, t)_{m \times n} = \sum_{i=1}^k D(u_i, t)_{m \times p_i} \cdot D(s, u_i)_{p_i \times n}$$

Here is the result in general.

Let $\mathcal{P} = \{P_i\}$ be the collection of s - t paths. Each path P_i has k_i edges ($k_i \geq 1$).



all P need not have the same k !

For brevity of notation, consider path P with k edges of the form

$$P: s = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k = t. \quad \text{Then}$$

$$D(s, t) = \sum_{P \in \mathcal{P}} D(v_{k-1}, t) \cdot D(v_{k-2}, v_{k-1}) \cdot \dots \cdot D(v_1, v_2) \cdot D(s, v_1)$$

Note the order of multiplication: you start with the term closest to s on the rightmost end of the product.

We can generalize the DP algorithms for node-to-node and loss-to-weight derivative calculations to the vector setting.