## Backpropagate $(y, A, \eta)$

- Given weight matrices $W^2, \ldots, W^L$ and bias vectors $\vec{b}^2, \ldots, \vec{b}^L$ at each layer,

- Declare vars:

  - $\vec{d}_{\vec{a}} C$, the derivative of cost w.r.t current layer's activations

  - $\vec{d}_{\vec{z}} C$, same w.r.t current layer's $\vec{z}$ values

  - $\Delta \vec{b}$, change in $\vec{b}$ at current layer

  - $\Delta W$, same for $W$

  - $\vec{a}$, current layer's activations

  - $\vec{a}'$, previous layer's activations

  - $\vec{y}$, vector corresponding to classification $y$

- $\vec{y} = \hat{e}_y$
- $\vec{a} = $ last element of $A$
- $\vec{d}_{\vec{a}} C = \vec{a} - \vec{y}$
- $\vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot \left[ \vec{a} \odot (\vec{1} - \vec{a}) \right]$

- For each layer $\ell$ from $L$ down to $2$:

  - $\vec{a} = \ell^{th}$ elem. of $A$
  - $\vec{a}' = (\ell-1)^{th}$ elem. of $A$

  - $\Delta \vec{b} = -\eta \, \vec{d}_{\vec{z}} C$
  - $\Delta W = -\eta \, \vec{d}_{\vec{z}} C \otimes \vec{a}'$

  - $\vec{d}_{\vec{a}} C = (W^{\ell})^t \, \vec{d}_{\vec{z}} C$
  - $\vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot \left[ \vec{a}' \odot (\vec{1} - \vec{a}') \right]$

  - $\vec{b}^{\ell} \mathrel{+}= \Delta \vec{b}$
  - $W^{\ell} \mathrel{+}= \Delta W$

## Train $(\vec{x}_1, \ldots, \vec{x}_n, y_1, \ldots, y_n, \eta, T)$

- Declare vars:
  - $A$, the list of activations returned by Feedforward

- Repeat $T$ times:
  - For each $\vec{x}, y$:
    - $A = $ Feedforward $(\vec{x})$
    - Backpropagate $(y, A, \eta)$