

Machine Learning for Intelligent Systems

Lecture 14: Backpropagation in Networks

Reading: UML 20.6

The slides are altered from class to use column matrices throughout.

Instructors: Nika Haghtalab (this time) and Thorsten Joachims

1

Multi Layer Neural Network

Input layer: $V_0 = \vec{x}$

Hidden layers: $V_1 = \vec{v}_1, V_2 = \vec{v}_2, \dots, V_{d-1} = \vec{v}_{d-1}$

Output: V_d

Depth: d

Width: k

Vector of weights going from layer i to the j^{th} node of layer $i + 1$: $\vec{w}_{i,j}$

Vector of values in layer $i + 1$, $\vec{v}_{i+1} = \begin{pmatrix} \sigma(\vec{v}_i \cdot \vec{w}_{i,1}) \\ \vdots \\ \sigma(\vec{v}_i \cdot \vec{w}_{i,k-1}) \end{pmatrix}$

Output: $v_{d,1}$

2

Common Activation Functions

Use a non-linear activation function on nodes of a hidden layer.

Name	Function	Gradient	Graph
Binary step	$\text{sign}(x)$	$\begin{cases} 0 & x \neq 0 \\ N/A & x = 0 \end{cases}$	
sigmoid	$\sigma(x) = \frac{1}{1 + \exp(-x)}$	$\sigma(x)(1 - \sigma(x))$	
Tanh	$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$	$1 - \tanh^2(x)$	
Rectified Linear (ReLU)	$\text{relu}(x) = \max(x, 0)$	$\begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$	

Sometimes, $\sigma(x)$ denotes the "generic" notion of activation function, not necessarily sigmoid.

3

Closer look at the nodes

$\vec{w}_{i,j}$: vector of weight from layer i to j^{th} node in the next layer.

Weighted sum of values from the previous layer: $s_{i+1,j} = \vec{w}_{i,j} \cdot \vec{v}_i$

The value of the node: $v_{i+1,j} = \sigma(s_{i+1,j})$

Vectorized form: $\vec{v}_{i+1} = \begin{pmatrix} \sigma(\vec{v}_i \cdot \vec{w}_{i,1}) \\ \vdots \\ \sigma(\vec{v}_i \cdot \vec{w}_{i,k-1}) \end{pmatrix}$

Bias term, has no incoming weights, so it's a fixed constant.

4

Simpler form

Vectorized form: $\vec{v}_{i+1} = \begin{pmatrix} \sigma(\vec{v}_i \cdot \vec{w}_{i,1}) \\ \vdots \\ \sigma(\vec{v}_i \cdot \vec{w}_{i,k-1}) \end{pmatrix}$

Bias term, has no incoming weights.

Presentation Trick:
 → When σ is sigmoid, $\sigma(\mathbf{0}) = 1/2$
 → We can assume that $\vec{w}_{i,k} = 0$
 → The network is fully connected.
 → $\vec{v}_{i+1} = \sigma(W_i^T \vec{v}_i)$

5

Prediction with Neural Networks

How do we use Neural Networks for predictions?

$$\sigma \left(W_{d-1}^T \dots \sigma \left(W_2^T \sigma \left(W_1^T \sigma \left(W_0^T \vec{v}_0 \right) \right) \dots \right) \right)$$

Forward Pass (Prediction)

Input: Neural Network with weight matrices W_0, W_1, \dots, W_{d-1} and instance (\vec{x}, y)

$\vec{v}_0 = \vec{x}$

For $\ell = 1, \dots, d$

- $\vec{s}_\ell = W_{\ell-1}^T \vec{v}_{\ell-1}$
- $\vec{v}_\ell = \sigma(\vec{s}_\ell)$

End For

Output \vec{v}_d

6

Learning the Weights

How can we learn weight matrices W_0, W_1, \dots, W_{d-1} given a sample set $S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)\}$.

Need to write the optimization:

$$\min_{W_0, \dots, W_{d-1}} \frac{1}{m} \sum_{i=1 \dots m} L(\text{forward pass}_{W_0, \dots, W_{d-1}}(\vec{x}_i), y_i)$$

For convenience, we don't use a regularizer in this lecture.

$L(y', y_i)$: Evaluates how good is y' as a prediction for y_i . We use $L(y', y_i) = \frac{1}{2}(y' - y_i)^2$

7

SGD for Neural Networks

How can we learn weight matrices W_0, W_1, \dots, W_{d-1} given a sample set $S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)\}$.

Stochastic Gradient Descent updates rule:

$$\min_{W_0, \dots, W_{d-1}} \frac{1}{m} \sum_{i=1 \dots m} L(\text{forward pass}_{W_0, \dots, W_{d-1}}(x), y_i)$$

Take a random (\vec{x}, y) from the samples, let ∂L denote the partial derivate for this sample with respect to the current weights

$$\text{For all } i = 0, \dots, d-1, \quad W_i \leftarrow W_i - \eta \frac{\partial L}{\partial W_i}$$

8

Derivative of the network

We need to compute the derivatives $\frac{\partial L}{\partial W_i}$

Chain rule in derivatives:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g(x)} \cdot \frac{\partial g(x)}{\partial h(x)} \cdot \frac{\partial h(x)}{\partial x}$$

(Element wise) Application of chain rule:

$$\begin{aligned} \frac{\partial L}{\partial w_{i,j}(z)} &= \frac{\partial L}{\partial v_{i+1,j}} \cdot \frac{\partial v_{i+1,j}}{\partial s_{i+1,j}} \cdot \frac{\partial s_{i+1,j}}{\partial w_{i,j}(z)} \\ &= \frac{\partial L}{\partial v_{i+1,j}} \sigma'(s_{i+1,j}) v_{i,z} \end{aligned}$$

9

Derivative wrt Node Values

We need to compute the derivatives $\frac{\partial L}{\partial v_{i,j}}$

For v_d , it's just the derivate of the loss with respect to the actual label y . For square loss:

$$\frac{\partial L}{\partial v_d} = v_d - y$$

Application of chain rule:

$$\begin{aligned} \frac{\partial L}{\partial v_{d-1,1}} &= \frac{\partial L}{\partial v_d} \cdot \frac{\partial v_d}{\partial s_d} \cdot \frac{\partial s_d}{\partial v_{d-1,1}} \\ &= (v_d - y) \cdot \sigma'(s_d) \cdot w_{d-1,1} \quad (1) \end{aligned}$$

10

Derivative wrt Node Values

More generally, recursively

$$\begin{aligned} \frac{\partial L}{\partial \vec{v}_i} &= \frac{\partial L}{\partial \vec{v}_{i+1}} \cdot \frac{\partial \vec{v}_{i+1}}{\partial \vec{s}_{i+1}} \cdot \frac{\partial \vec{s}_{i+1}}{\partial \vec{v}_i} \\ &= W_i \left(\sigma'(\vec{s}_{i+1}) \odot \frac{\partial L}{\partial \vec{v}_{i+1}} \right) \end{aligned}$$

Element-wise matrix multiplication: $\begin{pmatrix} a \\ b \end{pmatrix} \odot \begin{pmatrix} c \\ d \end{pmatrix} = \begin{pmatrix} ac \\ bd \end{pmatrix}$

11

All together

We can compute all the derivatives $\frac{\partial L}{\partial \vec{v}_i}$ before each round of SGD.

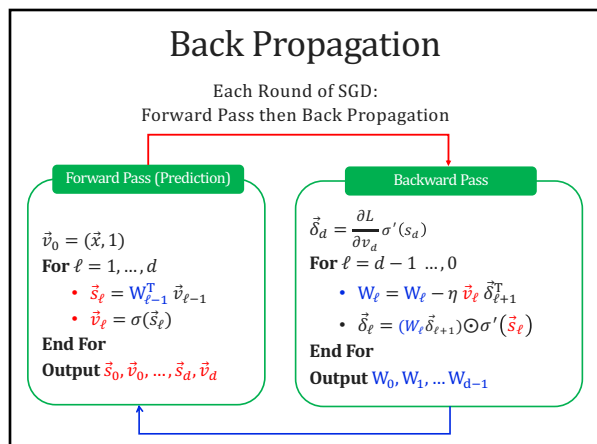
$$\frac{\partial L}{\partial v_d} = v_d - y$$

$$\frac{\partial L}{\partial \vec{v}_i} = W_i \left(\frac{\partial L}{\partial \vec{v}_{i+1}} \odot \sigma'(\vec{s}_{i+1}) \right)$$

When we start, we now have the derivatives with respect to weights.

$$\frac{\partial L}{\partial W_i} = \vec{v}_i \left(\frac{\partial L}{\partial \vec{v}_{i+1}} \odot \sigma'(\vec{s}_{i+1}) \right)^T$$

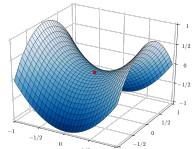
12



13

SGD on Non-Convex

The loss in neural networks is non-convex.



Many of the theoretical guarantees of SGD only hold for convex losses.

1. Initialization.
 - Don't start $W_0, W_1, \dots, W_{d-1} = 0$
 - Randomized your starting weights.
2. Run separate SGDs and take the best.

14