

Conclusions and Final Review

Thursday, December 06, 2001

William H. Hsu

Department of Computing and Information Sciences, KSU

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

Readings: Chapters 1-10, 13, Mitchell Chapters 14-21, Russell and Norvig



CIS 732: Machine Learning and Pattern Recognition

Kansas State University Department of Computing and Information Sciences

Lecture 0:

A Brief Overview of Machine Learning

- Overview: Topics, Applications, Motivation
- Learning = Improving with Experience at Some Task
 - Improve over task T,
 - with respect to performance measure *P*,
 - based on experience *E*.
- Brief Tour of Machine Learning
 - A case study
 - A taxonomy of learning
 - Intelligent systems engineering: *specification of learning problems*
- Issues in Machine Learning
 - Design choices
 - The performance element: intelligent systems
- Some Applications of Learning
 - Database mining, reasoning (inference/decision support), acting
 - Industrial usage of intelligent systems





Department of Computing and Information Sciences

Lecture 1: Concept Learning and Version Spaces

- Concept Learning as Search through *H*
 - Hypothesis space *H* as a state space
 - Learning: finding the correct hypothesis
- General-to-Specific Ordering over H
 - Partially-ordered set: Less-Specific-Than (More-General-Than) relation
 - Upper and lower bounds in H
- Version Space Candidate Elimination Algorithm
 - S and G boundaries characterize learner's uncertainty
 - Version space can be used to make predictions over unseen cases
- Learner Can Generate Useful Queries
- Next Lecture: When and Why Are Inductive Leaps Possible?



Lectures 2-3: Introduction to COLT

- The Need for Inductive Bias
 - Modeling inductive learners with equivalent deductive systems
 - Kinds of biases: preference (search) and restriction (language) biases
- Introduction to Computational Learning Theory (COLT)
 - Things COLT attempts to measure
 - Probably-Approximately-Correct (PAC) learning framework
- COLT: Framework Analyzing Learning Environments
 - Sample complexity of *C*, computational complexity of *L*, required expressive power of *H*
 - Error and confidence bounds (PAC: 0 < ? < 1/2, 0 < ? < 1/2)
- What PAC Prescribes
 - Whether to try to learn C with a known H
 - Whether to try to reformulate H (apply change of representation)
- Vapnik-Chervonenkis (VC) Dimension: Measures Expressive Power of H
- Mistake Bounds



Lectures 4-5: Decision Tree Induction

- Model ("Architecture"): Decision Trees (DTs)
- Algorithm *Build-DT*: Top Down Induction
- Entropy and Information Gain: *ID*3
 - Goal: to measure *uncertainty removed* by splitting on a candidate attribute A
 - ID3? Build-DT using Gain(•)
- ID3 as Hypothesis Space Search (in State Space of Decision Trees)
- Data Mining using *MLC*++ (Machine Learning Library in C++)
- Occam's Razor and Decision Trees
 - Preference biases versus language biases
 - Minimum Description Length (MDL) justification for Occam's Razor biases
- Overfitting
 - Problem: fitting training data too closely
 - Overfitting prevention, avoidance, and recovery techniques
- Other Ways to Make Decision Tree Induction More Robust



Lectures 6-7: ANNs (Perceptrons, MLPs, Backprop)

- Neural Networks: Parallel, Distributed Processing Systems
 - Biological and artificial (ANN) types
 - Perceptron (<u>linear threshold unit/gate</u>, aka LTU/LTG): model neuron
- Update Rules for Single-Layer Networks
 - Multiplicative (Hebbian, Winnow), additive (gradient: Perceptron, Delta Rule)
 - Batch versus incremental mode
- Advantages and Disadvantages of LTG
 - "Disadvantage" (tradeoff): simple and restrictive
 - "Advantage": perform well on many realistic problems (e.g., some text learning)
- Multi-Layer ANNs and Backprop
 - 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)
- Overfitting in ANNs
- Recurrent ANNs: Can Express Temporal <u>Depth</u> (<u>Non-Markovity</u>)

Department of Computing and Information Sciences

Lecture 8: Statistical Evaluation of Hypotheses

- Statistical Evaluation Methods for Learning: Three Questions
 - Generalization quality
 - How well does observed accuracy <u>estimate</u> generalization accuracy?
 - Estimation bias and variance
 - Confidence intervals
 - Comparing generalization quality
 - How certain are we that h_1 is better than h_2 ?
 - Confidence intervals for paired tests
 - Learning and statistical evaluation
 - What is the best way to make the most of limited data?
 - *k*-fold CV
- Tradeoffs: Bias versus Variance
- Next: Sections 6.1-6.5, Mitchell (Bayes's Theorem; ML; MAP)



Lectures 9-11: Bayesian Learning

- Introduction to Bayesian Learning: Bayes's Theorem
 - Definition of conditional (posterior) probability P(x | y)
 - Bayes's Theorem: $P(x | y) = P(y | x) \cdot P(x) / P(y)$
- <u>Maximum A Posteriori (MAP) and Maximum Likelihood (ML) Hypotheses</u>
 - Bayes's Rule and MAP
 - Uniform priors: allow use of MLE to generate MAP hypotheses
 - Relation to version spaces, candidate elimination
- <u>Minimum Description Length (MDL) Revisited</u>
 - <u>Bayesian Information Criterion (BIC)</u>: justification for Occam's Razor
- <u>Bayes Optimal Classifier (BOC)</u>
 - Using BOC as a "gold standard"
- Gibbs Classifier
 - Ratio bound
- Simple (Naïve) Bayes

 \mathbf{v}_{NB} ? arg max $\mathbf{P}_{i}^{2}\mathbf{v}_{j}^{2}$ $\mathbf{P}_{i}^{2}\mathbf{x}_{i}^{2}/\mathbf{v}_{j}^{2}$



Rationale for assumption; pitfalls

CIS 732: Machine Learning and Pattern Recognition

Kansas State University Department of Computing and Information Sciences

Lectures 12-13: Bayesian Belief Networks (BBNs)

- Graphical Models of Probability
 - Bayesian networks: introduction
 - Definition and basic principles
 - Conditional independence (causal Markovity) assumptions, tradeoffs
 - Inference and learning using Bayesian networks
 - Inference in polytrees (singly-connected BBNs)
 - Acquiring and applying CPTs: gradient algorithm *Train-BN*
- Structure Learning in Trees: MWST Algorithm *Learn-Tree-Structure*
- Reasoning under Uncertainty using BBNs
 - Learning, eliciting, applying CPTs
 - In-class exercise: *Hugin* demo; CPT elicitation, application
 - Learning BBN structure: <u>constraint-based</u> versus <u>score-based</u> approaches
 - K2, other scores and search algorithms
- Causal Modeling and Discovery: Learning Causality from Observations
- Incomplete Data: Learning and Inference (Expectation-Maximization)

Lecture 15:

EM, Unsupervised Learning, and Clustering

- Expectation-Maximization (EM) Algorithm
- Unsupervised Learning and Clustering
 - Types of unsupervised learning
 - Clustering, vector quantization
 - Feature extraction (typically, dimensionality reduction)
 - Constructive induction: unsupervised learning in support of supervised learning
 - Feature construction (aka feature extraction)
 - Cluster definition
 - Algorithms
 - EM: mixture parameter estimation (e.g., for AutoClass)
 - AutoClass: Bayesian clustering
 - Principal Components Analysis (PCA), factor analysis (FA)
 - <u>Self-Organizing Maps (SOM)</u>: projection of data; competitive algorithm
 - Clustering problems: <u>formation</u>, <u>segmentation</u>, <u>labeling</u>
- Next Lecture: Time Series Learning and Characterization



Lecture 16: Introduction to Time Series Analysis

- Introduction to Time Series
 - 3 phases of analysis: forecasting (prediction), modeling, characterization
 - Probability and time series: <u>stochastic processes</u>
 - Linear models: ARMA models, approximation with temporal ANNs
 - Time series understanding and learning
 - Understanding: <u>state-space reconstruction</u> by <u>delay-space embedding</u>
 - Learning: parameter estimation (e.g., using temporal ANNs)
- Further Reading
 - Analysis: Box et al, 1994; Chatfield, 1996; Kantz and Schreiber, 1997
 - Learning: Gershenfeld and Weigend, 1994
 - Reinforcement learning: next...
- Next Lecture: Policy Learning, Markov Decision Processes (MDPs)
 - Read Chapter 17, Russell and Norvig, Sections 13.1-13.2, Mitchell
 - Exercise: 16.1(a), Russell and Norvig (bring answers to class; don't peek!)



Lecture 17: Policy Learning and MDPs

- Making Decisions in Uncertain Environments
 - Framework: <u>Markov Decision Processes</u>, <u>Markov Decision Problems</u> (MDPs)
 - Computing policies
 - Solving MDPs by dynamic programming given a stepwise reward
 - Methods: value iteration, policy iteration
 - Decision-theoretic agents
 - Decision cycle, Kalman filtering
 - Sensor fusion (*aka* data fusion)
 - <u>Dynamic Bayesian networks</u> (DBNs) and <u>dynamic decision networks</u> (DDNs)
- Learning Problem
 - Mapping from observed actions and rewards to decision models
 - Rewards/penalties: reinforcements
- Next Lecture: Reinforcement Learning
 - Basic model: passive learning in a known environment
 - **Q** learning: policy learning by <u>a</u>daptive <u>dynamic programming (ADP)</u>



Lecture 18: Introduction to Reinforcement Learning

- Control Learning
 - Learning policies from <*state*, *reward*, *action*> observations
 - Objective: choose optimal actions given new percepts and incremental rewards
 - Issues
 - Delayed reward
 - Active learning opportunities
 - Partial observability
 - Reuse of sensors, effectors
- Q Learning
 - Action-value function Q: state? action? value (expected utility)
 - Training rule
 - Dynamic programming algorithm
 - Q learning for deterministic worlds
 - Convergence to true Q
 - Generalizing **Q** learning to nondeterministic worlds
- Next Week: More Reinforcement Learning (Temporal Differences)



Lecture 19: More Reinforcement Learning (TD)

- <u>Reinforcement Learning (RL)</u>
 - Definition: learning policies ? : state ? action from << state, action>, reward>
 - <u>Markov decision problems (MDPs</u>): finding control policies to choose optimal actions
 - <u>Q-learning</u>: produces action-value function Q: state? action? value (expected utility)
 - Active learning: <u>experimentation (exploration) strategies</u>
 - Exploration function: f(u, n)
 - Tradeoff: greed (u) preference versus novelty (1 / n) preference, aka curiosity
- <u>Temporal Diffference (TD) Learning</u>
 - ?: constant for blending alternative training estimates from multi-step lookahead
 - **TD(?)**: algorithm that uses recursive training rule with **?**-estimates
- Generalization in RL
 - Explicit representation: tabular representation of U, M, R, Q
 - <u>Implicit representation</u>: <u>compact</u> (*aka* <u>compressed</u>) representation



Lecture 20: Neural Computation

- Review: Feedforward <u>Artificial Neural Networks (ANNs</u>)
- Advanced ANN Topics
 - Models
 - Modular ANNs
 - Associative memories
 - Boltzmann machines
 - Applications
 - Pattern recognition and scene analysis (image processing)
 - Signal processing
 - Neural reinforcement learning
- Relation to Bayesian Networks and Genetic Algorithms (GAs)
 - Bayesian networks as a species of connectionist model
 - Simulated annealing and GAs: MCMC methods
 - Numerical ("subsymbolic") and symbolic AI systems: principled integration
- Next Week: Combining Classifiers (WM, Bagging, Stacking, Boosting)



Lecture 21: Combiners (WM, Bagging, Stacking)

- Combining Classifiers
 - Problem definition and motivation: improving accuracy in concept learning
 - General framework: collection of weak classifiers to be improved (data fusion)
- <u>Weighted Majority (WM)</u>
 - Weighting system for collection of algorithms
 - Weights each algorithm in proportion to its training set accuracy
 - Use this weight in performance element (and on test set predictions)
 - Mistake bound for WM
- <u>Bootstrap Aggregating (Bagging)</u>
 - Voting system for collection of algorithms
 - Training set for each member: sampled with replacement
 - Works for unstable inducers
- Stacked Generalization (aka Stacking)
 - Hierarchical system for combining inducers (ANNs or other inducers)
 - Training sets for "leaves": sampled with replacement; combiner: validation set
- Next Lecture: <u>Boosting</u> the Margin, <u>Hierarchical Mixtures of Experts</u>



Lecture 22:

More Combiners (Boosting, Mixture Models)

- Committee Machines *aka* Combiners
- Static Structures (Single-Pass)
 - Ensemble averaging
 - For improving weak (especially unstable) classifiers
 - e.g., weighted majority, bagging, stacking
 - Boosting the margin
 - Improve performance of any inducer: weight examples to emphasize errors
 - Variants: filtering (*aka* consensus), resampling (*aka* subsampling), reweighting
- Dynamic Structures (Multi-Pass)
 - Mixture of experts: training in combiner inducer (aka gating network)
 - Hierarchical mixtures of experts: hierarchy of inducers, combiners
- Mixture Model (aka Mixture of Experts)
 - Estimation of mixture coefficients (i.e., weights)
 - Hierarchical Mixtures of Experts (HME): multiple combiner (gating) levels
- Next Week: Intro to GAs, GP (9.1-9.4, Mitchell; 1, 6.1-6.5, Goldberg)



Lecture 23:

Introduction to Genetic Algorithms (GAs)

- Evolutionary Computation
 - Motivation: process of natural selection
 - Limited population; individuals compete for membership
 - Method for parallelizing and stochastic search
 - Framework for problem solving: search, optimization, learning
- Prototypical (Simple) <u>Genetic Algorithm (GA)</u>
 - Steps
 - Selection: reproduce individuals probabilistically, in proportion to fitness
 - Crossover: generate new individuals probabilistically, from pairs of "parents"
 - Mutation: modify structure of individual randomly
 - How to represent hypotheses as individuals in GAs
- An Example: <u>GA-Based Inductive Learning (GABIL)</u>
- Schema Theorem: Propagation of Building Blocks
- Next Lecture: Genetic Programming, The Movie





Lecture 24:

Introduction to Genetic Programming (GP)

- Genetic Programming (GP)
 - Objective: program synthesis
 - Application of evolutionary computation (especially genetic algorithms)
 - Search algorithms
 - Based on mechanics of natural selection, natural genetics
 - Design application
- Steps in GP Design
 - <u>Terminal set</u>: program variables
 - Function set: operators and macros
 - <u>Fitness cases</u>: evaluation environment (compare: validation tests in software engineering)
 - <u>Control parameters</u>: "runtime" configuration variables for GA (population size and organization, number of generations, syntactic constraints)
 - <u>Termination criterion and result designation</u>: when to stop, what to return
- Next Week: Instance-Based Learning (IBL)



Lecture 25: Instance-Based Learning (IBL)

- <u>Instance Based Learning (IBL)</u>
 - k-Nearest Neighbor (k-NN) algorithms
 - When to consider: few continuous valued attributes (low dimensionality)
 - Variants: distance-weighted k-NN; k-NN with attribute subset selection
 - Locally-weighted regression: function approximation method, generalizes *k*-NN
 - <u>Radial-Basis Function (RBF) networks</u>
 - Different kind of artificial neural network (ANN)
 - Linear combination of local approximation ? global approximation to f(?)
- <u>Case-Based Reasoning (CBR) Case Study: CADET</u>
 - Relation to IBL
 - CBR online resource page: <u>http://www.ai-cbr.org</u>
- Lazy and Eager Learning
- Next Week
 - Rule learning and extraction
 - <u>Inductive logic programming (ILP)</u>





Lecture 26: Rule Learning and Extraction

- Learning Rules from Data
- Sequential Covering Algorithms
 - Learning single rules by search
 - Beam search
 - Alternative covering methods
 - Learning rule sets
- First-Order Rules
 - Learning single first-order rules
 - Representation: first-order Horn clauses
 - Extending Sequential-Covering and Learn-One-Rule: variables in rule preconditions
 - FOIL: learning first-order rule sets
 - Idea: inducing logical rules from observed relations
 - Guiding search in FOIL
 - Learning recursive rule sets
- Next Time: Inductive Logic Programming (ILP)



Lecture 27: Inductive Logic Programming (ILP)

- Induction as Inverse of Deduction
 - Problem of induction revisited
 - Definition of induction
 - Inductive learning as specific case
 - Role of induction, deduction in automated reasoning
 - Operators for automated deductive inference
 - Resolution rule (and operator) for deduction
 - <u>First-order predicate calculus (FOPC) and resolution theorem proving</u>
 - Inverting resolution
 - Propositional case
 - First-order case (inverse entailment operator)
- <u>Inductive Logic Programming (ILP)</u>
 - Cigol: inverse entailment (very susceptible to combinatorial explosion)
 - Progol: sequential covering, general-to-specific search using inverse entailment
- Next Week: <u>Knowledge Discovery in Databases (KDD</u>), Final Review



Lecture 28: KDD and Data Mining

- Knowledge Discovery in Databases (KDD) and Data Mining
 - <u>Stages</u>: selection (filtering), processing, transformation, learning, inference
 - Design and implementation issues
- Role of Machine Learning and Inference in Data Mining
 - Roles of unsupervised, supervised learning in KDD
 - Decision support (information retrieval, prediction, policy optimization)
- Case Studies
 - Risk analysis, transaction monitoring (filtering), prognostic monitoring
 - Applications: business decision support (pricing, fraud detection), automation
- Resources Online
 - Microsoft DMX Group (Fayyad): <u>http://research.microsoft.com/research/DMX/</u>
 - KSU KDD Lab (Hsu): <u>http://ringil.cis.ksu.edu/KDD/</u>
 - CMU KDD Lab (Mitchell): <u>http://www.cs.cmu.edu/~cald</u>
 - KD Nuggets (Piatetsky-Shapiro): <u>http://www.kdnuggets.com</u>
 - NCSA Automated Learning Group (Welge)
 - ALG home page: <u>http://www.ncsa.uiuc.edu/STI/ALG</u>
 - NCSA D2K: <u>http://chili.ncsa.uiuc.edu</u>





Meta-Summary

- Machine Learning Formalisms
 - Theory of computation: PAC, mistake bounds
 - Statistical, probabilistic: PAC, confidence intervals
- Machine Learning Techniques
 - Models: version space, decision tree, perceptron, winnow, ANN, BBN, SOM, Q functions, GA/GP building blocks (<u>schemata</u>), GP building blocks
 - Algorithms: candidate elimination, *ID3*, backprop, MLE, Simple (Naïve) Bayes, *K2*,
 EM, SOM convergence, LVQ, ADP, *Q*-learning, TD(?), simulated annealing, sGA
- Final Exam Study Guide
 - Know
 - Definitions (terminology)
 - How to solve problems from Homeworks 1 and 3 (problem sets)
 - How algorithms in Homeworks 2, 4, and 5 (machine problems) work
 - Practice
 - Sample exam problems (handout)
 - Example runs of algorithms in Mitchell, lecture notes
 - Don't panic! «

