

Backpropagation

Goodfellow et al. 2016 (sec. 6.5)

Gradient descent

$$\widehat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i)$$

- SGD step requires computing the gradient of the loss

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha^{(t)} \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_t), y_t)$$

- Use chain rule

$$\frac{\partial}{\partial \boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}), y) = \ell'(f_{\boldsymbol{\theta}}(\mathbf{x}), y) \frac{\partial}{\partial \boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\mathbf{x})$$



BACKPROPAGANDA

History



S. Linnainmaa
1970



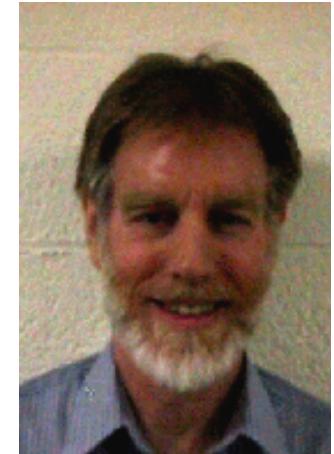
P. Werbos
1974



D. Rumelhart



G. Hinton



R. Williams

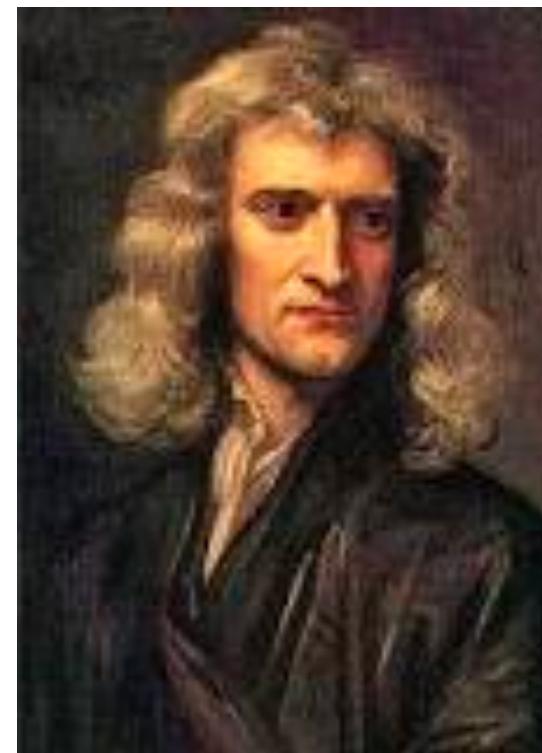
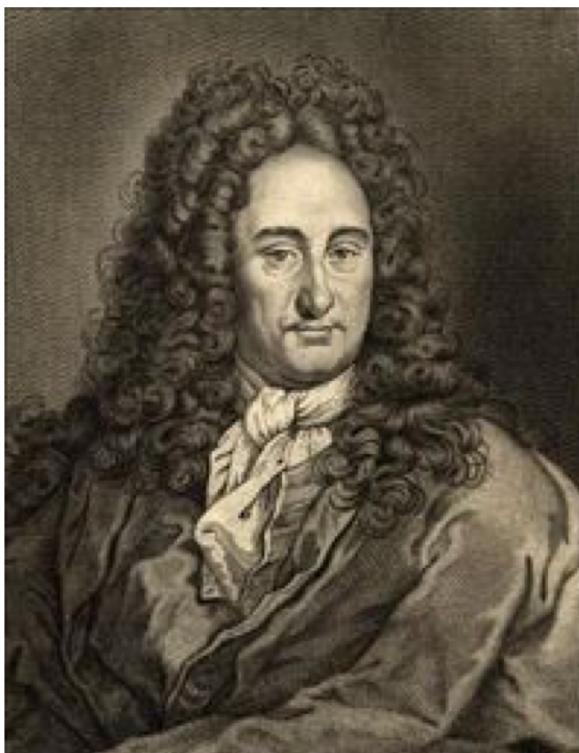


Letter | Published: 09 October 1986

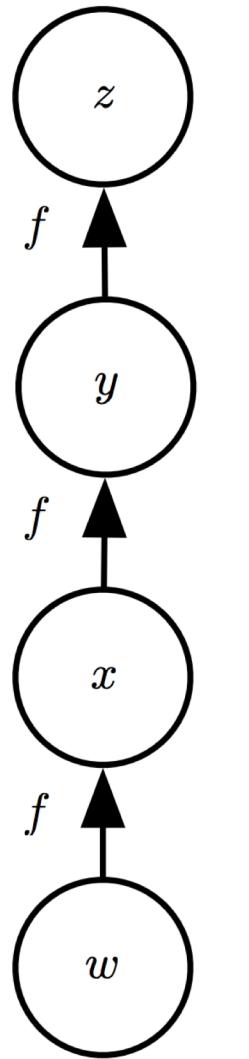
Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

History



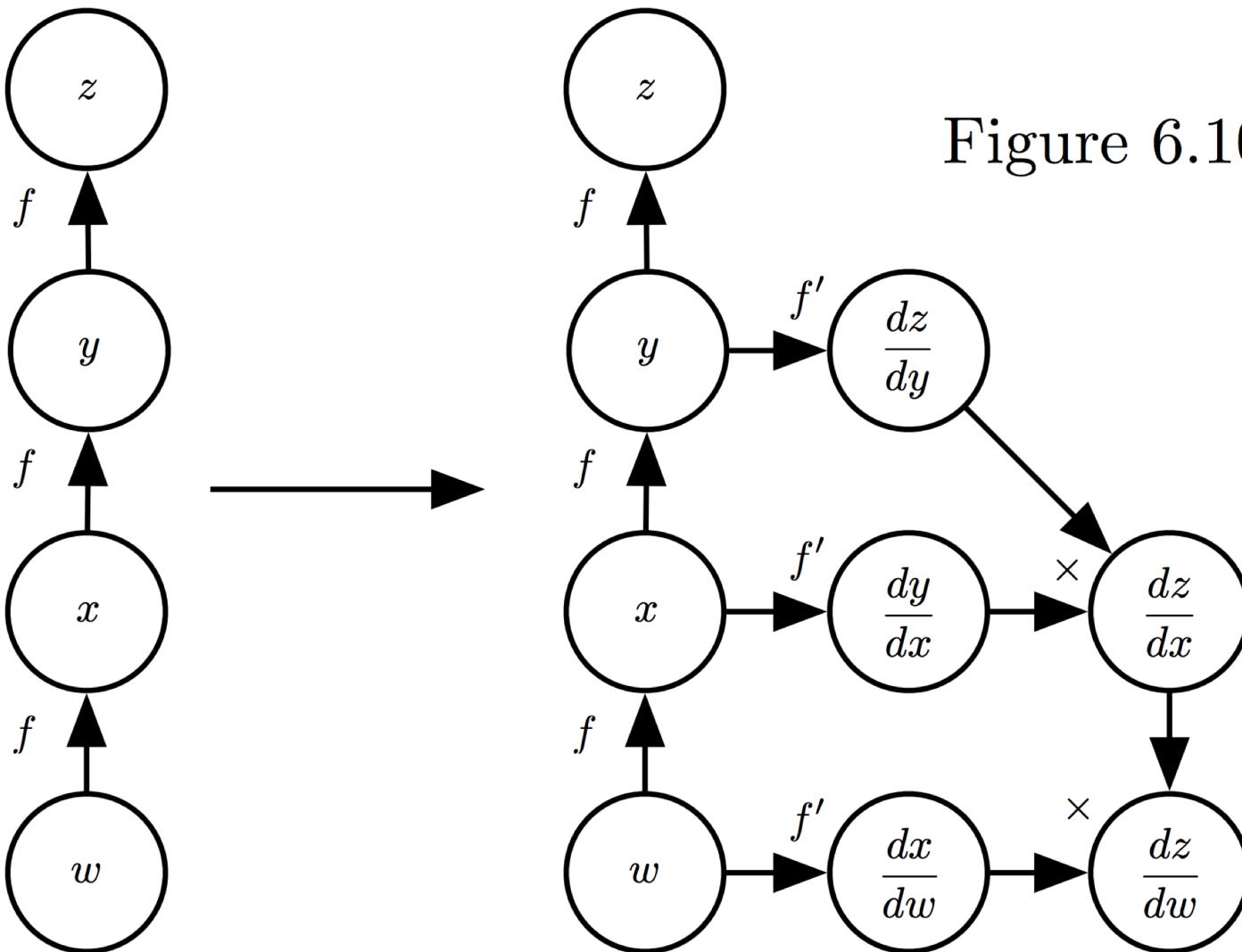
Scalar example



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

Back-prop avoids computing this twice

Scalar example



Vector case



$$\begin{array}{ll} \mathbf{x} \in \mathbb{R}^m & \\ \mathbf{z}' = \mathbf{Bx} & \mathbf{z}' \in \mathbb{R}^n \\ \mathbf{y}' = \varphi(\mathbf{z}') & \mathbf{y}' \in \mathbb{R}^n \\ \mathbf{z} = \mathbf{A}\mathbf{y}' & \mathbf{z} \in \mathbb{R}^k \\ \mathbf{y} = \psi(\mathbf{z}) & \mathbf{y} \in \mathbb{R}^k \\ l = \ell(\mathbf{y}) & l \in \mathbb{R} \end{array}$$

Slide credit: A. Bronstein

Chain rule

$$\frac{\partial l}{\partial \mathbf{A}} = \frac{\partial \mathbf{y}}{\partial \mathbf{A}} \frac{\partial l}{\partial \mathbf{y}} = \sum_{i=1}^k \frac{\partial y_i}{\partial \mathbf{A}} \frac{\partial l}{\partial y_i}$$

$$\delta \mathbf{A} = \frac{\partial \mathbf{y}}{\partial \mathbf{A}} \delta \mathbf{y} = \sum_{i=1}^k \frac{\partial y_i}{\partial \mathbf{A}} \delta y_i$$

$$\begin{aligned}\delta \mathbf{y} &= \frac{\partial l}{\partial \mathbf{y}} = \begin{pmatrix} \frac{\partial l}{\partial y_1} \\ \vdots \\ \frac{\partial l}{\partial y_k} \end{pmatrix} \\ \delta \mathbf{A} &= \frac{\partial l}{\partial \mathbf{A}} \\ &= \begin{pmatrix} \frac{\partial l}{\partial a_{11}} & \frac{\partial l}{\partial a_{1n}} \\ \vdots & \vdots \\ \frac{\partial l}{\partial a_{k1}} & \frac{\partial l}{\partial a_{kn}} \end{pmatrix}\end{aligned}$$

Chain rule

$$y_i = \psi \left(\sum_{j=1}^n a_{ij} y'_j \right) = \psi(z_i)$$

$$\frac{\partial y_i}{\partial \mathbf{A}} = \begin{pmatrix} - & \mathbf{-0-} & - \\ \psi'(z_i)y'_1 & \cdots & \psi'(z_i)y'_n \\ - & \mathbf{-0-} & - \end{pmatrix}$$

Chain rule

$$\frac{\partial y_i}{\partial \mathbf{A}} = \begin{pmatrix} - & -\mathbf{0}- & - \\ \psi'(z_i)y'_1 & \cdots & \psi'(z_i)y'_n \\ - & -\mathbf{0}- & - \end{pmatrix}$$

$$\delta \mathbf{A} = \sum_{i=1}^k \frac{\partial y_i}{\partial \mathbf{A}} \delta y_i = \begin{pmatrix} \delta y_1 \psi'(z_1) y'_1 & \cdots & \delta y_1 \psi'(z_1) y'_n \\ \vdots & \ddots & \vdots \\ \delta y_k \psi'(z_k) y'_1 & \cdots & \delta y_k \psi'(z_k) y'_n \end{pmatrix}$$

$$= \begin{pmatrix} \delta y_1 & & \\ & \ddots & \\ & & \delta y_k \end{pmatrix} \begin{pmatrix} \psi'(z_1) & & \\ & \ddots & \\ & & \psi'(z_k) \end{pmatrix} \begin{pmatrix} -\mathbf{y}'^T & - \\ \vdots & \\ -\mathbf{y}'^T & - \end{pmatrix}$$

Chain rule

$$\frac{\partial l}{\partial \mathbf{B}} = \frac{\partial \mathbf{y}'}{\partial \mathbf{B}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \frac{\partial l}{\partial \mathbf{y}} = \frac{\partial \mathbf{y}'}{\partial \mathbf{B}} \delta \mathbf{y}' = \sum_{i=1}^n \frac{\partial y'_i}{\partial \mathbf{B}} \delta y'_i$$

$$\delta \mathbf{y}' = \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \frac{\partial l}{\partial \mathbf{y}} = \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \delta \mathbf{y}$$

Chain rule

$$y'_i = \varphi \left(\sum_{j=1}^m b_{ij} x_j \right) = \varphi(z'_i)$$

$$\frac{\partial y'_i}{\partial \mathbf{B}} = \begin{pmatrix} - & -\mathbf{0}- & - \\ \varphi'(z'_i)x_1 & \cdots & \varphi'(z'_i)x_m \\ - & -\mathbf{0}- & - \end{pmatrix}$$

Chain rule

$$\frac{\partial y'_i}{\partial \mathbf{B}} = \begin{pmatrix} - & -\mathbf{0}- & - \\ \varphi'(z'_i)x_1 & \cdots & \varphi'(z'_i)x_m \\ - & -\mathbf{0}- & - \end{pmatrix}$$

$$\delta \mathbf{B} = \sum_{i=1}^n \frac{\partial y'_i}{\partial \mathbf{B}} \delta y'_i$$

$$= \begin{pmatrix} \delta y'_1 & & \\ & \ddots & \\ & & \delta y'_n \end{pmatrix} \begin{pmatrix} \varphi'(z'_1) & & \\ & \ddots & \\ & & \varphi'(z'_n) \end{pmatrix} \begin{pmatrix} -\mathbf{x}^T - \\ \vdots \\ -\mathbf{x}^T - \end{pmatrix}$$

Chain rule

$$y_i = \psi \left(\sum_{j=1}^n a_{ij} y'_j \right) = \psi(z_i)$$

$$\frac{\partial y_i}{\partial y'_j} = \psi'(z_i) a_{ij}$$

Chain rule

$$\frac{\partial y_i}{\partial y'_j} = \psi'(z_i) a_{ij}$$

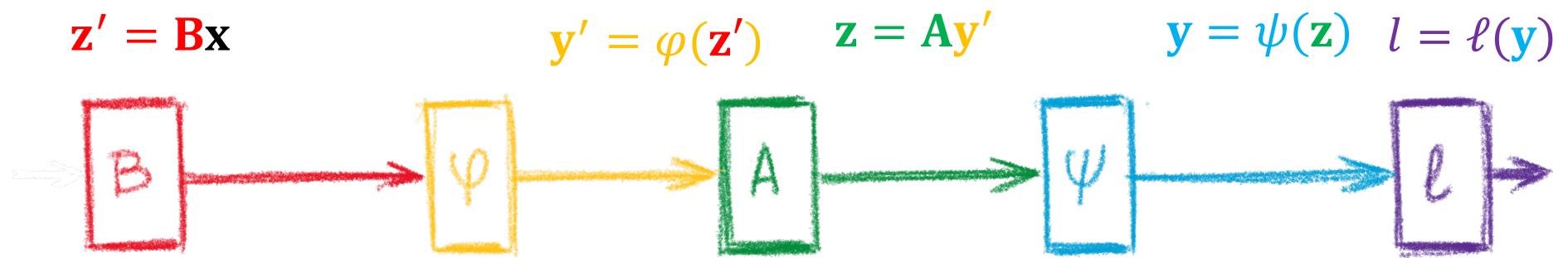
$$\frac{\partial \mathbf{y}}{\partial \mathbf{y}'} = \begin{pmatrix} \frac{\partial y_1}{\partial y'_1} & \dots & \frac{\partial y_k}{\partial y'_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial y'_n} & \dots & \frac{\partial y_k}{\partial y'_n} \end{pmatrix} = \begin{pmatrix} \psi'(z_1)a_{11} & \dots & \psi'(z_k)a_{k1} \\ \vdots & \ddots & \vdots \\ \psi'(z_1)a_{1n} & \dots & \psi'(z_k)a_{kn} \end{pmatrix}$$

Chain rule

$$\frac{\partial \mathbf{y}}{\partial \mathbf{y}'} = \begin{pmatrix} \psi'(z_1) & & \\ & \ddots & \\ & & \psi'(z_k) \end{pmatrix} \mathbf{A}^T$$

$$\delta \mathbf{y}' = \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \delta \mathbf{y} = \begin{pmatrix} \psi'(z_1) & & \\ & \ddots & \\ & & \psi'(z_k) \end{pmatrix} \mathbf{A}^T \delta \mathbf{y}$$

Chain rule



Chain rule

$$\delta \mathbf{y}' = \begin{pmatrix} \vdots \\ \psi'(z_1) \\ \vdots \end{pmatrix} \mathbf{A}^T \delta \mathbf{y}$$
$$\delta \mathbf{y} = \frac{\partial l}{\partial \mathbf{y}}$$
$$\delta \mathbf{B} = \begin{pmatrix} \vdots \\ \delta \mathbf{y}' \\ \vdots \end{pmatrix} \left(\begin{pmatrix} \vdots \\ \varphi'(\mathbf{z}) \\ \vdots \end{pmatrix} \right) \begin{pmatrix} -\mathbf{x}^T & - \\ \vdots & \end{pmatrix}$$
$$\delta \mathbf{A} = \begin{pmatrix} \vdots \\ \delta \mathbf{y} \\ \vdots \end{pmatrix} \left(\begin{pmatrix} \vdots \\ \psi'(\mathbf{z}) \\ \vdots \end{pmatrix} \right) \begin{pmatrix} -\mathbf{y}'^T & - \\ \vdots & \end{pmatrix}$$

Forward and backward pass

Input: $\mathbf{y}_0 = \mathbf{x}$

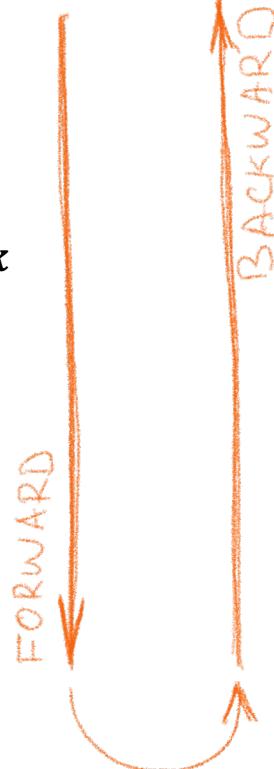
For $k = 1, \dots, L$

$$\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k$$

$$\mathbf{y}_k = \varphi_k(\mathbf{z}_k)$$

Output: $\mathbf{y} = \mathbf{y}_L$

Loss: $l = \ell(\mathbf{y})$



For $k = L, L - 1, \dots, 1$

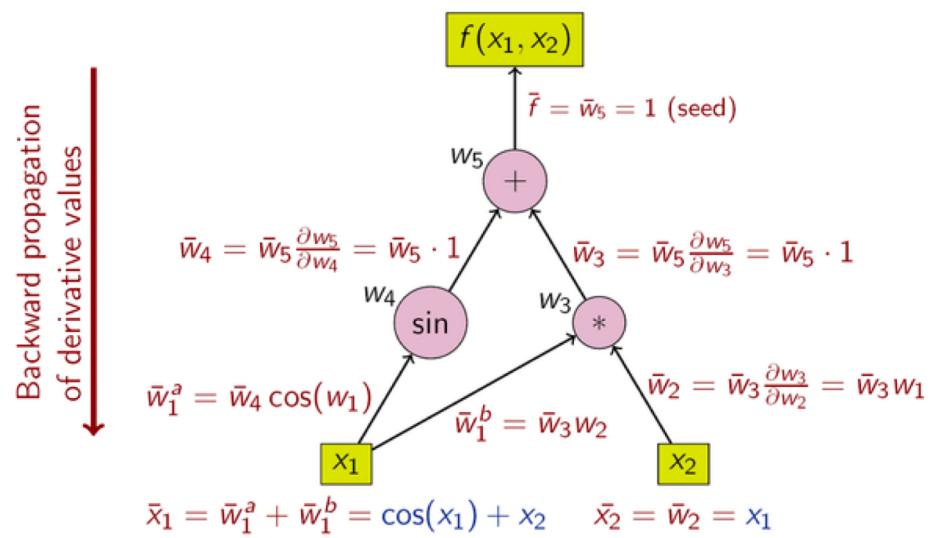
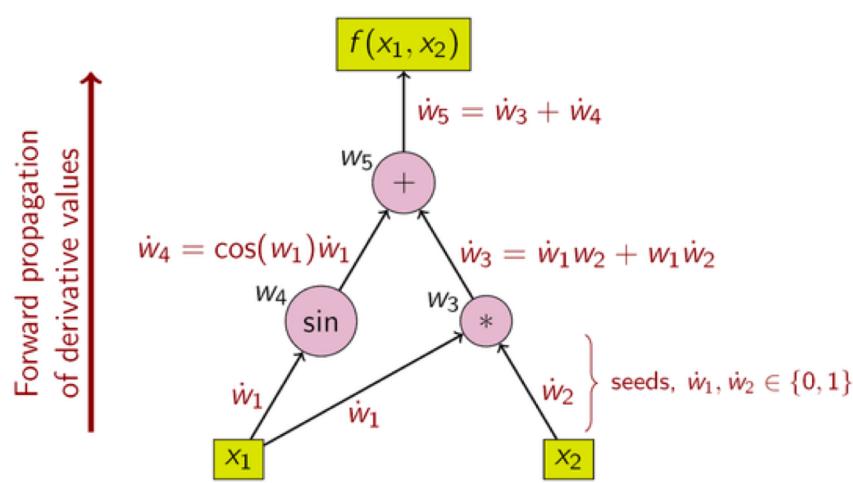
$$\delta \mathbf{W}_k = \begin{pmatrix} \vdots & \vdots \\ \delta \mathbf{y}_k & \varphi'_k(\mathbf{z}_k) \\ \vdots & \vdots \end{pmatrix} \mathbf{1} \mathbf{y}_{k-1}^T$$

$$\delta \mathbf{b}_k = \begin{pmatrix} \vdots & \vdots \\ \delta \mathbf{y}_k & \varphi'_k(\mathbf{z}_k) \\ \vdots & \vdots \end{pmatrix}$$

$$\delta \mathbf{y}_{k-1} = \begin{pmatrix} \vdots & \vdots \\ \varphi'_k(\mathbf{z}_k) & \mathbf{W}_k^T \delta \mathbf{y}_k \\ \vdots & \vdots \end{pmatrix}$$

Loss grad: $\delta \mathbf{y}_L = \nabla \ell(\mathbf{y}_L)$

Automatic differentiation



“Differentiable programming”