

Artificial Intelligence Programming

More Neural Networks

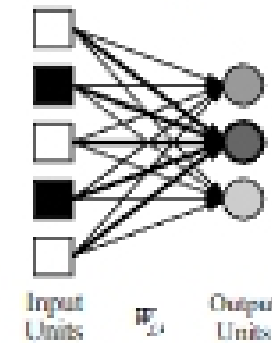
Chris Brooks

Department of Computer Science
University of San Francisco

Neural networks - refresher

- A network is composed of layers of nodes
 - input, hidden, output layers in feedforward nets
- An input is applied to the input units
- The resulting output is the value of the function the net computes.

Perceptrons - refresher



- Single-layer networks (perceptrons) are the simplest form of NN.
- Easy to understand, but computationally limited.
- Each input unit is directly connected to one or more output units.

Department of Computer Science, University of San Francisco, © 2011

Department of Computer Science, University of San Francisco, © 2011

Delta rule - refresher

- The appeal of perceptrons is the ability to automatically learn their weights in a supervised fashion.
- The weight updating rule is known as the **delta rule**.
- $\Delta w_{ij} = \alpha \sum_{d \in D} (t_d - o_d) x_{ij}^d$
- Where D is the training set, t_d is expected output and o_d is actual output.

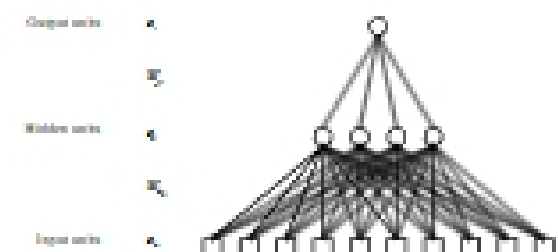
Department of Computer Science, University of San Francisco, © 2011

Multilayer Networks

- While perceptrons have the advantage of a simple learning algorithm, their computational limitations are a problem.
- What if we add another "hidden" layer?
- Computational power increases
 - With one hidden layer, can represent any continuous function
 - With two hidden layers, can represent any function
- Problem: How to find the correct weights for hidden nodes?

Department of Computer Science, University of San Francisco, © 2011

Multilayer Network Example



Department of Computer Science, University of San Francisco, © 2011

Backpropagation

- Backpropagation is an extension of the perceptron learning algorithm to deal with multiple layers of nodes.
- Nodes use sigmoid activation function, rather than the step function
 - $g(\text{input}_i) = \frac{1}{1 + e^{-\text{input}_i}}$.
 - $g'(\text{input}_i) = g(\text{input}_i)(1 - g(\text{input}_i))$ (good news here - to compute g' , we just need g)
- We will still "follow the gradient", where g' gives us the gradient.

Downloaded from <https://www.coursera.org/lecture/ai-introduction/backpropagation-1>

More on computing error

- Recall that our goal is to minimize the error function.
- To do this, we want to change the weights so as to reduce the error.
- We define error as a function of the weights like so:
 - $E(\mathbf{w}) = \text{expected} - \text{actual}$
 - $E(\mathbf{w}) = \text{expected} - g(\text{input})$
- So, to determine how to change the weights, we compute the derivative of error with respect to the weights.
 - This tells us the slope of the error curve at that point.

Downloaded from <https://www.coursera.org/lecture/ai-introduction/backpropagation-2>

More on computing error

- $E(\mathbf{w}) = \text{expected} - g(\text{input})$
- $E(\mathbf{w}) = \text{expected} - g(\mathbf{w} \cdot \mathbf{i})$
- $\frac{\partial E}{\partial \mathbf{w}} = 0 - g'(\text{input})\mathbf{i}$
- $\delta_{\mathbf{w}} = \frac{\partial E}{\partial \mathbf{w}} = -g'(\text{input}) \cdot (1 - g(\text{input})) \cdot \mathbf{i}$

Downloaded from <https://www.coursera.org/lecture/ai-introduction/backpropagation-3>

Updating hidden-output weights

- Each weight is updated by $\alpha \cdot \Delta_k$
- $W_{j,k} = W_{j,k} + \alpha \cdot a_j \cdot \Delta_k$

Downloaded from <https://www.coursera.org/lecture/ai-introduction/backpropagation-4>

Backpropagation

- Updating input-hidden weights:
- Idea: each hidden node is responsible for a fraction of the error in δ_k .
- Divide δ_k according to the strength of the connection between the hidden and output node.
- For each hidden node j
- $\delta_j = g'(\text{input}_j)(1 - g(\text{input}_j)) \sum_{k \in \text{output}} W_{j,k} \delta_k$
- Update rule for input-hidden weights:
- $W_{k,j} = W_{k,j} + \alpha \cdot \text{input}_k \cdot \delta_j$

Downloaded from <https://www.coursera.org/lecture/ai-introduction/backpropagation-5>

Backpropagation Algorithm

- The whole algorithm can be summed up as:
 - While not done:
 - for d in training set
 - Apply inputs of d , propagate forward.
 - for node i in output layer
 - $\delta_i = \text{output}_i + (1 - \text{output}_i) \cdot (t_{\text{exp}} - \text{output}_i)$
 - for each hidden node j
 - $\delta_j = \text{output}_j + (1 - \text{output}_j) \cdot \sum W_{j,i} \delta_i$
 - Adjust each weight
 - $W_{j,i} = W_{j,i} + \alpha \cdot \delta_j \cdot \text{input}_j$

Downloaded from <https://www.coursera.org/lecture/ai-introduction/backpropagation-6>

Theory vs Practice

- In the definition of backpropagation, a single update for all weights is computed for all data points at once.
 - Find the update that minimizes total sum of squared error.
- Guaranteed to converge in this case.
- Problem: This is often computationally space-intensive.
 - Requires creating a matrix with one row for each data point and inverting it.
- In practice, updates are done incrementally instead.

Downloaded from Study Drive on Wednesday, 20 September 2017

Stopping conditions

- Unfortunately, incremental updating is not **guaranteed** to converge.
- Also, convergence can take a long time.
- When to stop training?
 - Fixed number of iterations
 - Total error below a set threshold
 - Convergence - no change in weights

Downloaded from Study Drive on Wednesday, 20 September 2017

Comments on Backpropagation

- Also works for multiple hidden layers
- Backpropagation is only guaranteed to converge to a local minimum
 - May not find the absolute best set of weights
- Low initial weights can help with this
 - Makes the network act more linearly - fewer minima
- Can also use random restart - train multiple times with different initial weights.

Downloaded from Study Drive on Wednesday, 20 September 2017

Momentum

- Since backpropagation is a hillclimbing algorithm, it is susceptible to getting stuck in plateaus
 - Areas where local weight changes don't produce an improvement in the error function.
- A common extension to backpropagation is the addition of a momentum term.
 - Carries the algorithm through minima and plateaus.
- Idea: remember the "direction" you were going in, and by default keep going that way.
- Mathematically, this means using the second derivative.

Downloaded from Study Drive on Wednesday, 20 September 2017

Momentum

- Implementing momentum typically means remembering what update was done in the previous iteration.
- Our update rule becomes:
 - $\Delta w_{ji}(n) = \alpha \Delta_j x_{ji} + \beta \Delta w_{ji}(n-1)$
- To consider the effect, imagine that our new delta is zero (we haven't made any improvement)
- Momentum will keep the weights "moving" in the same direction.
- Also gradually increases step size in areas where gradient is unchanging.
 - This speeds up convergence, helps escape plateaus and local minima.

Downloaded from Study Drive on Wednesday, 20 September 2017

Design issues

- One difficulty with neural nets is determining how to **encode** your problem
 - Inputs must be 1 and 0, or else real-valued numbers.
 - Same for outputs
- Symbolic variables can be given binary encodings
- More complex concepts may require care to represent correctly.

Downloaded from Study Drive on Wednesday, 20 September 2017