## Lecture 10

### Artificial Neural Networks in Data Engineering: Overview

**Wednesday, February 9, 2000**

**William H. Hsu**

**Department of Computing and Information Sciences, KSU**

http://www.cis.ksu.edu/~bhsu

Readings:
Chapter 19, Russell and Norvig

**KSU**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Lecture Outline

- **Read Sections 4.5-4.9, Mitchell; Chapter 4, Bishop; Rumelhart *et al***
- **Multi-Layer Networks**
  - **Nonlinear transfer functions**
  - **Multi-layer networks of nonlinear units (sigmoid, hyperbolic tangent)**
- **Backpropagation of Error**
  - **The backpropagation algorithm**
    - **Relation to error gradient function for nonlinear units**
    - **Derivation of training rule for feedfoward multi-layer networks**
  - **Training issues**
    - **Local optima**
    - **Overfitting in ANNs**
- **Hidden-Layer Representations**
- **Examples: Face Recognition and Text-to-Speech**
- **Advanced Topics (Brief Survey)**
- **Next Week: Chapter 5 and Sections 6.1-6.5, Mitchell; Quinlan paper**

**KSU**

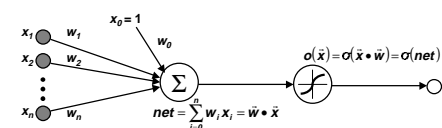CIS 830: Advanced Topics in Artificial Intelligence

---

## Multi-Layer Networks of Nonlinear Units

- **Nonlinear Units**
  - **Recall: activation function $sgn(w \bullet x)$**
  - **Nonlinear activation function: generalization of $sgn$**



- **Multi-Layer Networks**
  - **A specific type: Multi-Layer Perceptrons (MLPs)**
  - **Definition: a multi-layer feedforward network is composed of an input layer, one or more hidden layers, and an output layer**
  - **"Layers": counted in weight layers (e.g., 1 hidden layer ≡ 2-layer network)**
  - **Only hidden and output layers contain perceptrons (threshold or nonlinear units)**
- **MLPs in Theory**
  - **Network (of 2 or more layers) can represent any function (arbitrarily small error)**
  - **Training even 3-unit multi-layer ANNs is $NP$-hard (Blum and Rivest, 1992)**
- **MLPs in Practice**
  - **Finding or *designing* effective networks for arbitrary functions is difficult**
  - **Training is very computation-intensive even when structure is "known"**

**KSU**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Nonlinear Activation Functions



$$o(\vec{x}) = \sigma(\vec{x} \bullet \vec{w}) = \sigma(net)$$

$$net = \sum_{i=0}^{n} w_i x_i = \vec{w} \bullet \vec{x}$$

- **Sigmoid Activation Function**
  - **Linear threshold gate activation function: $sgn(w \bullet x)$**
  - **Nonlinear activation (*aka* transfer, squashing) function: generalization of $sgn$**
  - **$\sigma$ is the sigmoid function** $\sigma(net) = \dfrac{1}{1 + e^{-net}}$
  - **Can derive gradient rules to train**
    - **One sigmoid unit**
    - **Multi-layer, feedforward networks of sigmoid units (using backpropagation)**
- **Hyperbolic Tangent Activation Function** $\sigma(net) = \dfrac{\sinh(net)}{\cosh(net)} = \dfrac{e^{net} - e^{-net}}{e^{net} + e^{-net}}$

**KSU**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Backpropagation Algorithm

- **Intuitive Idea: Distribute *Blame* for Error to Previous Layers**
- **Algorithm *Train-by-Backprop* (D, r)**
  - **Each training example is a pair of the form <x, t(x)>, where x is the vector of input values and t(x) is the output value. r is the learning rate (e.g., 0.05)**
  - **Initialize all weights $w_i$ to (small) random values**
  - **UNTIL the termination condition is met, DO**
    **FOR each <x, t(x)> in D, DO**
    **Input the instance x to the unit and compute the output $o(x) = \sigma(net(x))$**
    **FOR each output unit k, DO**
    $$\delta_k = o_k(x)(1 - o_k(x))(t_k(x) - o_k(x))$$
    **FOR each hidden unit j, DO**
    $$\delta_j = h_j(x)(1 - h_j(x)) \sum_{k \in outputs} v_{jk} \delta_j$$
    **Update each $w = u_{i,j}(a = h_j)$ or $w = v_{j,k}(a = o_k)$**
    $$w_{start\text{-}layer,\, end\text{-}layer} \leftarrow w_{start\text{-}layer,\, end\text{-}layer} + \Delta w_{start\text{-}layer,\, end\text{-}layer}$$
    $$\Delta w_{start\text{-}layer,\, end\text{-}layer} \leftarrow r\, \delta_{end\text{-}layer}\, a_{end\text{-}layer}$$
  - **RETURN final u, v**



**KSU**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Backpropagation and Local Optima

- **Gradient Descent in Backprop**
  - **Performed over entire *network* weight vector**
  - **Easily generalized to arbitrary directed graphs**
  - **Property: Backprop on feedforward ANNs will find a *local* (not necessarily global) error minimum**
- **Backprop in Practice**
  - **Local optimization often works well (can *run multiple times*)**
  - **Often include weight momentum $\alpha$**
    $$\Delta w_{start\text{-}layer,\, end\text{-}layer}(n) = r\, \delta_{end\text{-}layer}\, a_{end\text{-}layer} + \alpha \Delta w_{start\text{-}layer,\, end\text{-}layer}(n-1)$$
  - **Minimizes error over training examples - generalization to subsequent instances?**
  - **Training often *very* slow: thousands of iterations over D (epochs)**
  - **Inference (applying network after training) typically very fast**
    - **Classification**
    - **Control**

**KSU**

CIS 830: Advanced Topics in Artificial Intelligence

## Feedforward ANNs: Representational Power and Bias

- **Representational (i.e., Expressive) Power**
  - Backprop presented for feedforward ANNs with single hidden layer (2-layer)
  - 2-layer feedforward ANN
    - Any <u>Boolean function</u> (simulate a 2-layer AND-OR network)
    - Any <u>bounded continuous function</u> (*approximate with arbitrarily small error*) [Cybenko, 1989; Hornik *et al*, 1989]
  - Sigmoid functions: set of <u>basis functions</u>; used to compose arbitrary functions
  - 3-layer feedforward ANN: any function (*approximate with arbitrarily small error*) [Cybenko, 1988]
  - Functions that ANNs are good at acquiring: <u>Network Efficiently Representable Functions</u> (NERFs) - how to characterize? [Russell and Norvig, 1995]
- **Inductive Bias of ANNs**
  - *n*-dimensional Euclidean space (<u>weight space</u>)
  - Continuous (error function smooth with respect to weight parameters)
  - Preference bias: "smooth interpolation" among positive examples
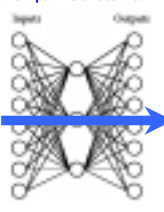  - Not well understood yet (known to be computationally hard)

CIS 830: Advanced Topics in Artificial Intelligence

**Kansas State University**
Department of Computing and Information Sciences

---

## Learning Hidden Layer Representations

- **Hidden Units and <u>Feature Extraction</u>**
  - Training procedure: hidden unit representations that minimize error E
  - Sometimes backprop will define new hidden features that are not explicit in the input representation *x*, but which capture properties of the input instances that are most relevant to learning the target function *t(x)*
  - Hidden units express *newly constructed features*
  - *Change of representation* to linearly separable *D'*
- **A Target Function (<u>Sparse</u> *aka* <u>1-of-C</u>, Coding)**

| Input | Hidden Values | Output |
|---|---|---|
| 1 0 0 0 0 0 0 0 → | 0.89 0.04 0.08 → | 1 0 0 0 0 0 0 0 |
| 0 1 0 0 0 0 0 0 → | 0.01 0.11 0.88 → | 0 1 0 0 0 0 0 0 |
| 0 0 1 0 0 0 0 0 → | 0.01 0.97 0.27 → | 0 0 1 0 0 0 0 0 |
| 0 0 0 1 0 0 0 0 → | 0.99 0.97 0.71 → | 0 0 0 1 0 0 0 0 |
| 0 0 0 0 1 0 0 0 → | 0.03 0.05 0.02 → | 0 0 0 0 1 0 0 0 |
| 0 0 0 0 0 1 0 0 → | 0.22 0.99 0.99 → | 0 0 0 0 0 1 0 0 |
| 0 0 0 0 0 0 1 0 → | 0.80 0.01 0.98 → | 0 0 0 0 0 0 1 0 |
| 0 0 0 0 0 0 0 1 → | 0.60 0.94 0.01 → | 0 0 0 0 0 0 0 1 |

  - Can this be learned? (Why or why not?)

CIS 830: Advanced Topics in Artificial Intelligence

**Kansas State University**
Department of Computing and Information Sciences

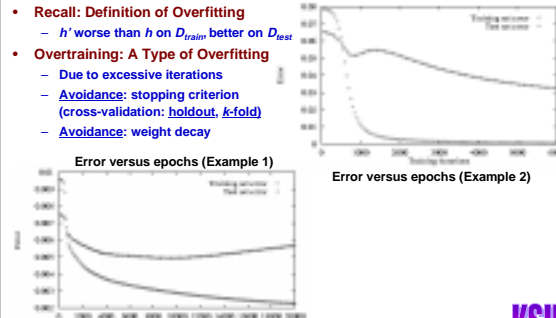---

## Convergence of Backpropagation

- **No Guarantee of Convergence to Global Optimum Solution**
  - Compare: perceptron convergence (to best $h \in H$, *provided $h \in H$*; i.e., LS)
  - Gradient descent to some local error minimum (perhaps not global minimum…)
  - Possible improvements on backprop (<u>BP</u>)
    - Momentum term (BP variant with slightly different weight update rule)
    - <u>Stochastic gradient descent</u> (BP algorithm variant)
    - Train multiple nets with different initial weights; find a good <u>mixture</u>
  - Improvements on feedforward networks
    - <u>Bayesian learning</u> for ANNs (e.g., <u>simulated annealing</u>) - later
    - Other global optimization methods that integrate over multiple networks
- **Nature of Convergence**
  - Initialize weights near zero
  - Therefore, initial network near-linear
  - Increasingly non-linear functions possible as training progresses

CIS 830: Advanced Topics in Artificial Intelligence

**Kansas State University**
Department of Computing and Information Sciences

---

## Overtraining in ANNs

- **Recall: Definition of Overfitting**
  - *h'* worse than *h* on $D_{train}$, better on $D_{test}$
- **Overtraining: A Type of Overfitting**
  - Due to excessive iterations
  - <u>Avoidance</u>: stopping criterion (cross-validation: <u>holdout</u>, *k-fold*)
  - <u>Avoidance</u>: weight decay

**Error versus epochs (Example 1)**

**Error versus epochs (Example 2)**



CIS 830: Advanced Topics in Artificial Intelligence

**Kansas State University**
Department of Computing and Information Sciences
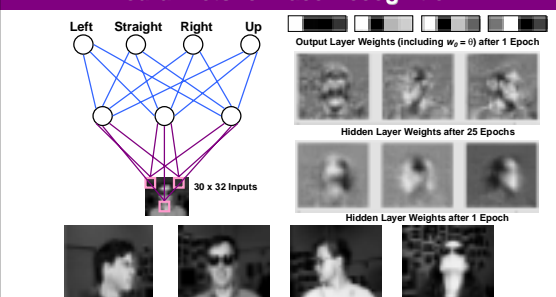
---

## Overfitting in ANNs

- **Other Causes of Overfitting Possible**
  - Number of hidden units sometimes set in advance
  - Too few hidden units ("underfitting")
    - ANNs with no growth
    - Analogy: underdetermined linear system of equations (more unknowns than equations)
  - Too many hidden units
    - ANNs with no pruning
    - Analogy: fitting a quadratic polynomial with an approximator of degree >> 2
- **Solution Approaches**
  - <u>Prevention</u>: <u>attribute subset selection</u> (using pre-filter or wrapper)
  - <u>Avoidance</u>
    - Hold out cross-validation (CV) set or split *k* ways (when to stop?)
    - Weight decay: decrease each weight by some factor on each epoch
  - <u>Detection/recovery</u>: <u>random restarts</u>, <u>addition and deletion</u> of weights, units

CIS 830: Advanced Topics in Artificial Intelligence

**Kansas State University**
Department of Computing and Information Sciences

---

## Example: Neural Nets for Face Recognition



Left   Straight   Right   Up

30 x 32 Inputs

Output Layer Weights (including $w_0 = \theta$) after 1 Epoch

Hidden Layer Weights after 25 Epochs

Hidden Layer Weights after 1 Epoch

- 90% Accurate Learning Head Pose, Recognizing 1-of-20 Faces
- http://www.cs.cmu.edu/~tom/faces.html

CIS 830: Advanced Topics in Artificial Intelligence

**Kansas State University**
Department of Computing and Information Sciences

## Example: *NetTalk*

- **Sejnowski and Rosenberg, 1987**
- **Early Large-Scale Application of Backprop**
  - Learning to convert text to speech
    - Acquired model: *a mapping from letters to phonemes and stress marks*
    - Output passed to a speech synthesizer
  - Good performance after training on a vocabulary of ~1000 words
- **Very Sophisticated Input-Output Encoding**
  - Input: 7-letter window; determines the phoneme for the center letter and context on each side; distributed (i.e., sparse) representation: 200 bits
  - Output: units for articulatory modifiers (e.g., "voiced"), stress, closest phoneme; distributed representation
  - 40 hidden units; 10000 weights total
- **Experimental Results**
  - Vocabulary: trained on 1024 of 1463 (informal) and 1000 of 20000 (dictionary)
  - 78% on informal, ~60% on dictionary
- **http://www.boltz.cs.cmu.edu/benchmarks/nettalk.html**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Recurrent Networks

- **Representing Time Series with ANNs**
  - Feedforward ANN: $y(t+1) = net\ (x(t))$
  - Need to capture temporal relationships



- **Solution Approaches**
  - Directed cycles
  - Feedback
    - Output-to-input [Jordan]
    - Hidden-to-input [Elman]
    - Input-to-input
  - Captures time-lagged relationships
    - Among $x(t' \leq t)$ and $y(t+1)$
    - Among $y(t' \leq t)$ and $y(t+1)$
  - Learning with recurrent ANNs
    - Elman, 1990; Jordan, 1987
    - Principe and deVries, 1992
    - Mozer, 1994; Hsu and Ray, 1998

CIS 830: Advanced Topics in Artificial Intelligence

---

## Some Current Issues and Open Problems in ANN Research

- **Hybrid Approaches**
  - Incorporating knowledge and analytical learning into ANNs
    - Knowledge-based neural networks [Flann and Dietterich, 1989]
    - Explanation-based neural networks [Towell *et al*, 1990; Thrun, 1996]
  - Combining uncertain reasoning and ANN learning and inference
    - Probabilistic ANNs
    - Bayesian networks [Pearl, 1988; Heckerman, 1996; Hinton *et al*, 1997] - later
- **Global Optimization with ANNs**
  - Markov chain Monte Carlo (MCMC) [Neal, 1996] - e.g., simulated annealing
  - Relationship to genetic algorithms - later
- **Understanding ANN Output**
  - Knowledge extraction from ANNs
    - Rule extraction
    - Other decision surfaces
  - Decision support and KDD applications [Fayyad *et al*, 1996]
- **Many, Many More Issues (Robust Reasoning, Representations, etc.)**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Some ANN Applications

- **Diagnosis**
  - Closest to pure concept learning and classification
  - Some ANNs can be post-processed to produce probabilistic diagnoses
- **Prediction and Monitoring**
  - *aka* prognosis (sometimes forecasting)
  - Predict a continuation of (typically numerical) data
- **Decision Support Systems**
  - *aka* recommender systems
  - Provide assistance to human "subject matter" experts in making decisions
    - Design (manufacturing, engineering)
    - Therapy (medicine)
    - Crisis management (medical, economic, military, computer security)
- **Control Automation**
  - Mobile robots
  - Autonomic sensors and actuators
- **Many, Many More (ANNs for Automated Reasoning, etc.)**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Terminology

- **Multi-Layer ANNs**
  - Focused on one species: (feedforward) multi-layer perceptrons (MLPs)
  - Input layer: an implicit layer containing $x_i$
  - Hidden layer: a layer containing input-to-hidden unit weights and producing $h_j$
  - Output layer: a layer containing hidden-to-output unit weights and producing $o_k$
  - $n$-layer ANN: an ANN containing $n$ - 1 hidden layers
  - Epoch: one training iteration
  - Basis function: set of functions that span $H$
- **Overfitting**
  - Overfitting: $h$ does better than $h'$ on training data and worse on test data
  - Overtraining: overfitting due to training for too many epochs
  - Prevention, avoidance, and recovery techniques
    - Prevention: attribute subset selection
    - Avoidance: stopping (termination) criteria (CV-based), weight decay
- **Recurrent ANNs: Temporal ANNs with Directed Cycles**

CIS 830: Advanced Topics in Artificial Intelligence

---

## Summary Points

- **Multi-Layer ANNs**
  - Focused on feedforward MLPs
  - Backpropagation of error: distributes penalty (loss) function throughout network
  - Gradient learning: takes derivative of error surface with respect to weights
    - Error is based on difference between desired output (*t*) and actual output (*o*)
    - Actual output (*o*) is based on activation function
    - Must take partial derivative of $\sigma \Rightarrow$ choose one that is easy to differentiate
    - Two $\sigma$ definitions: sigmoid (*aka* logistic) and hyperbolic tangent (*tanh*)
- **Overfitting in ANNs**
  - Prevention: attribute subset selection
  - Avoidance: cross-validation, weight decay
- **ANN Applications: Face Recognition, Text-to-Speech**
- **Open Problems**
- **Recurrent ANNs: Can Express Temporal Depth (Non-Markovity)**
- **Next: Neural Reinforcement Learning**

CIS 830: Advanced Topics in Artificial Intelligence