

Chain Rule Review

Scalar valued functions of scalars

$$y(x) = f(g(h(x)))$$

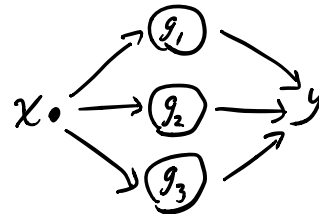
$$y'(x) = f'(g(h(x))) g'(h(x)) h'(x)$$

Scalar output

$$y(x) = f(g_1(x), g_2(x), g_3(x))$$

$$y'(x) = \frac{\partial f}{\partial g_1} \frac{dg_1}{dx} + \frac{\partial f}{\partial g_2} \frac{dg_2}{dx} + \frac{\partial f}{\partial g_3} \frac{dg_3}{dx}$$

Influence Diagram



Scalar Output Multivariate input

$$f(x_1, x_2, x_3)$$

$$\nabla f = (\partial_{x_1} f, \partial_{x_2} f, \partial_{x_3} f)$$

https://www.youtube.com/watch?v=wG_nF1awSSY

Finite differences

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$
$$\frac{\partial f}{\partial x_i} \approx \frac{f(x + he_i) - f(x)}{h}$$

Requires $O(n)$ evaluations:
 $e_1, e_2, e_3, \dots, e_n$

2:37 / 14:24

What is Automatic Differentiation?

What are the benefits/challenges with computing gradients by finite differences?

Benefit

Only need
black box
access to f

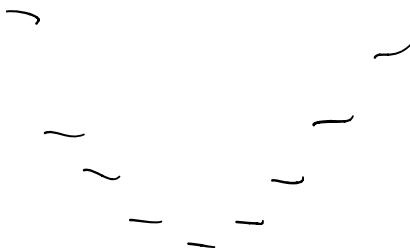
Challenge

Choose value of h
Numerical instability if h small
inaccurate if h large

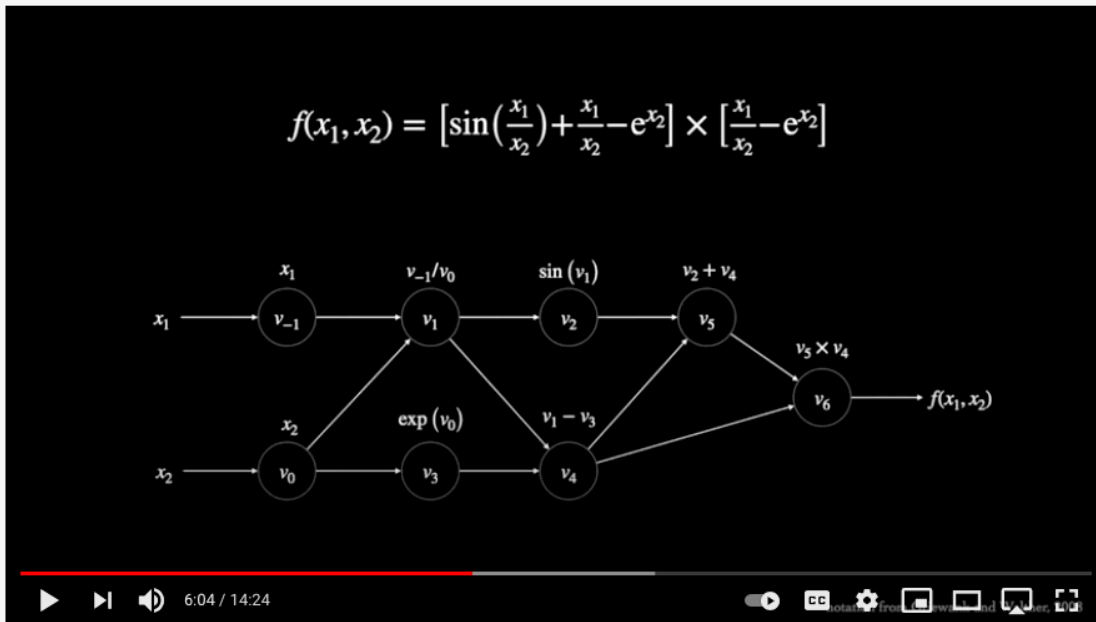
To find $\nabla_x f$, compute

$\frac{\partial f}{\partial x_i}$ for $i=1 \dots n$

$n+1$ forward passes of f



Computational Graphs + Forward Mode Auto. Diff.



What is Automatic Differentiation?

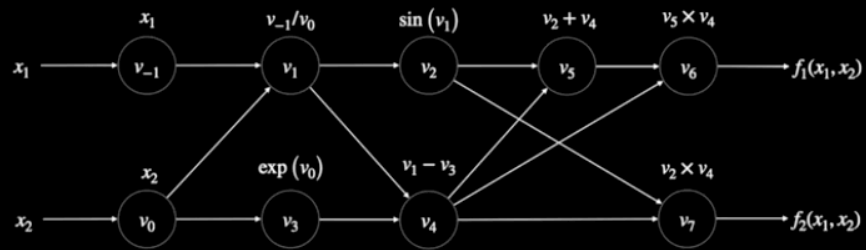
$$\frac{\partial f}{\partial x}$$

set based on which partial is computed

Primals		Tangents	
$v_{-1} = x_1$	= 1.500	\dot{v}_{-1}	= 1.000
$v_0 = x_2$	= 0.500	\dot{v}_0	= 0.000
$v_1 = v_{-1}/v_0$	= 3.000	$\dot{v}_1 = (v_0\dot{v}_{-1} - v_{-1}\dot{v}_0)/v_0^2$	= 2.000
$v_2 = \sin(v_1)$	= 0.141	$\dot{v}_2 = \cos(v_1) \times \dot{v}_1$	= -1.980
$v_3 = \exp(v_0)$	= 1.649	$\dot{v}_3 = v_3 \times \dot{v}_0$	= 0.000
$v_4 = v_1 - v_3$	= 1.351	$\dot{v}_4 = \dot{v}_1 - \dot{v}_3$	= 2.000
$v_5 = v_2 + v_4$	= 1.492	$\dot{v}_5 = \dot{v}_2 + \dot{v}_4$	= 0.020
$v_6 = v_5 \times v_4$	= 2.017	$\dot{v}_6 = \dot{v}_5 v_4 - \dot{v}_4 v_5$	= 3.012
$f(x_1, x_2) = v_6$	= 2.017	$\frac{\partial f}{\partial x_1} = \dot{v}_6$	= 3.012

6:34 / 14:24

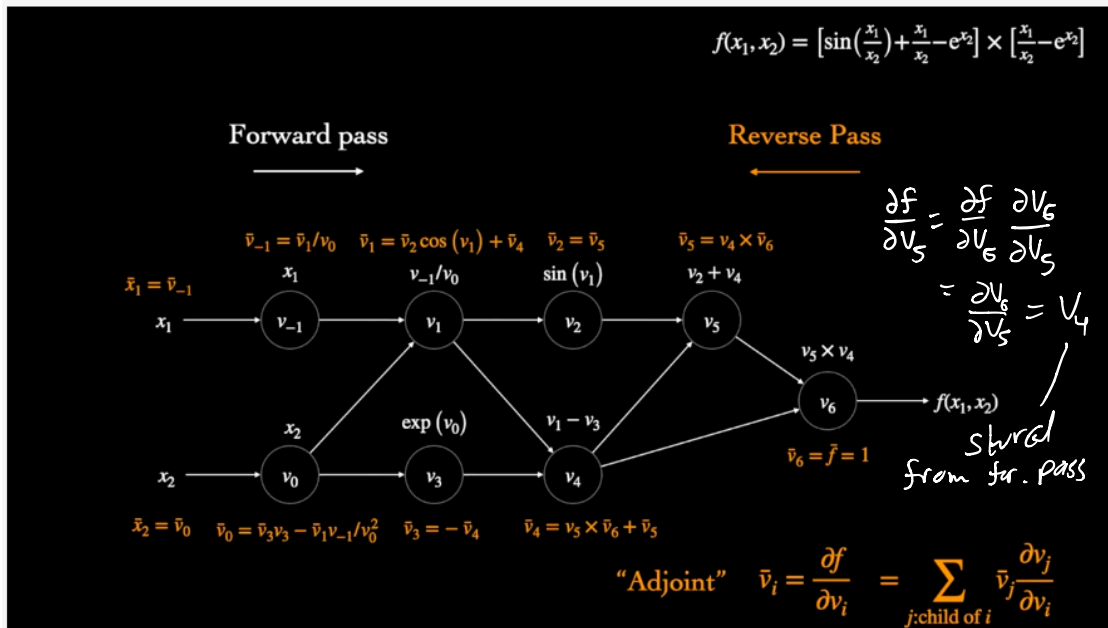
What is Automatic Differentiation?



compute $\frac{\partial f_1}{\partial x_1}$ and $\frac{\partial f_2}{\partial x_1}$ in a single forward pass

What is Automatic Differentiation?

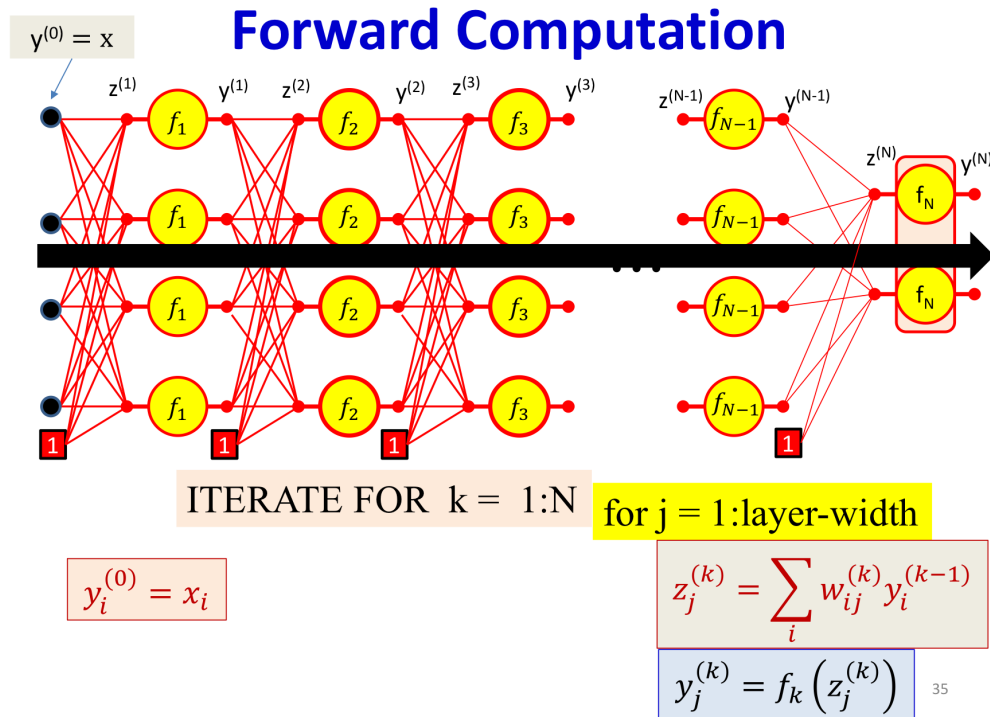
Benefits of Forward Mode Auto. Diff.



What is Automatic Differentiation?

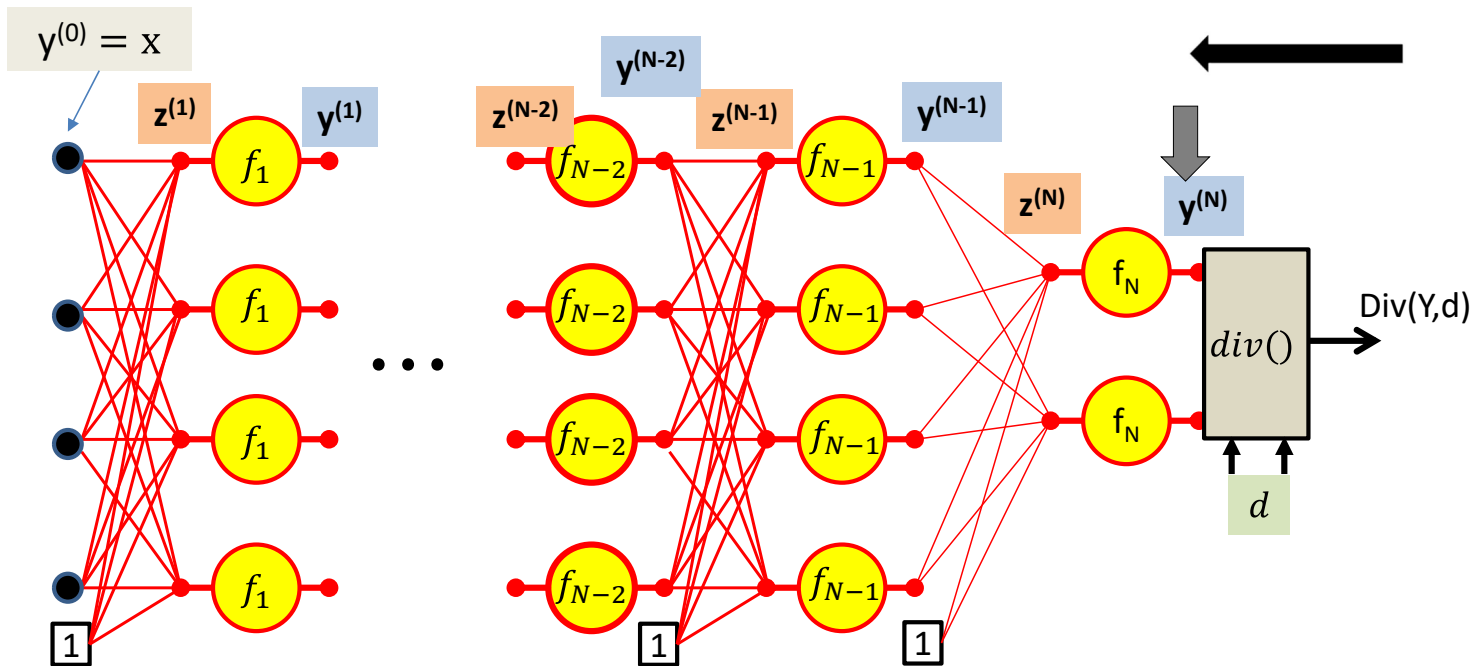
In forward pass, calculate activation of all neurons. We store those activations. In the backward pass, we will use these values to compute partial derivatives.

Reverse Mode Auto. Diff.



<http://deeplearning.cs.cmu.edu/S21/index.html>

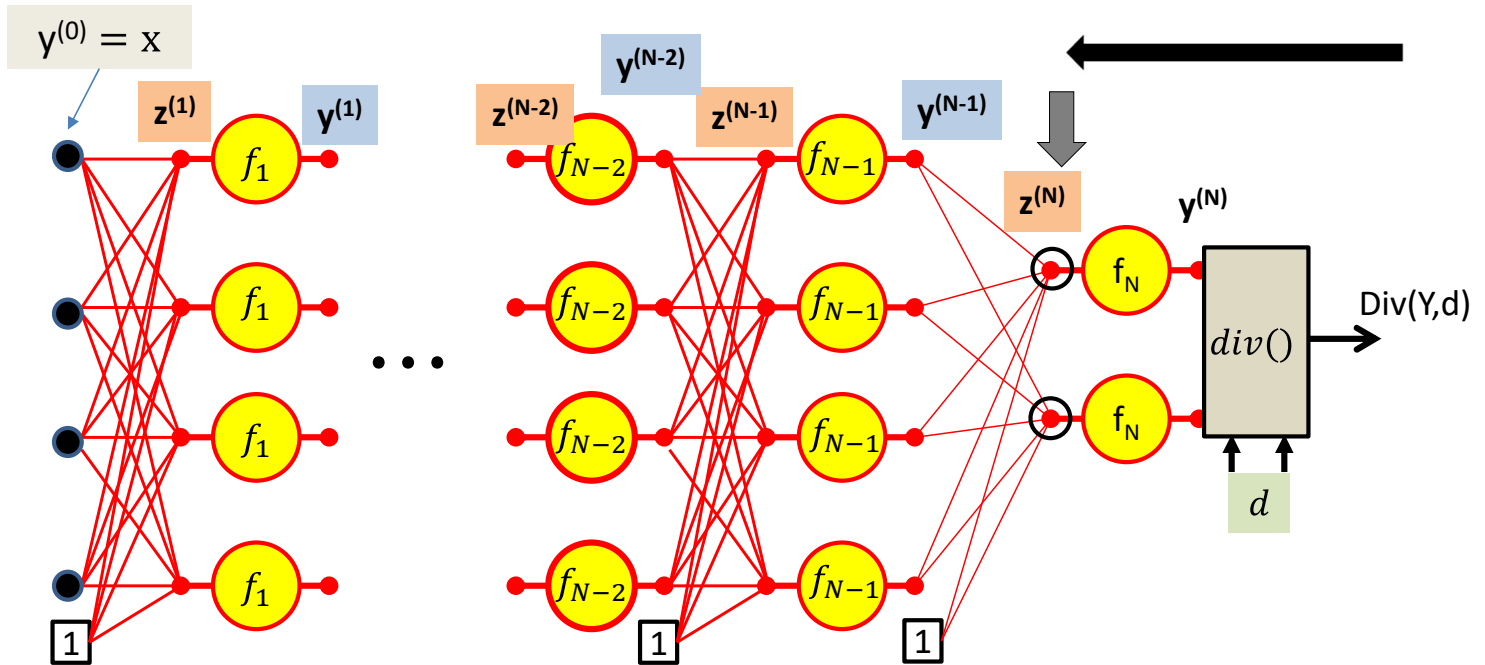
Computing derivatives



The derivative w.r.t the actual output of the final layer of the network is simply the derivative w.r.t to the output of the network

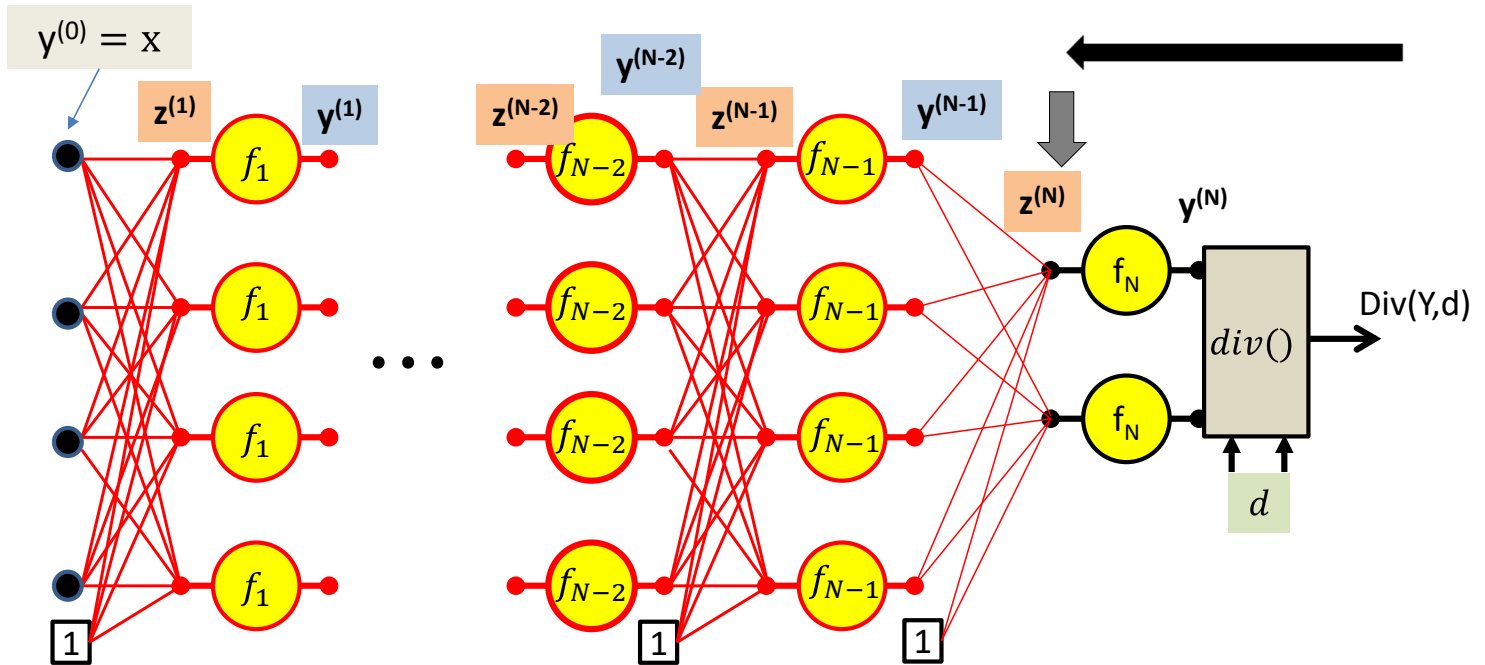
$$\frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}} = \frac{\partial \text{Div}(Y, d)}{\partial y_i}$$

Computing derivatives



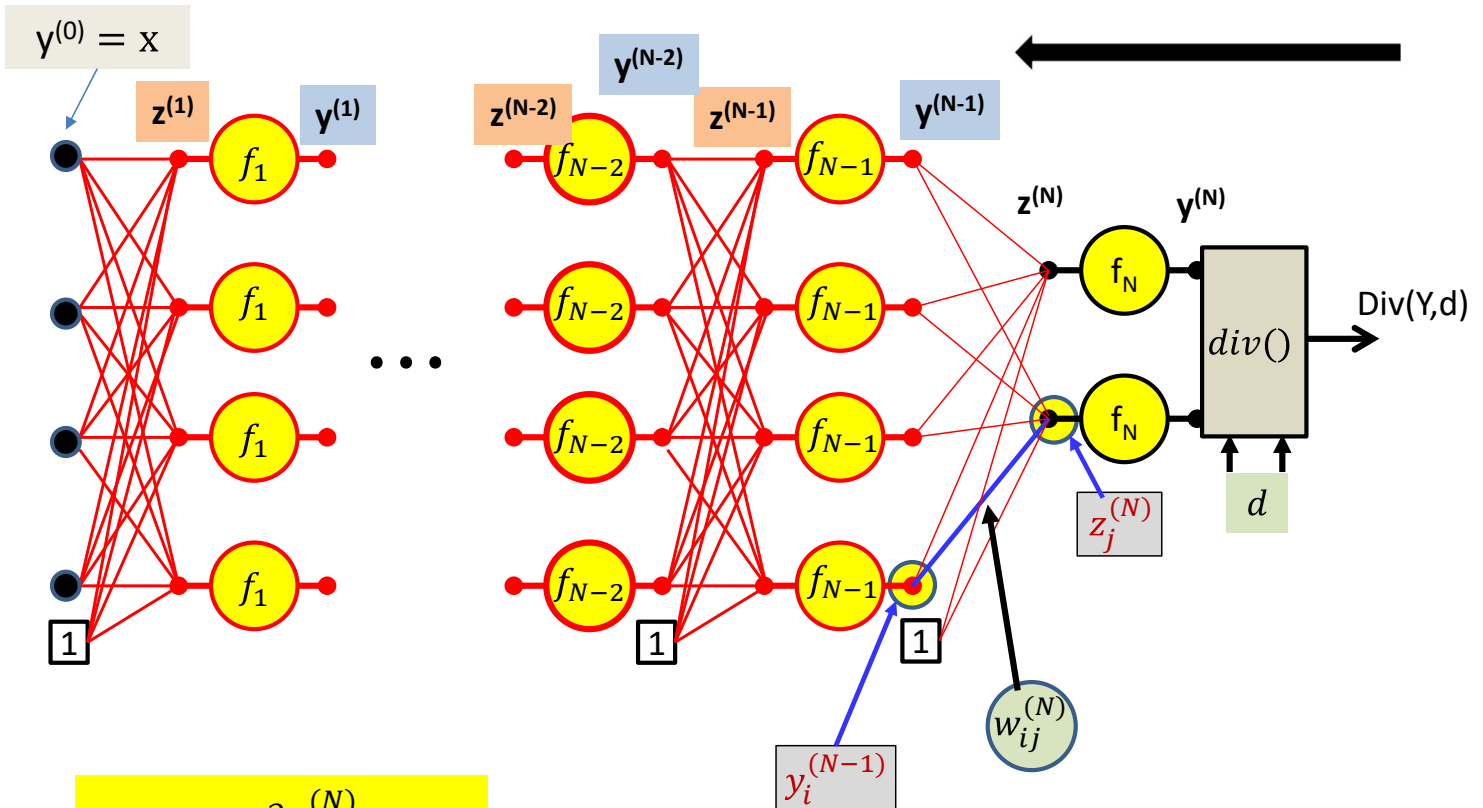
$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

Computing derivatives



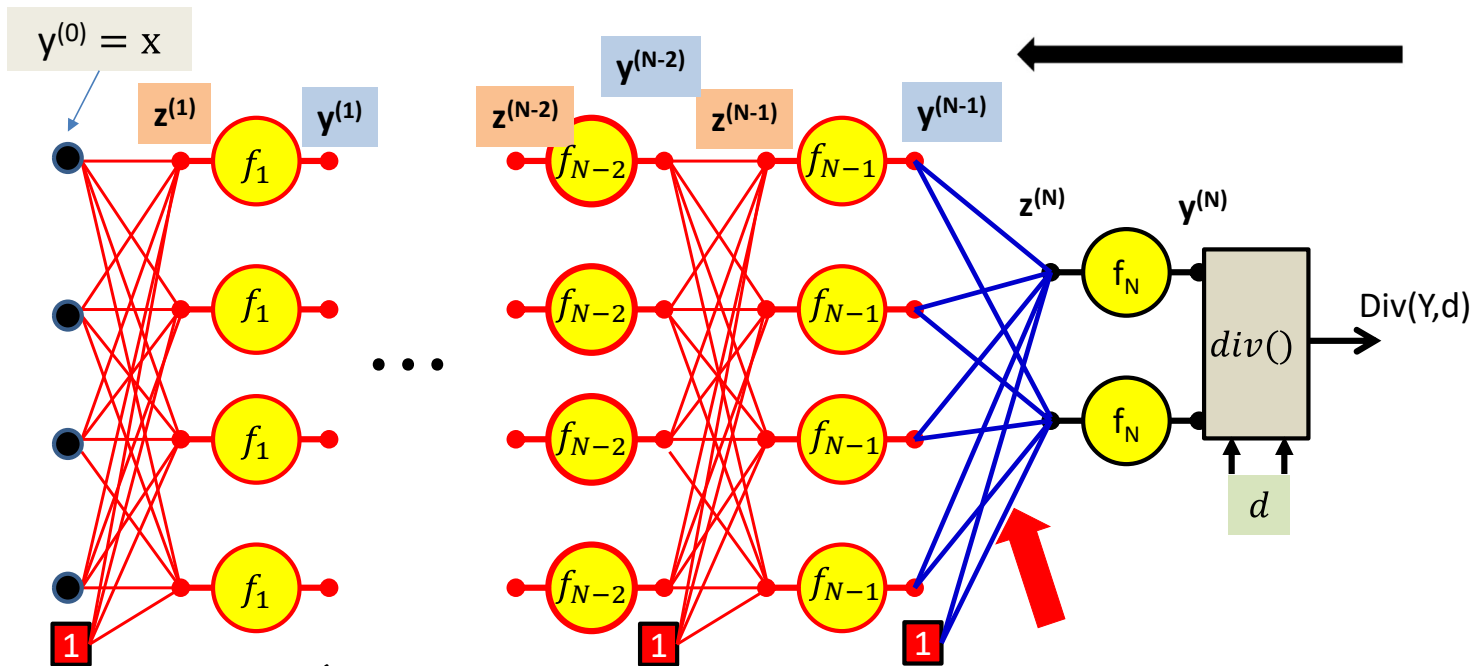
$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

Computing derivatives



$$\frac{\partial Div}{\partial w_{ij}^{(N)}} = \frac{\partial z_j^{(N)}}{\partial w_{ij}^{(N)}} \frac{\partial Div}{\partial z_j^{(N)}}$$

Computing derivatives



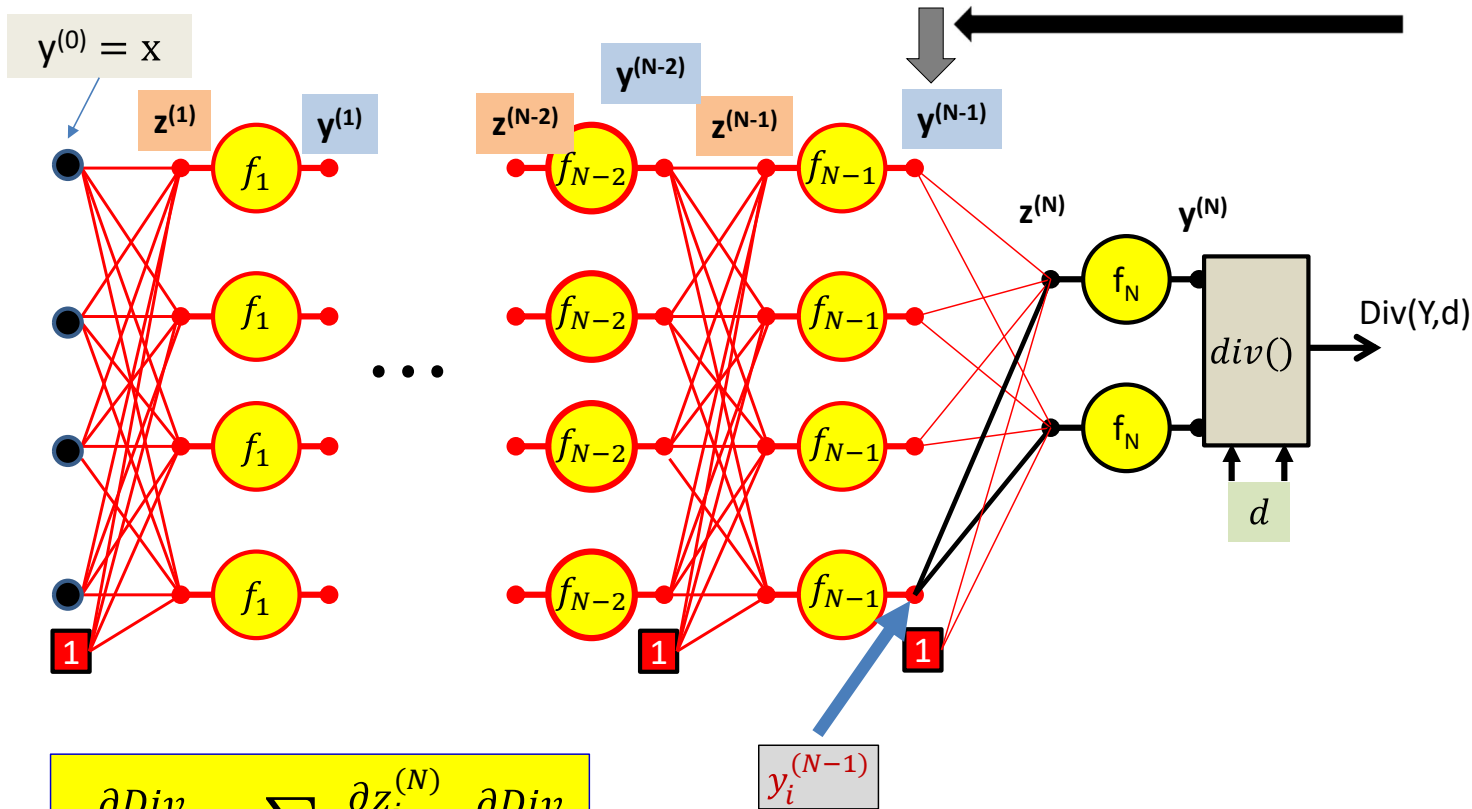
actis. value

$$\frac{\partial}{\partial w}(w \cdot y_i)$$

$$\frac{\partial \text{Div}}{\partial w_{ij}^{(N)}} = y_i^{(N-1)} \frac{\partial \text{Div}}{\partial z_j^{(N)}}$$

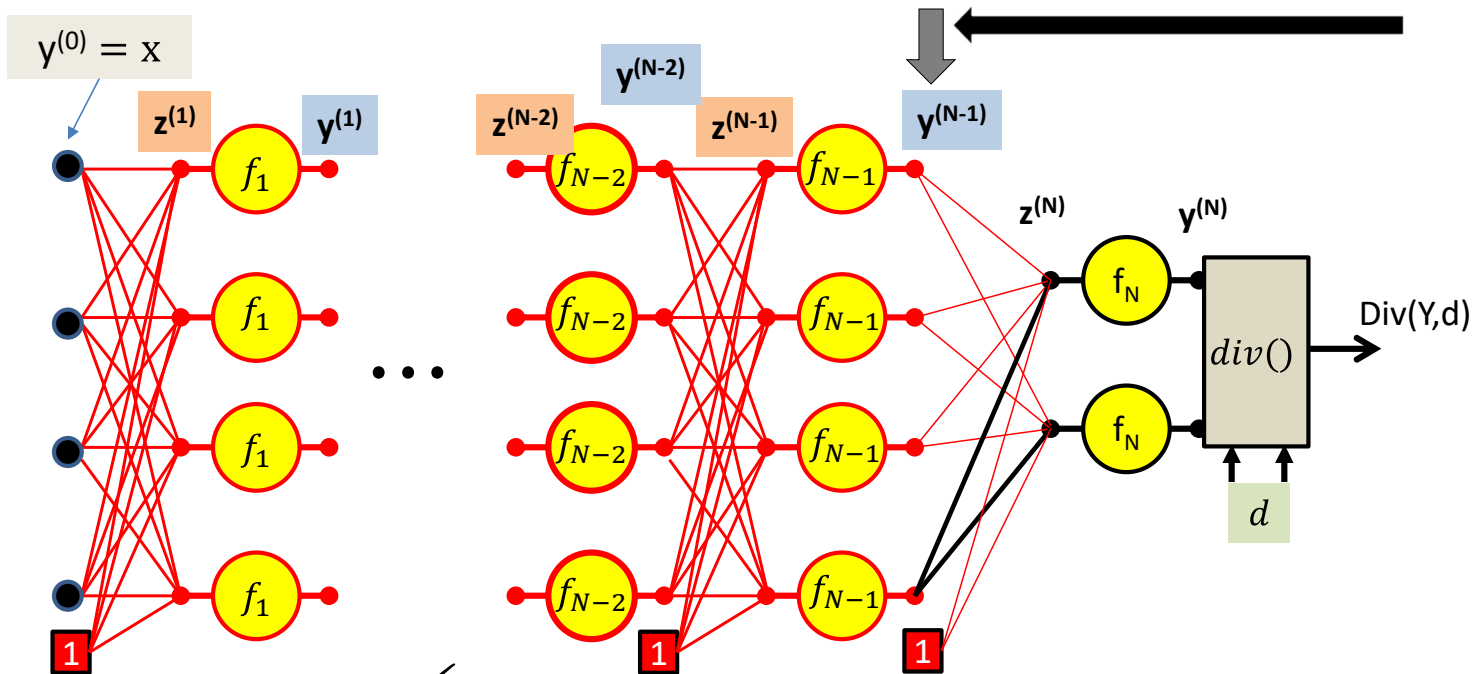
For the bias term $y_0^{(N-1)} = 1$

Computing derivatives



$$\frac{\partial Div}{\partial y_i^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}}$$

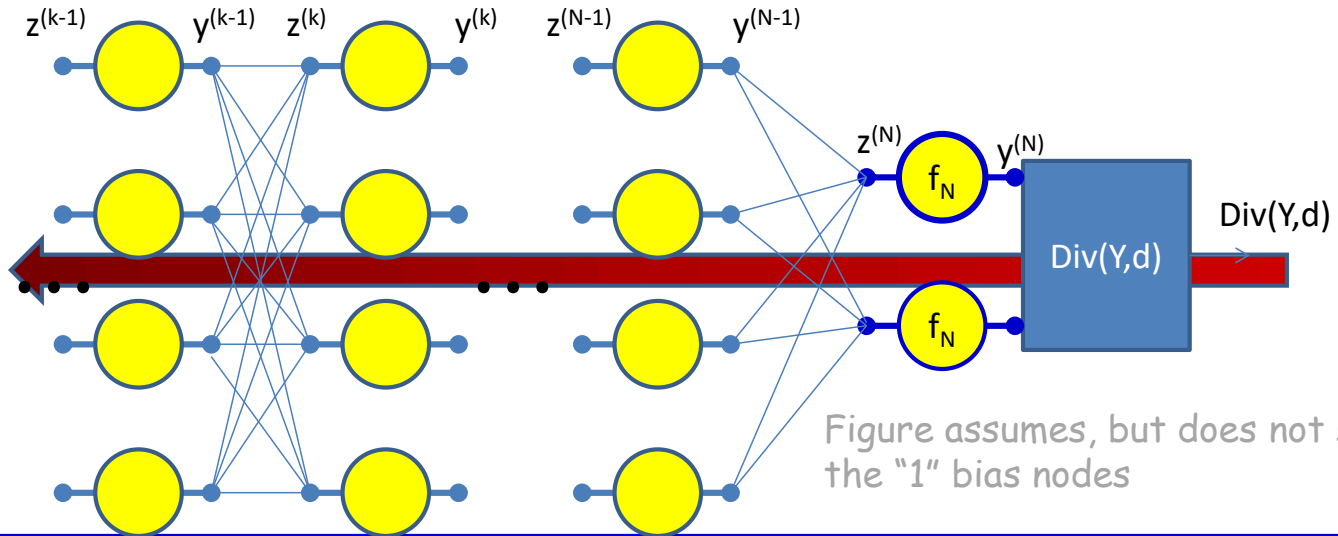
Computing derivatives



$$\frac{\partial Div}{\partial y_i^{(N-1)}} = \sum_j w_{ij}^{(N)} \frac{\partial Div}{\partial z_j^{(N)}}$$

$$\partial_y (w \cdot y)$$

Gradients: Backward Computation



Initialize: Gradient w.r.t network output

$$\frac{\partial Div}{\partial y_i^{(N)}} = \frac{\partial Div(Y, d)}{\partial y_i}$$

$$\frac{\partial Div}{\partial z_i^{(N)}} = f'_k(z_i^{(N)}) \frac{\partial Div}{\partial y_i^{(N)}}$$

For $k = N - 1..0$

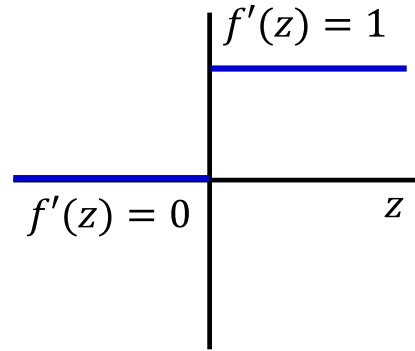
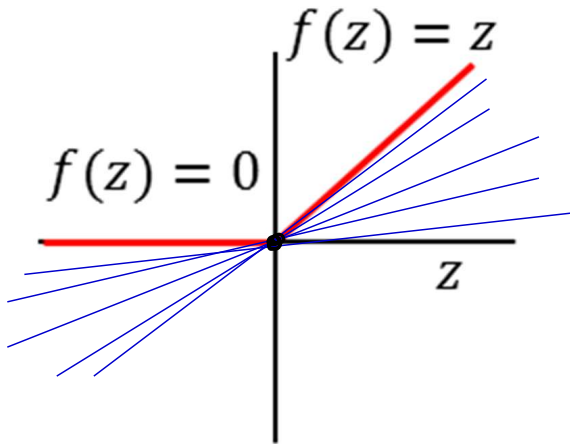
For $i = 1: \text{layer width}$

$$\frac{\partial Div}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \frac{\partial Div}{\partial z_j^{(k+1)}}$$

$$\frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$$

$$\forall j \frac{\partial Div}{\partial w_{ij}^{(k+1)}} = y_i^{(k)} \frac{\partial Div}{\partial z_j^{(k+1)}}$$

Subgradients and the RELU



$$f'(z) = \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \end{cases}$$

- The *subderivative* of a RELU is the slope of any line that lies entirely under it
 - The subgradient is a generalization of the subderivative
 - At the differentiable points on the curve, this is the same as the gradient
- Can use any subgradient at 0
 - Typically, will use the equation given

Consider a network from $\overset{\text{input}}{\mathbb{R}^2} \rightarrow \mathbb{R}^{1000} \rightarrow \mathbb{R}^{1000} \rightarrow \mathbb{R}$

Would computing the gradient with respect to the network weights be orders of magnitude slower than computing the gradient with respect to the input?

Must you compute the derivative wrt weights in order to compute derivative wrt the input?

No. The boxes in red below don't need the values in the yellow box

	<p>For $k = N - 1..0$ For $i = 1:\text{layer width}$</p>	
How many flops? →	$\frac{\partial Div}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \frac{\partial Div}{\partial z_j^{(k+1)}} \quad \frac{\partial Div}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial Div}{\partial y_i^{(k)}}$	$\sim 10^6$ flops
How many flops? →	$\forall j \frac{\partial Div}{\partial w_{ij}^{(k+1)}} = y_i^{(k)} \frac{\partial Div}{\partial z_j^{(k+1)}}$	$\sim 10^6$ flops

Multivariate input & output

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$J_f \in \mathbb{R}^{m \times n}$$

$$J_{f, i_j} =$$

Chain rule - multivariate input & output

$$y(x) = f(g(x))$$

$$J_y(x) = \underbrace{J_f(g(x))}_{\mathbb{R}^{m \times n}} \underbrace{J_g(x)}_{\mathbb{R}^{n \times p}} = \underbrace{\quad}_{\mathbb{R}^{m \times p}}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$g: \mathbb{R}^p \rightarrow \mathbb{R}^n$$

$$f \circ g: \mathbb{R}^p \rightarrow \mathbb{R}^m$$

Example

$$z(y) = \underbrace{W}_{\mathbb{R}^{m \times n}} \underbrace{y}_{\mathbb{R}^n} + \underbrace{b}_{\mathbb{R}^m}$$

$$J_z(y) = W$$

Example 0

Let $f: \mathbb{R}^m \rightarrow \mathbb{R}$ $W \in \mathbb{R}^{m \times n}$ $y \in \mathbb{R}^n$

$$D = f(\underbrace{Wy + b}_z)$$

$$J_D(y) = J_f(z)W$$

$$J_D(b) = J_f(z)$$

$$J_D(W) = J_f^t(z) y^t$$

