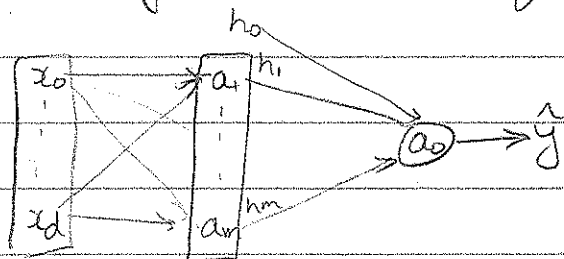


4/16/15

①

# NEURAL NETS : BACKPROP

① Forward Propagation or "Forward Prop" or Forward Pass  
→ Making predictions from a network



Matrix Mult { (a)  $\vec{a} = \begin{bmatrix} \vec{w}_1^T \vec{x} \\ \vdots \\ \vec{w}_n^T \vec{x} \end{bmatrix} = W_n \vec{x}$       $W_n = \begin{bmatrix} \leftarrow w_{11} \rightarrow \\ \vdots \\ \leftarrow w_{n1} \rightarrow \end{bmatrix}$

Non-linearity { (b)  $\vec{h} = g(\vec{a}) = \begin{bmatrix} h_0 = 1 \\ g(a_1) \\ \vdots \\ g(a_n) \end{bmatrix}$

Matrix Multiplication { (c)  $a_0 = \vec{w}_0^T \vec{h}$

Non-linearity { (d)  $\hat{y} = g(a_0)$

## (2) Parameter learning in Neural Nets

[Assume fixed structure/architecture: 3-layer NN]

Weights to learn:  $\left. \begin{array}{l} \{\vec{w}_1, \dots, \vec{w}_m\} \in \mathbb{R}^{d+1} \\ \vec{w}_0 \in \mathbb{R}^{m+1} \end{array} \right\} \vec{\theta}$

Dataset:  $\{(\vec{x}_1, y_1) \dots (\vec{x}_n, y_n)\}$

Objective / Loss  $L(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{Loss}(y_i, \hat{y}(\theta)) + \text{Reg}(\vec{\theta})$

many choices based on problem types

Say Regression

$$L(\vec{\theta}) = \underbrace{\frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y_i - \hat{y}(\theta))^2}_{L_2 \text{-loss "Least Squares"}} + \underbrace{\lambda \|\theta\|_2^2}_{L_2 \text{-norm}}$$

Say Classification

$$L(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N -\log P(y_i | \vec{x}_i, \vec{\theta}) + \lambda \|\theta\|_2^2$$

maximizing log-likelihood

$$-\log \left[ \frac{1}{1 + e^{-w_0 \vec{x}}} \right] \quad (\text{if } y_i = 1)$$

Algorithm: Gradient Descent:

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \eta \frac{\partial L}{\partial \vec{\theta}}$$

How do we compute gradient?  
Backprop!

But first, let's study a simpler case:

2.5 Detour: Assume no hidden units  
 $g(a) = a$  [Linear Regression]  
 $\hat{y} = \vec{w}_0^T \vec{x}$

$$L(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y_i - \vec{w}_0^T \vec{x}_i)^2 \quad \text{[Assume no Reg for simplicity]}$$

$$\frac{\partial L}{\partial \vec{w}_0} = \frac{1}{N} \sum_{i=1}^N \underbrace{(y_i - \vec{w}_0^T \vec{x}_i)}_{\text{error}} \cdot \underbrace{\vec{x}_i^T}_{\text{pattern/feature}}$$

Batch-Gradient Descent

$$\vec{w}_0^{(t+1)} = \vec{w}_0^{(t)} + \eta \frac{1}{N} \sum_{i=1}^N (y_i - \vec{w}_0^T \vec{x}_i) \cdot \vec{x}_i$$

Stochastic Gradient Descent

→ Sample  $(\vec{x}, y)$

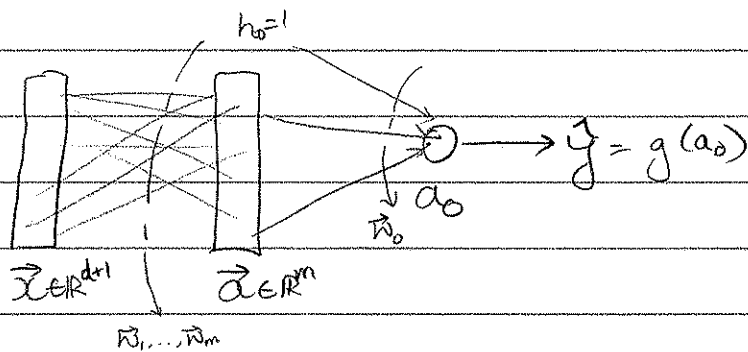
$$\vec{w}_0^{(t+1)} = \vec{w}_0^{(t)} + \eta (y - \vec{w}_0^T \vec{x}) \cdot \vec{x}$$

→ Similar expression as MLE in LR!

### ③ Backward Propagation or BackProp

→ algorithm for computing gradients in NN with Chain Rule

BackProp  $\equiv$  Chain Rule + Caching



→ Assume (for simplicity of math/notation)

→ Regression:  $(y - \hat{y})^2$

→ No regularization:

→ single data-point  $(x, y)$

$$\rightarrow L(\Theta) = L(\vec{w}_1, \dots, \vec{w}_m, \vec{w}_0) = \frac{1}{2} (y - \hat{y}(\Theta))^2$$

→ Goal: compute  $\frac{\partial L}{\partial \vec{w}_0}, \left\{ \frac{\partial L}{\partial \vec{w}_i} \right\}$

→ Approach: we will compute gradients backwards w.r.t  $a_0, \vec{w}_0, a_i, \vec{w}_i$

$$1) \frac{\partial L}{\partial a_0} = \frac{\partial}{\partial a_0} \left[ \frac{1}{2} (y - g(a_0))^2 \right] = -(y - g(a_0)) \cdot g'(a_0)$$

$$2) \frac{\partial L}{\partial \vec{w}_0} = \frac{\partial L}{\partial a_0} \cdot \frac{\partial a_0}{\partial \vec{w}_0} = \frac{\partial L}{\partial a_0} \cdot \frac{\partial (\vec{w}_0^T \vec{h})}{\partial \vec{w}_0} = \frac{\partial L}{\partial a_0} \cdot \vec{h}^T$$

[Chain Rule]      already computed above      error      pattern

3)  $\frac{\partial L}{\partial a_i} = \frac{\partial L}{\partial a_0} \cdot \frac{\partial a_0}{\partial a_i}$  [Chain Rule]

already computed

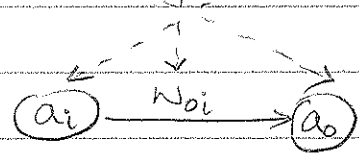
$$a_0 = \vec{w}_0^T h$$

$$= \sum_{l=0}^m w_{0l} h_l$$

$$= \sum_{i=0}^m w_{0i} g(a_i)$$

$$\frac{\partial a_0}{\partial a_i} = w_{0i} \cdot g'(a_i)$$

$\Rightarrow \frac{\partial L}{\partial a_i} = \frac{\partial L}{\partial a_0} \cdot w_{0i} \cdot g'(a_i)$



4)  $\frac{\partial L}{\partial \vec{w}_i} = \frac{\partial L}{\partial a_i} \cdot \frac{\partial a_i}{\partial \vec{w}_i}$

already computed

$$a_i = \vec{w}_i^T X$$

$$\frac{\partial a_i}{\partial \vec{w}_i} = \vec{X}^T$$

$\Rightarrow \frac{\partial L}{\partial \vec{w}_i} = \frac{\partial L}{\partial a_i} \cdot \vec{X}^T$

error pattern

