

① March 13

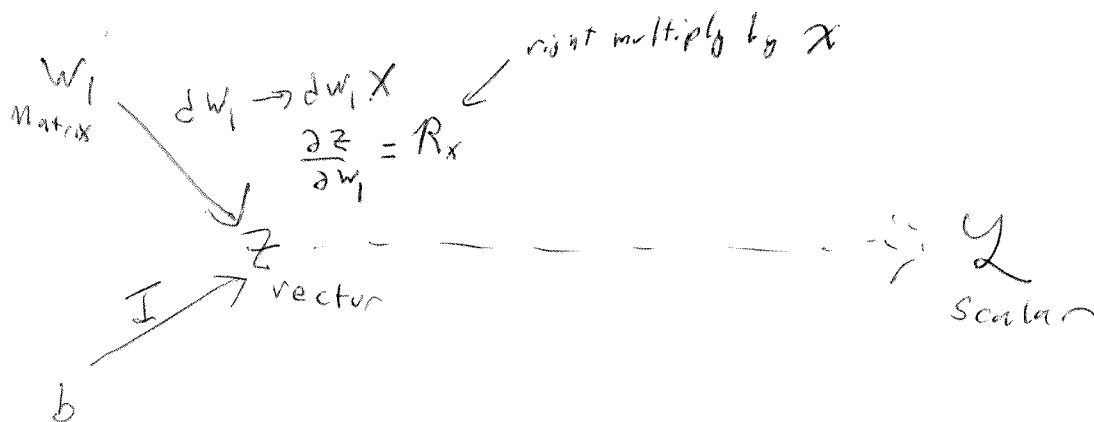
CR lecture 10
Neural Network

You've seen this
but it's worth going through again
+ in depth

x_i inputs (no variation)
black in CR's diagram

Forward Pass

$$\begin{cases} z = W_1 x + b_1 \\ h = \sigma_o(z) \\ y = W_2 h + b_2 \\ \mathcal{L} = \frac{1}{2} \|y - t\|^2 \end{cases}$$



$$\frac{\partial \mathcal{L}}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial W_1} \leftarrow \text{This gets confusing}$$

but remember $\nabla_{W_1} \mathcal{L}$ is a matrix (same shape as W_1)

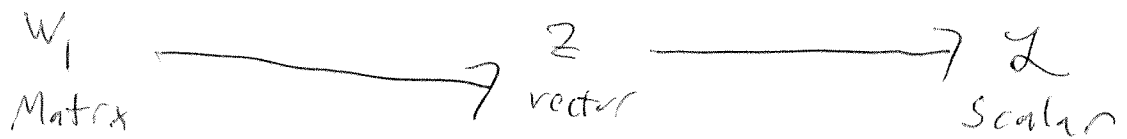
$$\frac{\partial \mathcal{L}}{\partial W_1} = \nabla_{W_1} \mathcal{L} \cdot \frac{\partial z}{\partial W_1} \quad (\text{directional derivative})$$

$$= \text{tr}((\nabla_{W_1} \mathcal{L})^T \frac{\partial z}{\partial W_1})$$

$$(\text{tr } A^T B = A \cdot B = \frac{\text{sum}(A \cdot B)}{\text{math}} = \frac{\text{sum}(A \cdot B)}{\text{Julia}})$$

$$\frac{\partial \mathcal{L}}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial z} [\frac{\partial z}{\partial W_1}]$$

(2) March 13



Combined

$$W_1 \xrightarrow[\nabla_{W_1} \mathcal{L} \circ dW_1]{dW_1 \rightarrow \frac{d\mathcal{L}}{dW_1} [dW_1]} \mathcal{L} \text{ scalar}$$

$$\text{tr}((\nabla_{W_1} \mathcal{L})^T dW_1)$$

$$\text{tr} = [\mathbf{I} \otimes (\nabla_{W_1} \mathcal{L})^T]$$

Any nice scalar function of a matrix

The pieces

$$W_1 \xrightarrow[dW_1 \rightarrow dW_1 X]{\quad} z$$

$R_X \leftarrow \text{right mult operator}$

$X^T \otimes \mathbf{I}$

Why no gradient?

$$z \xrightarrow[\begin{matrix} g^T dz \\ g \circ dz \end{matrix}]{\nabla_z \mathcal{L} \circ dz} \mathcal{L}$$

$g = \text{gradient}$

$$\frac{d\mathcal{L}}{dz} dz$$

row vector

$$\text{tr}((\nabla_{W_1} \mathcal{L})^T dW_1) = (\nabla_z \mathcal{L})^T dW_1 X \quad \text{cyclic prop}$$

$$= \text{tr}(X (\nabla_z \mathcal{L})^T dW_1)$$

$$\Rightarrow (\nabla_{W_1} \mathcal{L}) = (\nabla_z \mathcal{L}) X^T$$

(3)

The reverse mode people are good at back propagating gradients

First they make the notation less daunting

$$\nabla_{\text{any}} \mathcal{L} = \overline{\text{any}}$$

so

$$\boxed{\overline{W_1} = \overline{z} x^T}$$

pushback
note that
we go from \overline{z} to $\overline{W_1}$

Let's try another example one step away

$$\overline{z} \xrightarrow{D = \text{diag}(\sigma'(z))} \underset{\text{vector}}{h} \xrightarrow{\overline{h}^T} \underset{\text{scalar}}{y}$$

$$dy = \overline{h}^T (D dz) = (D \overline{h})^T dz$$

$$\boxed{\overline{z}} = D \overline{h} = \boxed{\sigma'(z) \circ \overline{h}}$$

$$y = w_2 h + b_2$$

$$\underset{\text{vector}}{h} \xrightarrow{dy = w_2 dh} \underset{\text{vector}}{y} \xrightarrow{\overline{y}^T} \underset{\text{scalar}}{y}$$

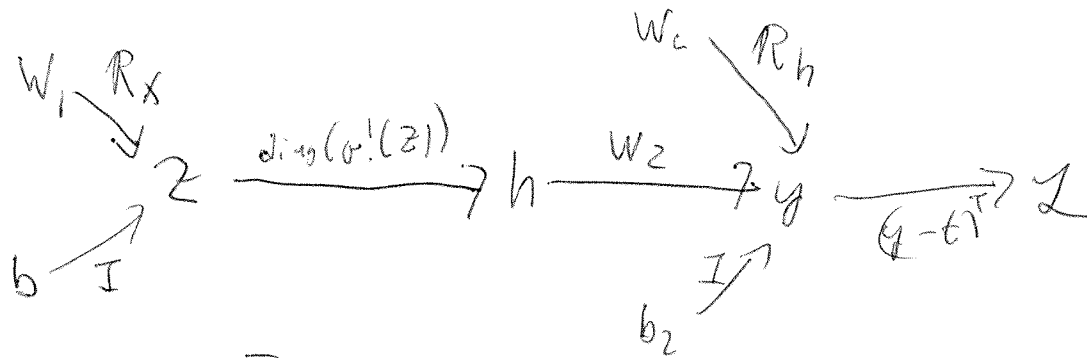
$$\overline{h}^T dy = \overline{y}^T w_2 dh$$

$$= (w_2^T \overline{y})^T dh$$

$$\overline{h} = w_2^T \overline{y}$$

(4)

By now perhaps you can see the pattern and not follow the slow way but do it the reverse way



$$\bar{y} = 1$$

$$\bar{y} = \bar{z} (y - t) \quad (\text{which is } y - t)$$

$$\bar{w}_2 = \bar{y} h^T$$

$$\bar{b}_2 = \bar{y}$$

$$\bar{h} = w_2^T \bar{y}$$

$$\bar{z} = \bar{h} \cdot \sigma'_z(z) \quad (\text{same as } \text{diag}^T)$$

$$\bar{w}_1 = \bar{z} x^T \quad (\bar{z} R_{x^T})$$

$$\bar{b}_1 = \bar{z}$$

(5)

Why is this called an adjoint?

In general if L is a linear operator from a vector space X to Y

we look for an adjoint L^* (I may use L^T)

$$\text{s.t. } \langle y, Lx \rangle = \langle L^T y, x \rangle \quad \text{for all } x \in X, y \in Y.$$

For us $\langle -, - \rangle =$ vector or matrix dot products

$$\text{e.g. } (R_A)^T = R_{A^T} \text{ proof}$$

$$\langle y, Ax \rangle = y^T Ax = \cancel{y^T A} \langle A^T y, x \rangle$$

e.g. $x \rightarrow v \cdot x$ is self-adjoint

$$\langle y, v \cdot x \rangle = \sum_i y_i (v_i x_i) = \langle v \cdot y, x \rangle$$

So you can take adjoints

from end to start to

compute gradients

(7)

really it's just

$$x_1 \xrightarrow{J_1} x_2 \xrightarrow{J_2} x_3 \xrightarrow{J_3 = g^T} \text{scalar}$$

$$\underbrace{(J_1^T J_2^T g)^T}_{\text{reverse}} = \left(g^T \underbrace{(J_2 J_1)}_{\substack{\uparrow \\ \text{forward}}} \right)$$

$$g^T J = (g^T J_2) \underbrace{J_1}_{\text{reverse}}$$

vjp - vector transpose Jacobian

What is $\left(\frac{d}{dx} \right)^T$?