General rules for backprop:

input — module — output

- How does output $\times$ signal from $\leftarrow$ how module
  change wrt module    output       changes wrt
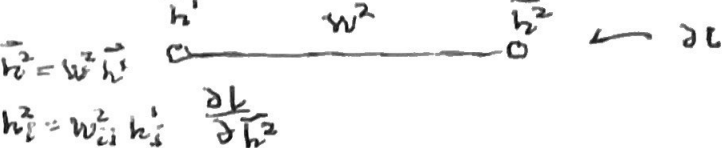                                    signal

* How does output $\times$ signal from $\leftarrow$ how signal
  change wrt input     output       propagates
                                    backwards

* Note that the "multiplication" rule is just a rule. For general interaction. It's a personal responsibility to determine the order of multiplication for shape compatability.

- Rules for shape

Let $s$ be a scalar, $\vec{v} \in \mathbb{R}^{k \times 1}$, $\vec{y} \in \mathbb{R}^{m \times 1}$, $\vec{x} \in \mathbb{R}^{n \times 1}$, and $W \in \mathbb{R}^{m \times n}$

then $\frac{\partial s}{\partial \vec{v}} \in \mathbb{R}^{1 \times k}$, $\frac{\partial \vec{y}}{\partial \vec{x}} \in \mathbb{R}^{m \times n}$, and $\frac{\partial s}{\partial W} \in \mathbb{R}^{m \times n}$
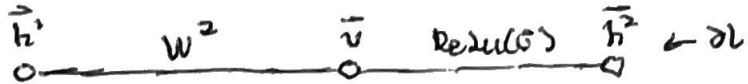
- Dealing with pure linear modules

$$\vec{h}^1 \quad\underset{W^2}{\rule{4cm}{0.4pt}}\quad \vec{h}^2 \quad\leftarrow \partial L$$

$\vec{h}^2 = W^2 \vec{h}^1$

$h_i^2 = W_{ij}^2 h_j^1 \qquad \frac{\partial L}{\partial \vec{h}^2}$

$\frac{\partial L}{\partial W^2} = \frac{\partial \vec{h}^2}{\partial W^2} \frac{\partial L}{\partial \vec{h}^2} \rightarrow \frac{\partial L}{\partial W_{os}^2} = \frac{\partial L}{\partial h_i^2} \times \frac{\partial h_i^2}{\partial W_{oj}^2} = \frac{\partial L}{\partial h_i^2} \cdot h_j^1 = \left(\frac{\partial L}{\partial \vec{h}^2}\right)^T \left(\vec{h}^1\right)^T$

$\frac{\partial L}{\partial \vec{h}_1} = \frac{\partial \vec{h}_2}{\partial \vec{h}_1} \frac{\partial L}{\partial \vec{h}_2} = \frac{\partial L}{\partial \vec{h}_2} W^2$  * switch for compatability

Dealing with functional and linear modules

$$\vec{h}^1 \quad\underset{W^2}{\rule{2cm}{0.4pt}}\quad \vec{v} \quad \text{ReLU}(\vec{v}) \quad \vec{h}^2 \leftarrow \partial L$$

$\frac{\partial L}{\partial \vec{h}^2}$

* Use $\vec{v}$ as an intermediary variable

$\vec{v} = W^2 \vec{h}^1$

$\vec{h}^2 = \text{ReLU}(\vec{v})$

$v_i = W_{ij}^2 h_j^1$

$h_i^2 = \text{ReLU}(W_{ij}^2 h_j^1)$

$h_i^2 = \text{ReLU}(v_i)$

$\frac{\partial h_i^2}{\partial v_j} = \delta_{ij} \mathbb{1}(v_i > 0)$

$\frac{\partial L}{\partial \vec{v}} = \frac{\partial \vec{h}^2}{\partial \vec{v}} \frac{\partial L}{\partial \vec{h}^2} = \text{ReLU}'(\vec{v}) \frac{\partial L}{\partial \vec{h}^2}$

ensure
row vector $\rightarrow \frac{\partial L}{\partial v_k} = \frac{\partial h_j^2}{\partial v_k} \frac{\partial L}{\partial h_j^2} = \frac{\partial L}{\partial h_i^2} \frac{\partial h_i^2}{\partial v_k} = \frac{\partial L}{\partial \vec{h}^2} \text{ReLU}'(\vec{v})$
shape

$\frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial \vec{v}} \frac{\partial \vec{v}}{\partial W^2} \rightarrow \frac{\partial L}{\partial v_i} \frac{\partial v_i}{\partial W_{ij}^2} \rightarrow \frac{\partial L}{\partial v_i} \cdot h_j^1$   ensure matrix shape

$= \left(\frac{\partial L}{\partial \vec{v}}\right)^T \left(\vec{h}^1\right)^T = \left(\frac{\partial L}{\partial \vec{h}^2} \cdot \text{diag}(\mathbb{1}(\vec{v} > 0))\right)^T \left(\vec{h}^1\right)^T = \text{diag}(\mathbb{1}(\vec{v} > 0))\left(\frac{\partial L}{\partial \vec{h}^2}\right)^T \left(\vec{h}^1\right)^T$

$\frac{\partial L}{\partial \vec{h}^1} = \frac{\partial L}{\partial \vec{v}} \frac{\partial \vec{v}}{\partial \vec{h}^1} = \frac{\partial L}{\partial v_i} \frac{\partial v_i}{\partial h_j^1} = \frac{\partial L}{\partial h_k^2} \frac{\partial h_k^2}{\partial v_i} \frac{\partial v_i}{\partial h_j^1} \sim (1 \times k) \times (k \times i) \times (i \times j) \sim (1 \times j)$

$= \frac{\partial L}{\partial \vec{h}^2} \cdot \text{ReLU}'(\vec{v}) W^2 = \frac{\partial L}{\partial \vec{h}^2} \cdot \text{diag}(\mathbb{1}(\vec{v} > 0)) \cdot W^2$

For linear modules (and functional modules)

signal from output $\approx$ how output changes wrt input $=$ how signal propalgates back

how output changes wrt module $\approx$ signal from output $=$ how module changes wrt module

Bonus : updating bias

$$\vec{y} = W\vec{x} + \vec{b}$$

$\overset{\vec{x}}{\circ}\!\!\!\rule[0.5ex]{4cm}{0.4pt}\!\!\!\overset{\vec{y}}{\circ} \leftarrow \partial L$

$$y_i = w_{ij} x_j + b_i \qquad \frac{\partial L}{\partial \vec{y}}$$

$$\frac{\partial L}{\partial \vec{b}} = \frac{\partial L}{\partial \vec{y}} \frac{\partial \vec{y}}{\partial \vec{b}} \rightarrow \frac{\partial L}{\partial b_k} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial b_k} = \frac{\partial L}{\partial y_i} \delta_{ik}$$

* Note that bias follows a different multiplication ordering

* Best to use general rule and use shape compatability to understand ordering