

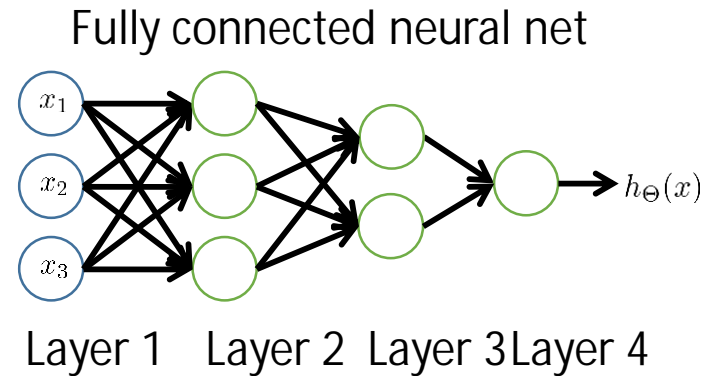
Neural Networks: Learning

Backpropagation

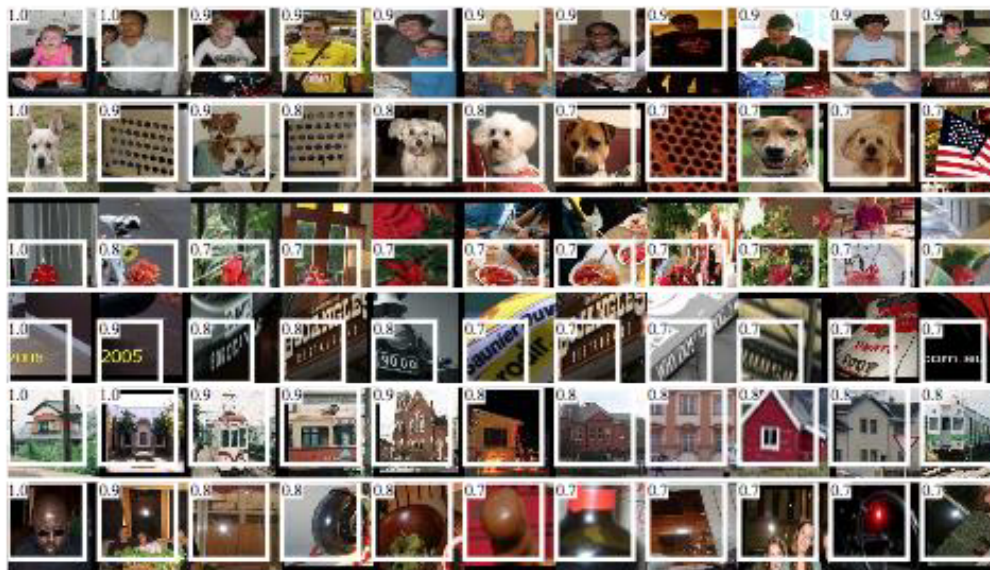
Last time

Neural network chains together many layers of “neurons” such as logistic units

Hidden neurons learn more and more abstract non-linear features



The Unreasonable Effectiveness of Deep Features



Maximal activations of pool_5 units

[R-CNN]

How can we learn parameters of a multilayer network?

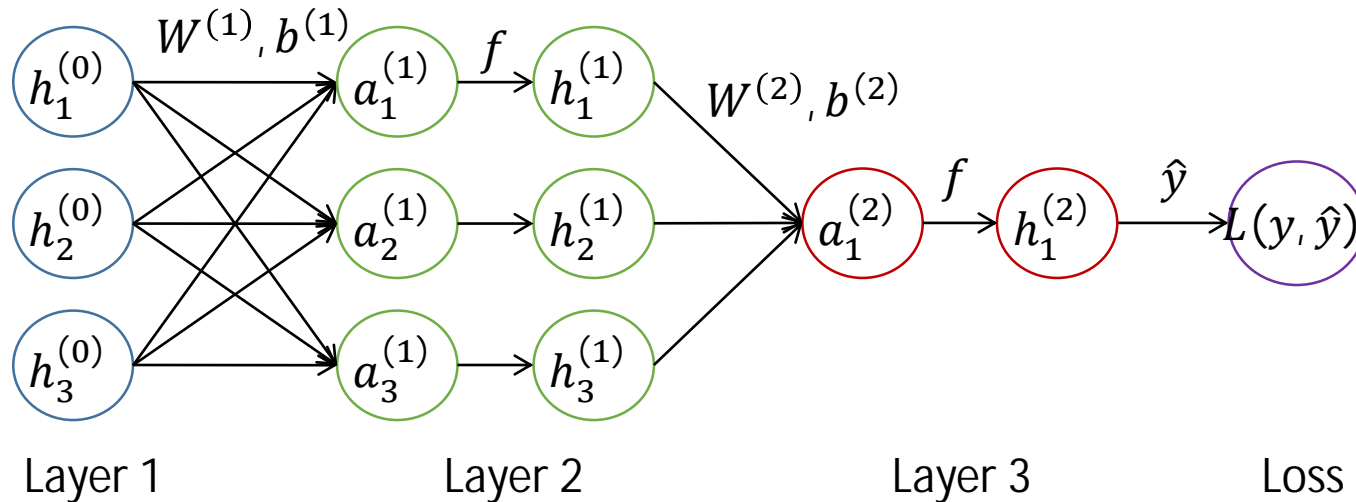
- Use gradient descent
- Compute gradient of loss using chain rule
- Efficient application of chain rule ==
backpropagation

Backpropagation

- [Slides from Stanford course](#)

Backpropagation for a neural net

- Example for a simple 1-hidden layer neural net
- What is the computational graph?



Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} , the input to process

Require: \mathbf{y} , the target output

$\mathbf{h}^{(0)} = \mathbf{x}$

for $k = 1, \dots, l$ do

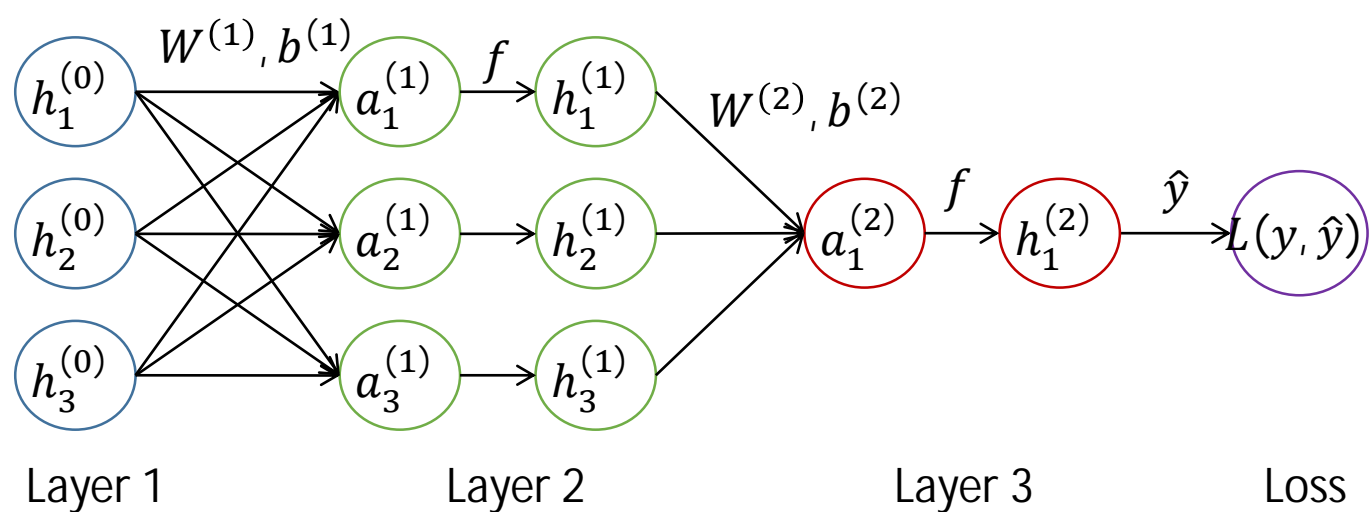
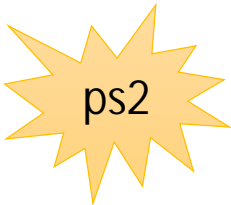
$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$

$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$

end for

$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$

$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$



After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

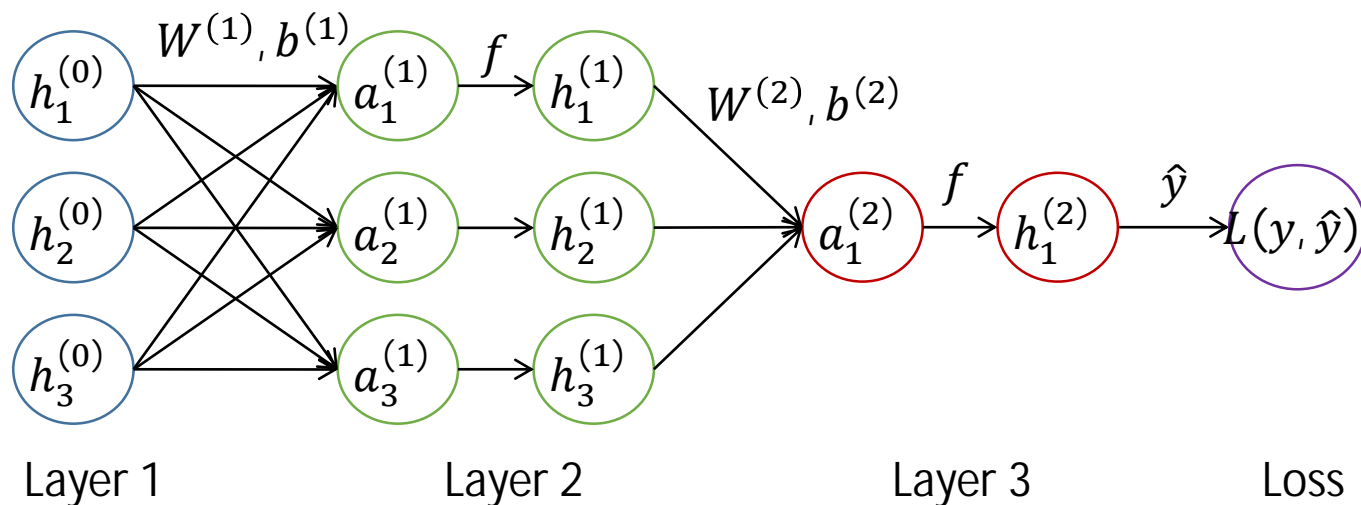
$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for



Artificial Neural Network:

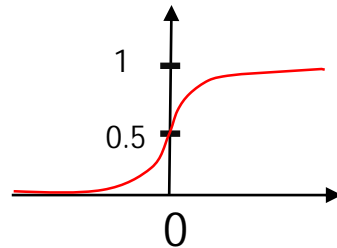
general notation

input $x = \begin{bmatrix} x_1 \\ \dots \\ x_5 \end{bmatrix}$

hidden layer activations

$$h^i = g(\Theta^{(i)} x)$$

$$g(z) = \frac{1}{1 + \exp(-z)}$$



output

$$h_{\Theta}(x) = g(\Theta^{(2)} a)$$

weights

$$\Theta^{(1)} = \begin{pmatrix} \theta_{11} & \dots & \theta_{15} \\ \vdots & \ddots & \vdots \\ \theta_{31} & \dots & \theta_{35} \end{pmatrix}$$

$$\Theta^{(2)} = \begin{pmatrix} \theta_{11} & \dots & \theta_{13} \\ \vdots & \ddots & \vdots \\ \theta_{31} & \dots & \theta_{33} \end{pmatrix}$$

