# Neural networks

Training neural networks - backpropagation algorithm

# MACHINE LEARNING

**Topics:** stochastic gradient descent (SGD)

- Algorithm that performs updates after each example

  ‣ initialize $\boldsymbol{\theta}$  ( $\boldsymbol{\theta} \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \ldots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$ )

  ‣ for N iterations

   - for each training example $(\mathbf{x}^{(t)}, y^{(t)})$

     ✓ $\Delta = -\nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) - \lambda \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$

     ✓ $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \, \Delta$

   $\left.\phantom{\begin{array}{c}a\\a\\a\end{array}}\right\}$ training epoch
   =
   iteration over **all** examples

- To apply this algorithm to neural network training, we need

  ‣ the loss function $l(\mathbf{f}(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$

  ‣ a procedure to compute the parameter gradients $\nabla_{\boldsymbol{\theta}} l(\mathbf{f}(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$

  ‣ the regularizer $\Omega(\boldsymbol{\theta})$  (and the gradient $\nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$ )

  ‣ initialization method

# BACKPROPAGATION

**Topics:** backpropagation algorithm

• This assumes a forward propagation has been made before

‣ compute output gradient (before activation)

$$\nabla_{\mathbf{a}^{(L+1)}(\mathbf{x})} - \log f(\mathbf{x})_y \quad \Longleftarrow \quad -(\mathbf{e}(y) - \mathbf{f}(\mathbf{x}))$$

‣ for $k$ from $L+1$ to $1$

- compute gradients of hidden layer parameter

$$\nabla_{\mathbf{W}^{(k)}} - \log f(\mathbf{x})_y \quad \Longleftarrow \quad \left(\nabla_{\mathbf{a}^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y\right) \ \mathbf{h}^{(k-1)}(\mathbf{x})^\top$$

$$\nabla_{\mathbf{b}^{(k)}} - \log f(\mathbf{x})_y \quad \Longleftarrow \quad \nabla_{\mathbf{a}^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y$$

- compute gradient of hidden layer below

$$\nabla_{\mathbf{h}^{(k-1)}(\mathbf{x})} - \log f(\mathbf{x})_y \quad \Longleftarrow \quad \mathbf{W}^{(k)^\top} \left(\nabla_{\mathbf{a}^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y\right)$$
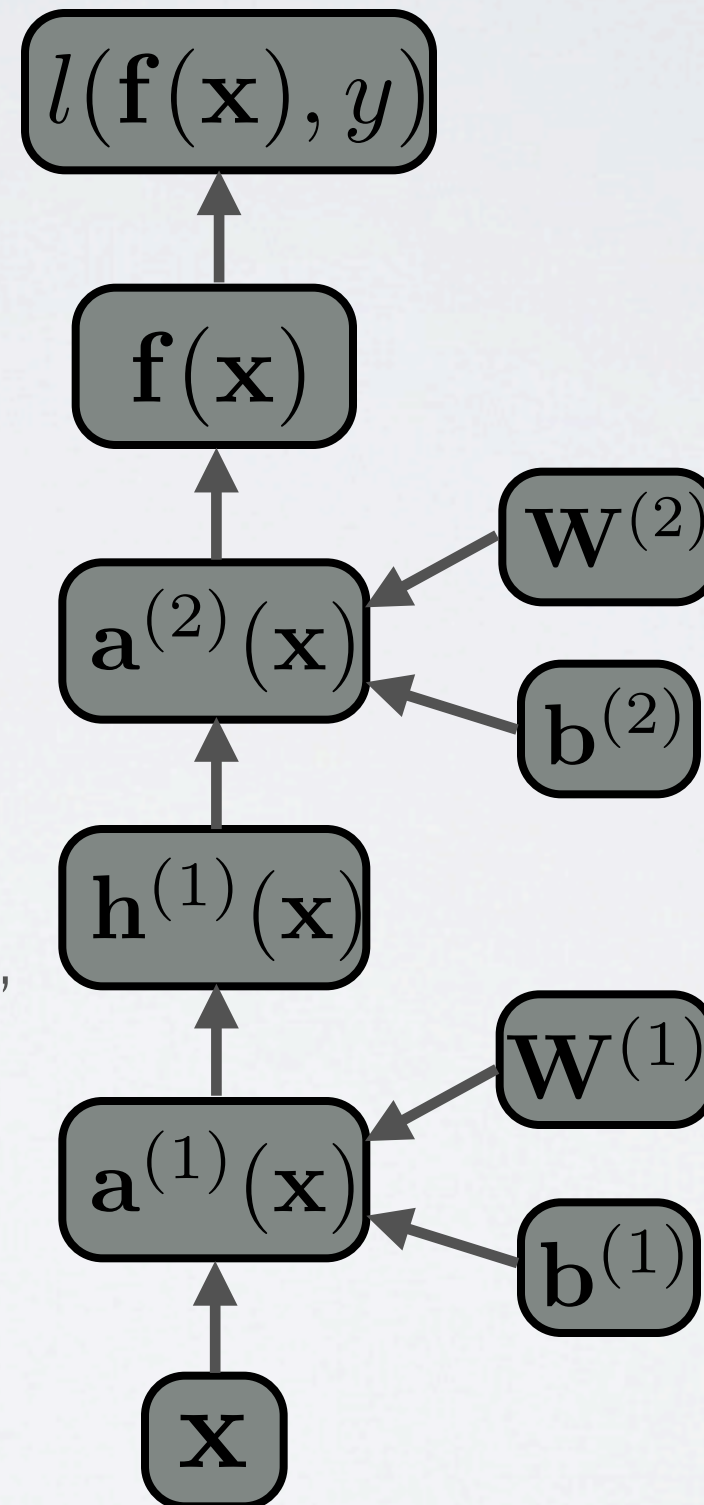
- compute gradient of hidden layer below (before activation)

$$\nabla_{\mathbf{a}^{(k-1)}(\mathbf{x})} - \log f(\mathbf{x})_y \quad \Longleftarrow \quad \left(\nabla_{\mathbf{h}^{(k-1)}(\mathbf{x})} - \log f(\mathbf{x})_y\right) \odot [\ldots, g'(a^{(k-1)}(\mathbf{x})_j), \ldots]$$
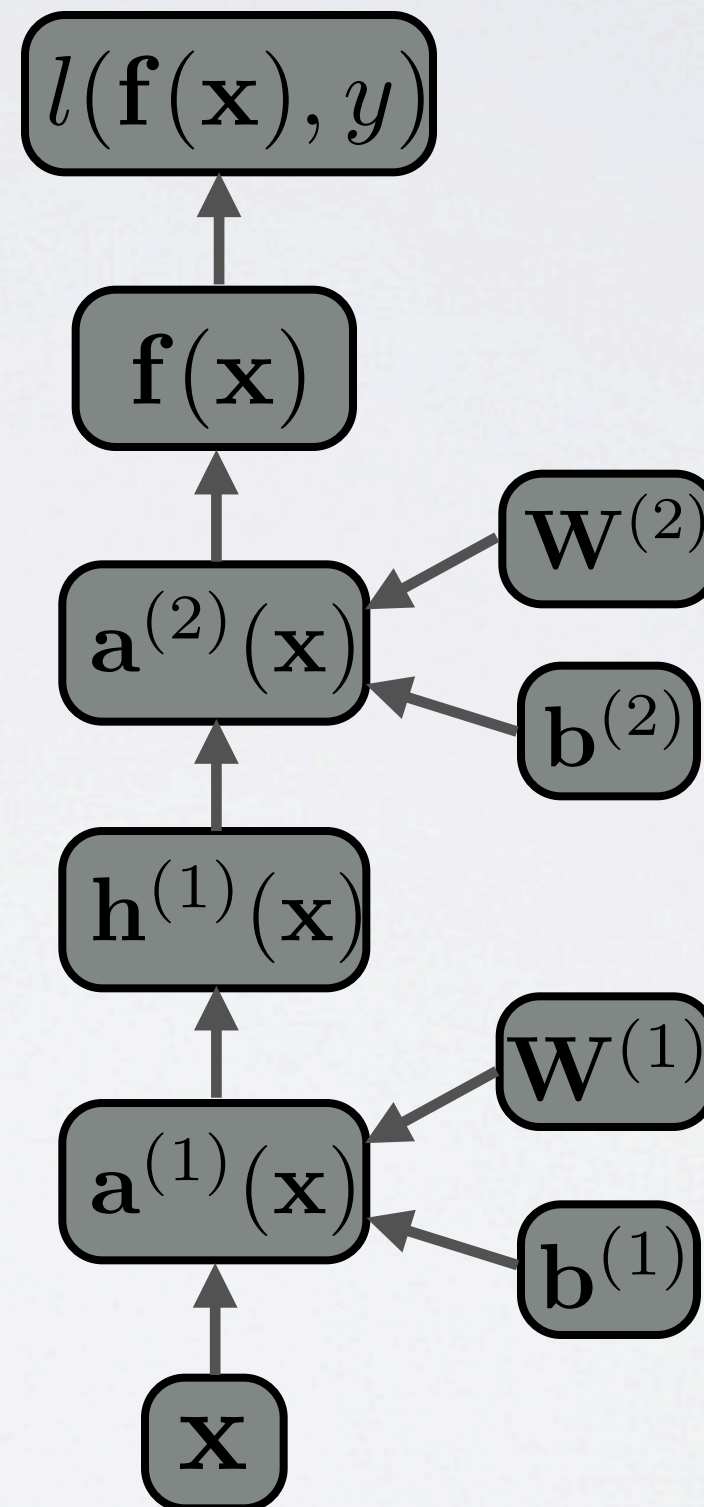
# FLOW GRAPH

**Topics:** flow graph

- Forward propagation can be represented as an acyclic flow graph

- It's a nice way of implementing forward propagation in a modular way

  ‣ each box could be an object with an fprop method, that computes the value of the box given its children

  ‣ calling the fprop method of each box in the right order yield forward propagation

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{W}^{(2)}$$

$$\mathbf{a}^{(2)}(\mathbf{x})$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{W}^{(1)}$$

$$\mathbf{a}^{(1)}(\mathbf{x})$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH

**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{W}^{(2)}$$

$$\mathbf{a}^{(2)}(\mathbf{x})$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{W}^{(1)}$$

$$\mathbf{a}^{(1)}(\mathbf{x})$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH
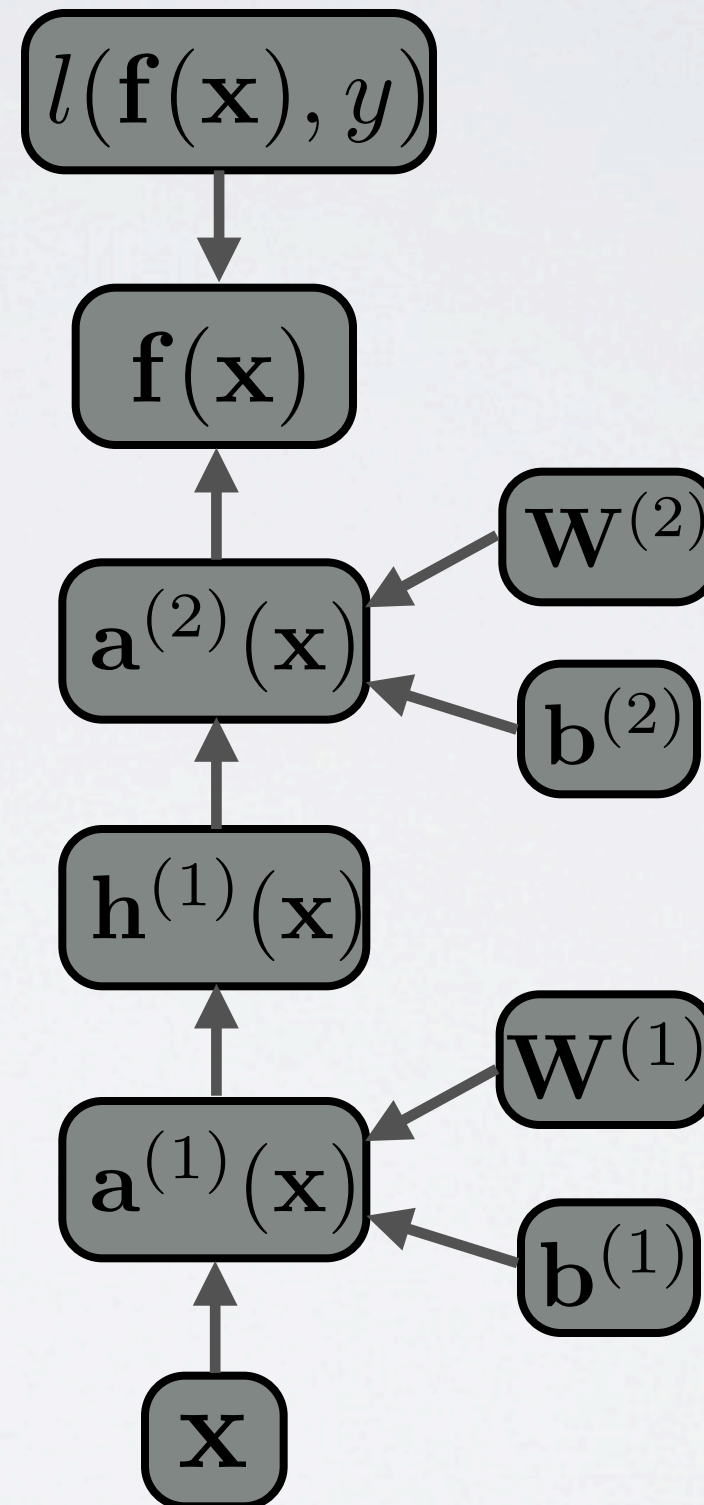
**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x}) \quad \mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{a}^{(1)}(\mathbf{x}) \quad \mathbf{W}^{(1)}$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH
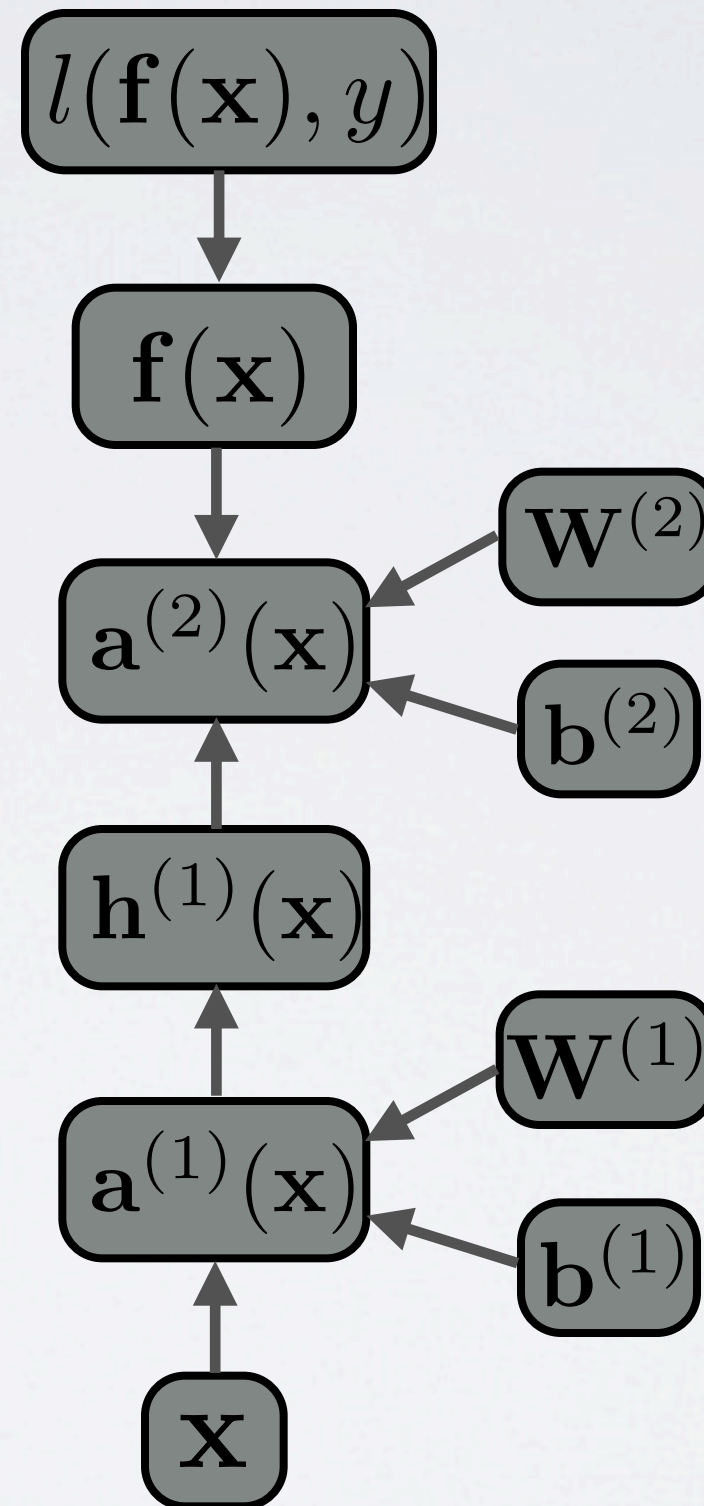
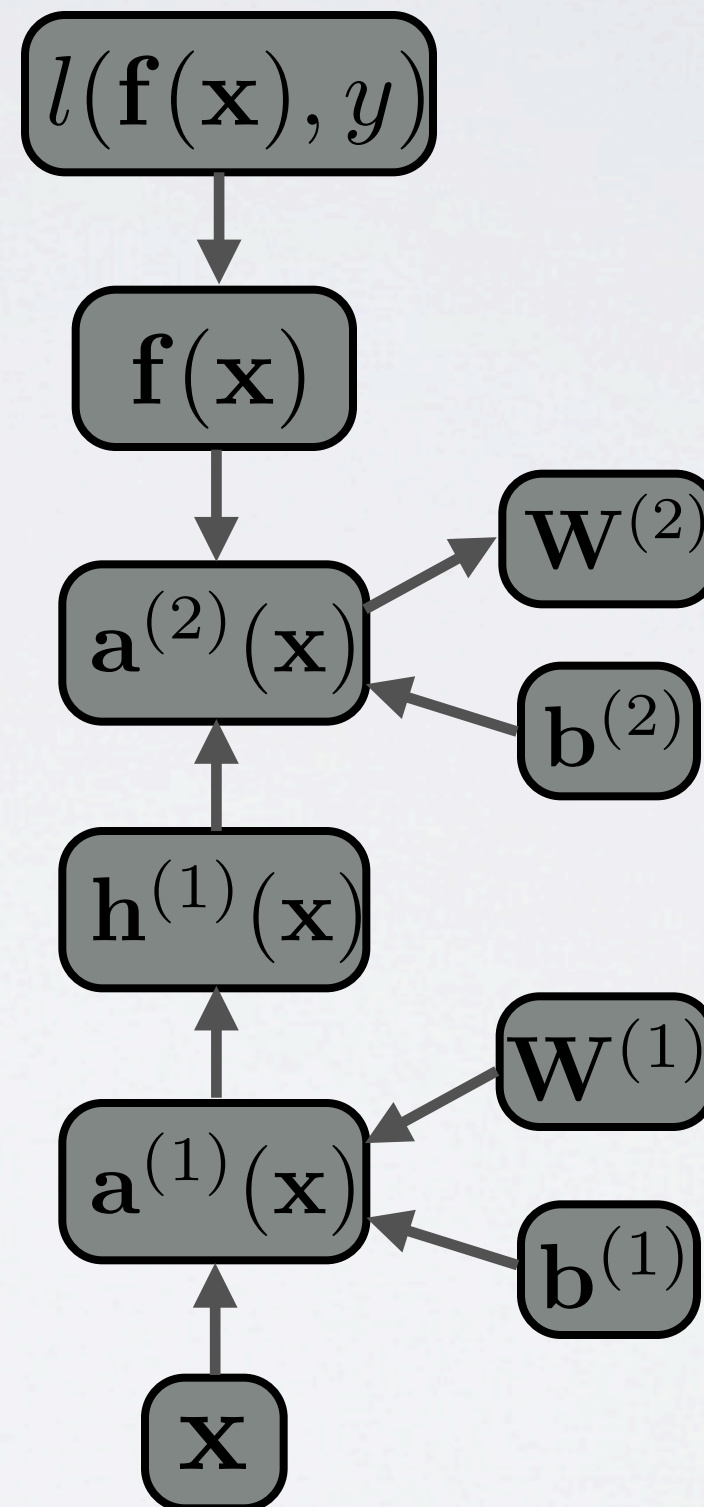**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x}) \quad \mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{a}^{(1)}(\mathbf{x}) \quad \mathbf{W}^{(1)}$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH

**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters
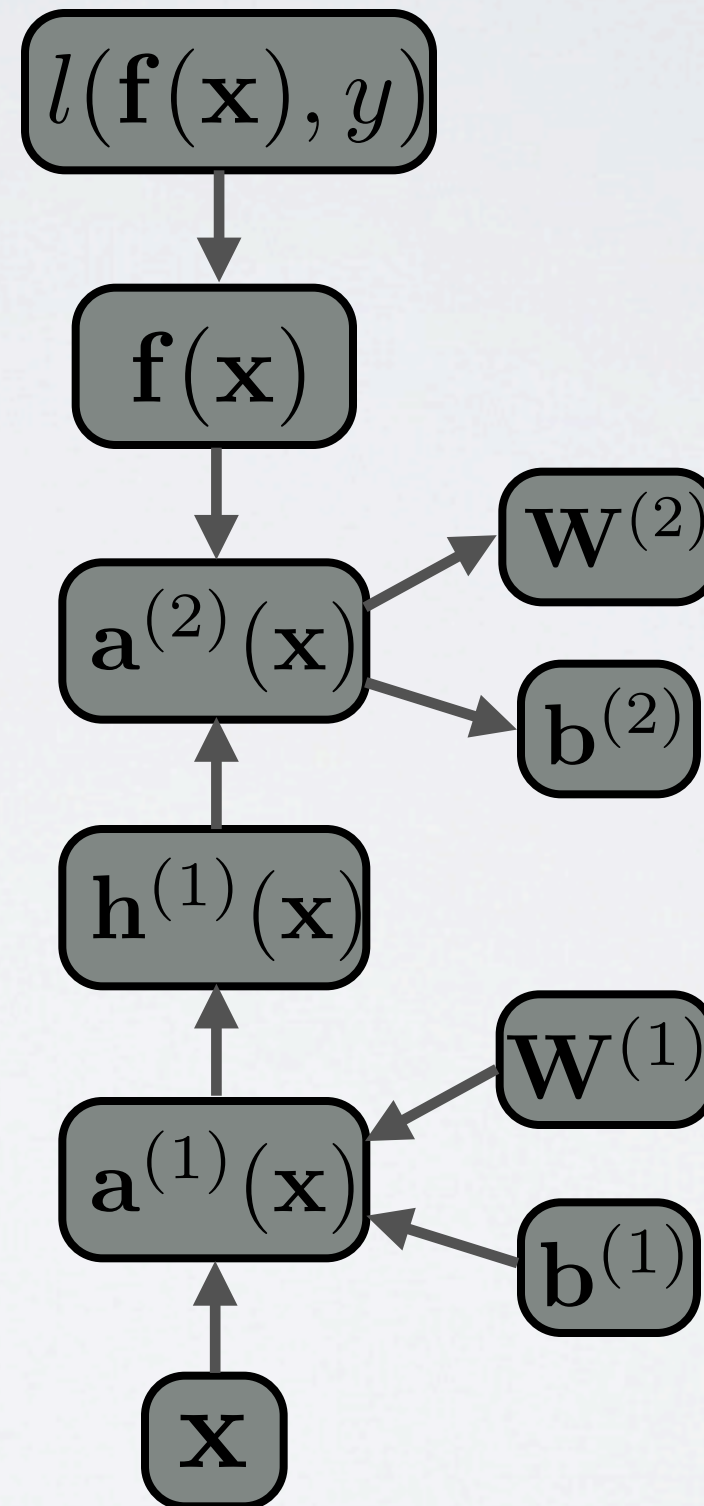
# FLOW GRAPH

**Topics:** automatic differentiation

• Each object also has a bprop method

‣ it computes the gradient of the loss with respect to each children

‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

• By calling bprop in the reverse order, we get backpropagation

‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x})$$

$$\mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{a}^{(1)}(\mathbf{x})$$

$$\mathbf{W}^{(1)}$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

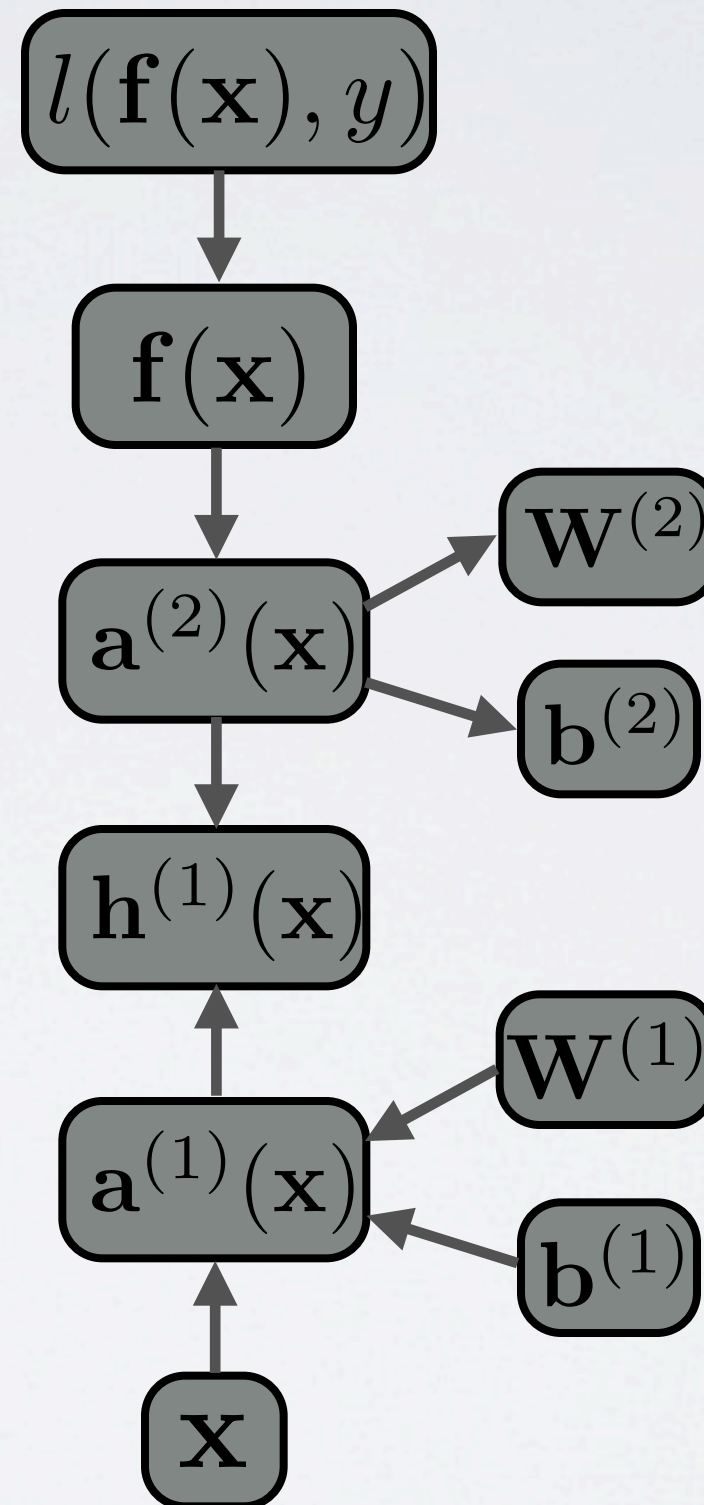# FLOW GRAPH

**Topics:** automatic differentiation

• Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

• By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x}) \qquad \mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{W}^{(1)}$$

$$\mathbf{a}^{(1)}(\mathbf{x})$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH
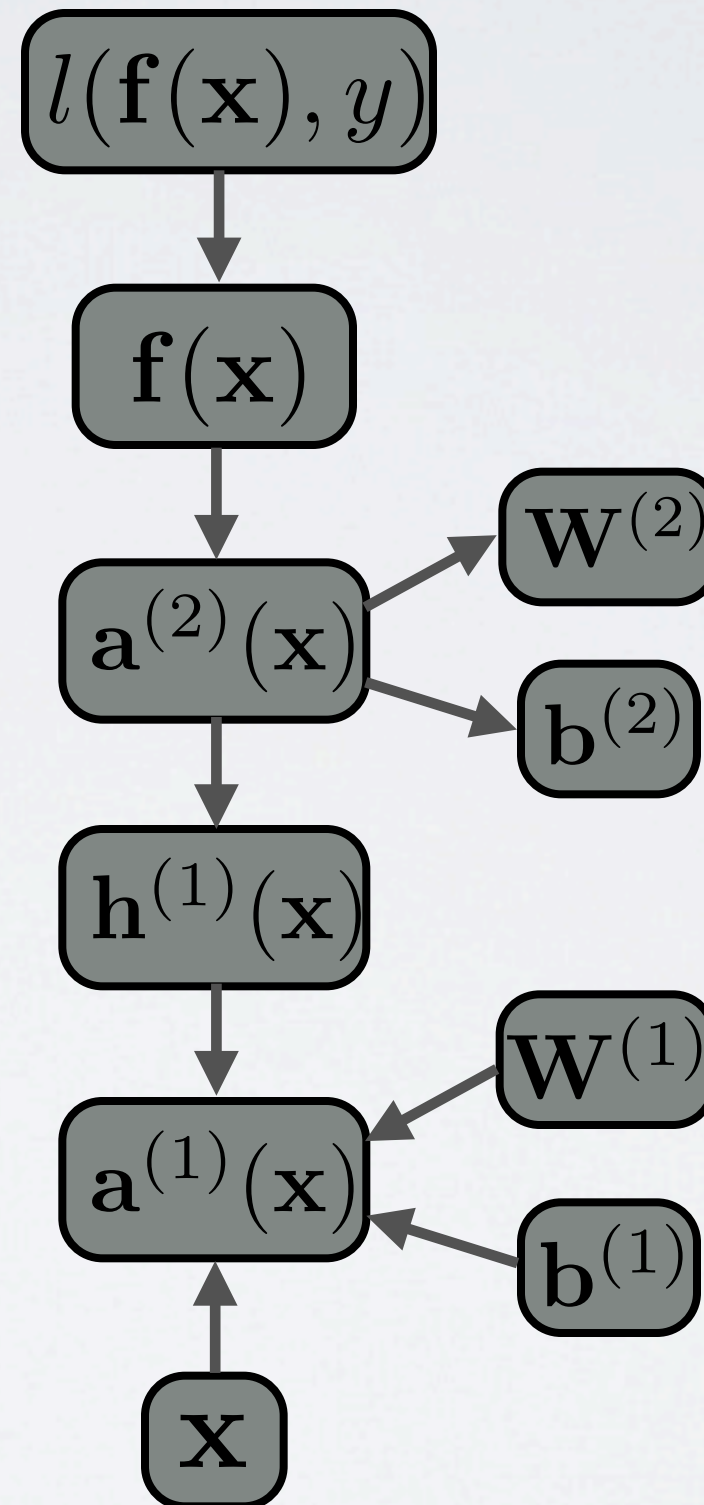
**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x}) \qquad \mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{W}^{(1)}$$

$$\mathbf{a}^{(1)}(\mathbf{x})$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH
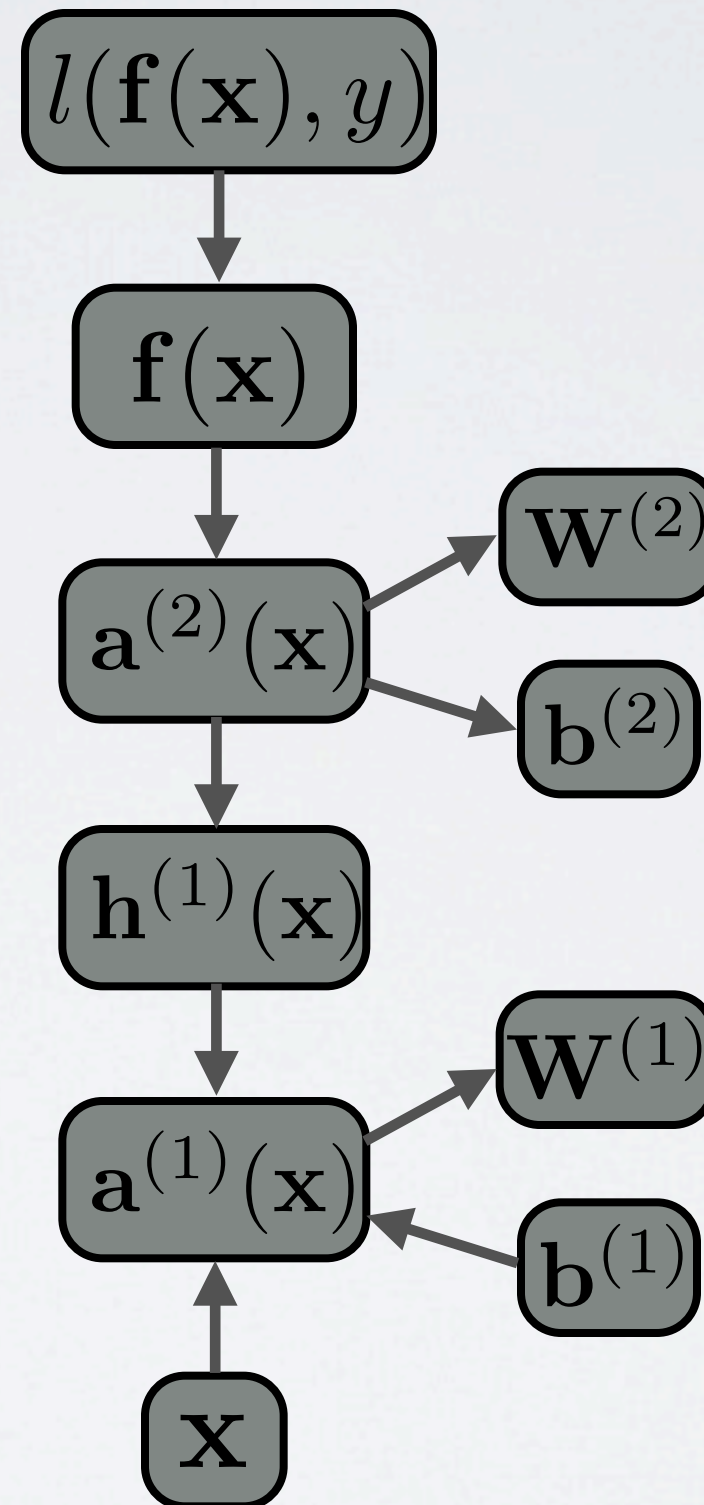
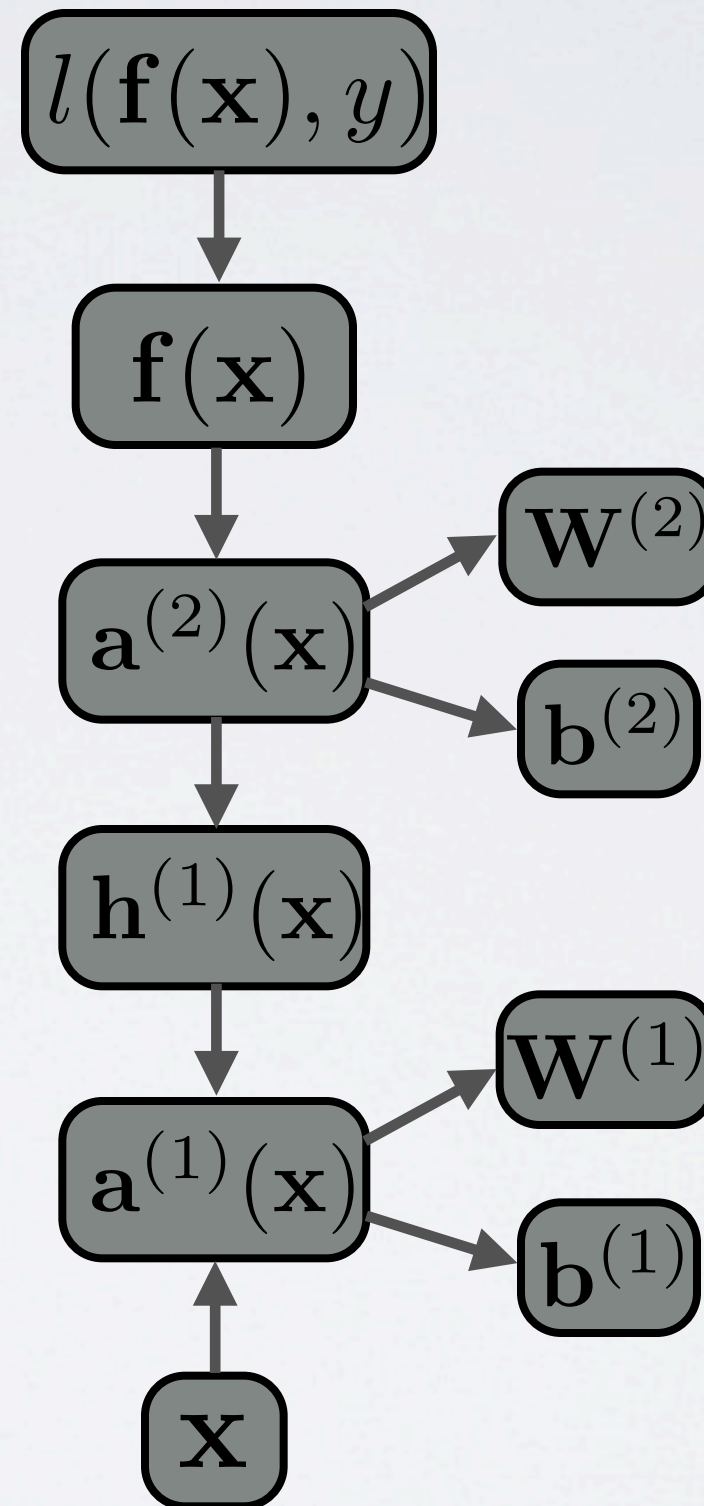**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x})$$

$$\mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{a}^{(1)}(\mathbf{x})$$

$$\mathbf{W}^{(1)}$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# FLOW GRAPH

**Topics:** automatic differentiation

- Each object also has a bprop method

  ‣ it computes the gradient of the loss with respect to each children

  ‣ fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents

- By calling bprop in the reverse order, we get backpropagation

  ‣ only need to reach the parameters

$$l(\mathbf{f}(\mathbf{x}), y)$$

$$\mathbf{f}(\mathbf{x})$$

$$\mathbf{a}^{(2)}(\mathbf{x}) \qquad \mathbf{W}^{(2)}$$

$$\mathbf{b}^{(2)}$$

$$\mathbf{h}^{(1)}(\mathbf{x})$$

$$\mathbf{a}^{(1)}(\mathbf{x}) \qquad \mathbf{W}^{(1)}$$

$$\mathbf{b}^{(1)}$$

$$\mathbf{x}$$

# GRADIENT CHECKING

**Topics:** finite difference approximation

• To debug your implementation of fprop/bprop, you can compare with a finite-difference approximation of the gradient

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$$

‣ $f(x)$ would be the loss

‣ $x$ would be a parameter

‣ $f(x + \epsilon)$ would be the loss if you add $\epsilon$ to the parameter

‣ $f(x - \epsilon)$ would be the loss if you subtract $\epsilon$ to the parameter