

Automatic Speech Recognition (CS753)

Lecture 9: Brief Introduction to Neural Networks

Instructor: Preethi Jyothi Feb 2, 2017

Final Project Landscape



Feed-forward Neural Network



Feed-forward Neural Network Brain Metaphor



Image from: <u>https://upload.wikimedia.org/wikipedia/commons/1/10/Blausen_0657_MultipolarNeuron.png</u>

Feed-forward Neural Network Parameterized Model



If \mathbf{x} is a 2-dimensional vector and the layer above it is a 2-dimensional vector \mathbf{h} , a fully-connected layer is associated with:

$$\mathbf{h} = \mathbf{x}\mathbf{W} + \mathbf{b}$$

where w_{ij} in **W** is the weight of the connection between ith neuron in the input row and jth neuron in the first hidden layer and **b** is the bias vector

Feed-forward Neural Network Parameterized Model



$$a_{5} = g(w_{35} \cdot a_{3} + w_{45} \cdot a_{4})$$

= $g(w_{35} \cdot (g(w_{13} \cdot a_{1} + w_{23} \cdot a_{2})) + w_{45} \cdot (g(w_{14} \cdot a_{1} + w_{24} \cdot a_{2})))$

The simplest neural network is the perceptron:

$$Perceptron(x) = xW + b$$

A 1-layer feedforward neural network has the form: $MLP(x) = g(xW_1 + b_1)W_2 + b_2$

Common Activation Functions (g)

Sigmoid: $\sigma(x) = 1/(1 + e^{-x})$



Common Activation Functions (g)

Sigmoid: $\sigma(x) = 1/(1 + e^{-x})$ Hyperbolic tangent (tanh): $tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$



Common Activation Functions (g)

Sigmoid: $\sigma(x) = 1/(1 + e^{-x})$ **Hyperbolic tangent (tanh)**: $tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$ **Rectified Linear Unit (ReLU):** RELU(x) = max(0, x)



Optimization Problem

- To train a neural network, define a loss function $L(y,\tilde{y})$: a function of the true output y and the predicted output \tilde{y}
- $L(y,\tilde{y})$ assigns a non-negative numerical score to the neural network's output, \tilde{y}
- The parameters of the network are set to minimise L over the training examples (i.e. a sum of losses over different training samples)
- L is typically minimised using a gradient-based method

Stochastic Gradient Descent (SGD)

SGD Algorithm

Inputs:

Function NN(x; θ), Training examples, $x_1 \dots x_n$ and outputs, $y_1 \dots y_n$ and Loss function L.

do until stopping criterion Pick a training example x_i , y_i Compute the loss $L(NN(x_i; \theta), y_i)$ Compute gradient of L, ∇L with respect to θ $\theta \leftarrow \theta - \eta \nabla L$

done

Return: θ

Training a Neural Network

Define the **Loss function** to be minimised as a node *L*

Goal: Learn weights for the neural network which minimise *L*

Gradient Descent: Find $\partial L/\partial w$ for every weight *w*, and update it as $w \leftarrow w - \eta \partial L/\partial w$

How do we efficiently compute $\partial L/\partial w$ for all *w*?

Will compute $\partial L/\partial u$ for every node *u* in the network!

 $\partial L/\partial w = \partial L/\partial u \cdot \partial u/\partial w$ where *u* is the node which uses *w*

Training a Neural Network

New goal: compute $\partial L/\partial u$ for every node u in the network

Simple algorithm: Backpropagation

Key fact: Chain rule of differentiation

If *L* can be written as a function of variables $v_1, ..., v_n$, which in turn depend (partially) on another variable *u*, then

 $\partial L/\partial u = \sum_i \partial L/\partial v_i \cdot \partial v_i/\partial u$

Backpropagation

If *L* can be written as a function of variables $v_1,..., v_n$, which in turn depend (partially) on another variable *u*, then



Then, the chain rule gives

 $\partial L/\partial u = \Sigma_v \in \Gamma(u) \ \partial L/\partial v \cdot \partial v/\partial u$

Backpropagation

```
\partial L/\partial u = \Sigma_v \in \Gamma(u) \ \partial L/\partial v \cdot \partial v/\partial u
```

Backpropagation

Base case: $\partial L / \partial L = 1$

For each *u* (top to bottom):

For each $v \in \Gamma(u)$:

Inductively, have computed $\partial L / \partial v$

Directly compute $\partial v / \partial u$

Compute $\partial L/\partial u$

Compute $\partial L/\partial w$ where $\partial L/\partial w = \partial L/\partial u \cdot \partial u/\partial w$



Forward Pass

First compute all values of *u* given an input, in a forward pass (The values of each node will be needed during backprop)

Where values computed in the forward pass may be needed

Neural Network Acoustic Models

- Input layer takes a window of acoustic feature vectors
- Output layer corresponds to classes (e.g. monophone labels, triphone states, etc.)



Phone posteriors

Image adapted from: Dahl et al., "Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition", TASL'12

Neural Network Acoustic Models

- Input layer takes a window of acoustic feature vectors
- Hybrid NN/HMM systems: replace GMMs with outputs of NNs



Image from: Dahl et al., "Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition", TASL'12