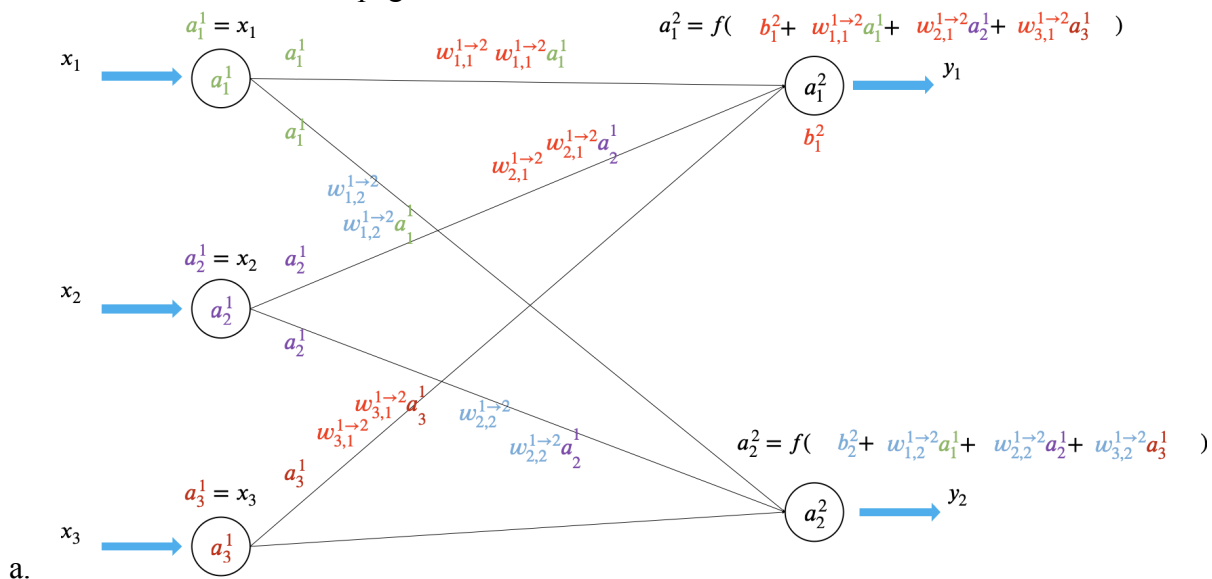


Supervised Learning VI – Neural Networks (cont.)

1. Neural Networks: Forward Propagation



2. NN as Computation Graph

a. Given a “batch” of examples X and the ground truth output Y gf

i. We know forward equations:

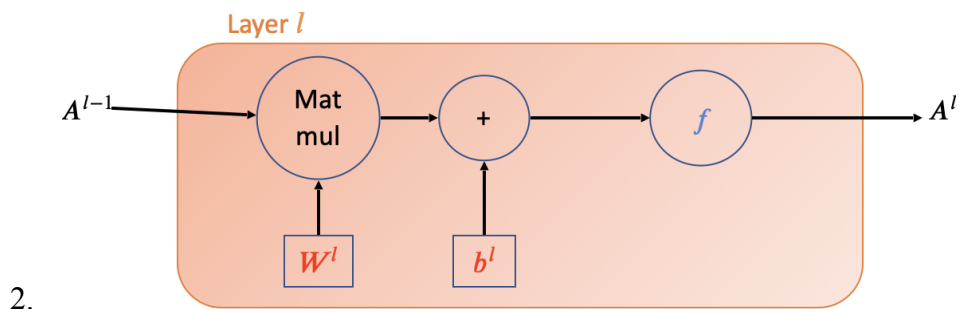
$$\mathbf{Z}^l = \mathbf{A}^l \mathbf{W}^{l-1 \rightarrow l} + \mathbf{b}^l$$

$$\mathbf{A}^l = f(\mathbf{Z}^l)$$

ii. Let's view a layer as a computation graph

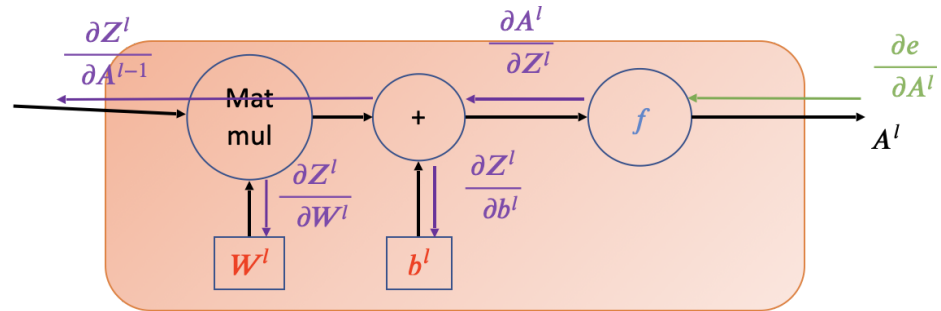
1. Layer has no idea where it is in the NN

a. That's ok! Chain rule to the rescue!



3. NN Computation Graph Backprop

- a. Each forward edge gets a backwards edge for chain rule!
 - i. We know these equations!



ii.

$$\frac{\partial A^l}{\partial Z^l} = f'(Z^l)$$

$$\frac{\partial Z^l}{\partial b^l} = \mathbf{1}^T$$

$$\frac{\partial Z^l}{\partial W^l} = (A^{l-1})^T$$

iii. $\frac{\partial Z^l}{\partial A^{l-1}} = (W^l)^T$

- b. Combine pieces of chain rule to get the derivatives we want

$$\frac{\partial e}{\partial Z^l} = f'(Z^l) \frac{\partial e}{\partial A^l} \quad \frac{\partial e}{\partial b^l} = \mathbf{1}^T \frac{\partial e}{\partial Z^l} \quad \frac{\partial e}{\partial W^l} = \frac{\partial e}{\partial Z^l} (A^{l-1})^T \quad \frac{\partial e}{\partial A^{l-1}} = (W^l)^T \frac{\partial e}{\partial Z^l}$$

- c. If we assume we have a “line graph”



- d. procedure forward(Matrix X) -> Matrix:

for(layer in layers):

X = layer.forward(X)

return X // now final predictions

procedure backwards(Matrix X, Matrix de_dA) -> Matrix:

// need to do a forward pass and cache the intermediary values

Stack[Matrix] cache = [X]

for(layer in layers):

X = layer.forward(X)

if layer is not last layer then: cache.push(X)

// now start at the last layer and walk backwards to the front

for(layer in reversed(layers)):

de_dA = layer.backwards(cache.pop(), de_dA)

return de_dA // now de_dX at this point

e.