

Intro Neural Networks

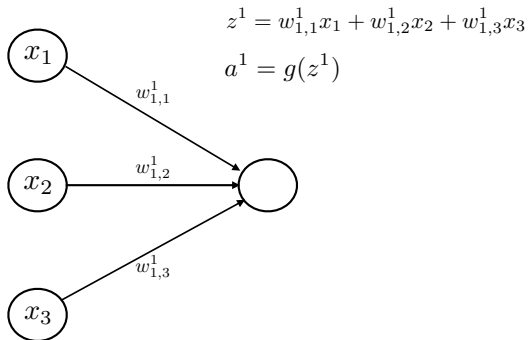
Matt Olson

September 8, 2017

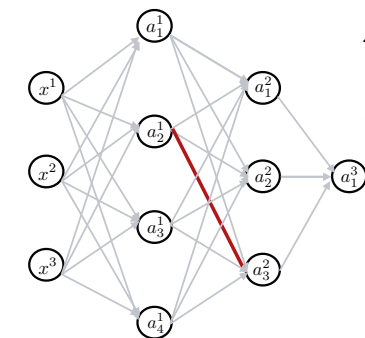
Overview

- Q: What is a Deep Neural Network?
 - ▶ A bunch of logistic regressions jammed together
- Training has recently become feasible: Hinton (2006), GPUs
- We will cover: notation, back prop, examples
- Outline / discussion study group organization

The Perceptron (Logistic Regression)



Larger Networks

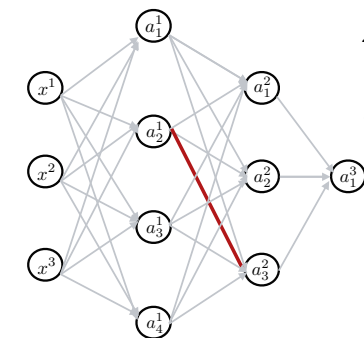


$$z_j^\ell = \sum_k w_{j,k}^l a_k^{l-1} \quad \text{for } l = 1, \dots, 3$$
$$a_j^\ell = g(z_j^\ell)$$

Layer 0 Layer 1 Layer 2 Layer 3

- $w_{j,k}^\ell$ connects node k in layer $\ell - 1$ to node j in layer ℓ .
- The red edge is $w_{3,2}^2$.

Larger Networks Ctd.



$$z_j^\ell = \sum_k w_{j,k}^l a_k^{l-1} \quad \text{for } l = 1, \dots, 3$$
$$a_j^\ell = g(z_j^\ell)$$

Layer 0 Layer 1 Layer 2 Layer 3

- (vectorized) $z^\ell = W_\ell a^\ell$ and $a^\ell = g(z^\ell)$
- Prediction at a point x : $\hat{y}(x) = g(W_3 g(W_2 g(W_1 x)))$

The Learning Problem

Observe $\{(x_1, y_1), \dots, (x_n, y_n)\}$ with $y_i \in \{0, 1\}$. Fix a network architecture with L hidden layers. Find weights W_1, \dots, W_L that minimize

$$\underset{W_1, \dots, W_L}{\text{minimize}} \quad 1/n \sum_{i=1}^n \mathcal{L}(y_i, \hat{y}(x_i))$$

where \mathcal{L} is a loss function (say, cross-entropy) and

$$\hat{y}(x_i) = g(W_L g(W_{L-1} g(\dots g(W_1 x))))$$

- This minimization problem is attacked through gradient descent
- Need to compute $\nabla_w \mathcal{L}(y_i, \hat{y}(x_i))$

The Chain Rule

Recall that for functions h_1, \dots, h_L (of conformable dimensions)

$$f = h_L \circ h_{L-1} \circ \dots \circ h_1$$

$$Df = Dh_L \circ Dh_{L-1} \circ \dots \circ Dh_1$$

$$\nabla f = (Dh_1)^T \circ (Dh_2)^T \circ \dots \circ (Dh_L)^T$$

Gradient Calculations

For simplicity, let's return to the network with two hidden layers from a previous slide and take $\mathcal{L}(y, \theta) = 1/2 (y - \theta)^2$.

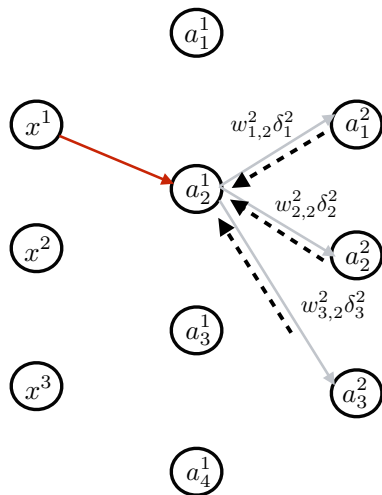
- Feed x_i through the network to get activations a^1, a^2, a^3 and prediction $\hat{y}(x_i)$
- Form $G(\ell) = \text{diag}([g'(z_1^\ell), \dots, g'(z_{L_\ell}^\ell)])$ for $\ell = 1, 2, 3$
- “Trick”: $\frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell} = \frac{\partial \mathcal{L}}{\partial z_j^\ell} \frac{\partial z_j^\ell}{\partial w_{j,k}^\ell}$
 - ▶ $\frac{\partial \mathcal{L}}{\partial z^3} = G(3)(y_i - \hat{y}(x_i))$
 - ▶ $\frac{\partial \mathcal{L}}{\partial z^2} = G(2)W_3^T G(3)(y_i - \hat{y}(x_i))$
 - ▶ $\frac{\partial \mathcal{L}}{\partial z^1} = G(1)W_2^T G(2)W_3^T G(3)(y_i - \hat{y}(x_i))$

Backpropagation

An efficient algorithm for computing $\frac{\partial \mathcal{L}}{\partial w_{j,k}^l}$ at a point x . For a neural network with $L + 1$ layers:

- (1) Feed x through the network and compute a^1, \dots, a^L and prediction $\hat{y}(x)$.
- (2) $G(\ell) = \text{diag}([g'(z_1^\ell), \dots, g'(z_{L_\ell}^\ell)])$ for $\ell = 1, \dots, L$.
- (3) $\delta^L = G(L)(y - \hat{y}(x))$
- (4) $\delta^\ell = G(\ell)W_{\ell+1}^T \delta^{\ell+1}$ for $\ell = 1, \dots, L - 1$.
- (5) $\frac{\partial \mathcal{L}}{\partial w_{j,k}^\ell} = a_k^{\ell-1} \delta_j^\ell$

Zooming In



$$\delta_2^1 = (w_{1,2}^2 \delta_1^2 + w_{2,2}^2 \delta_2^2 + w_{3,2}^2 \delta_3^2) g'(z_2^1)$$

$$\frac{\partial \mathcal{L}}{\partial w_{2,1}^1} = \delta_2^1 x^1$$