

# Principles of Learning in Humans and Machines

Alexander Petrov

Carnegie Mellon University  
<http://www.andrew.cmu.edu/~apetrov/>

This course provides a forest-level overview of the field of machine learning, highlighting its deep relationship with statistics, psychology, and philosophy of science. All these disciplines address complementary aspects of the same inductive problem—how to extract knowledge from the environment and use it to improve future performance. Various learning algorithms are presented with an emphasis on the underlying ideas rather than technical rigor. The algorithms are illustrated with applications from cognitive modeling, robotics, and artificial intelligence. These diverse learning systems are compared in an attempt to extract principles that apply to them all and thus characterize learning in general. The implications of these principles for psychology (the nature-nurture problem) and philosophy of science (the problem of induction) are discussed.

## **Class 1: Foundations: How and when is learning possible? Candidate elimination and statistical estimation**

This first class gets the ball rolling by considering two simple learning systems: one based on logic and one on statistics. The main objective is to illustrate the inductive problem and to introduce the concepts of hypothesis space, instance space, and inductive bias in relation to the problem of induction in philosophy of science.

- Version spaces and the candidate elimination algorithm
- Inductive bias as a necessary precondition for generalization
- The role of prior knowledge in learning
- Can we know the universe?
- Statistical estimation: linear regression
- Terminology and map of the terrain

*Required readings:* (See end of syllabus for the exact references)

Mitchell, T. (1997). Concept learning and the general-to-specific ordering. Chapter 2 in *Machine learning* (pp. 20-51).

Sagan, C. (1974). Can we know the universe? Reflections on a grain of salt. Chapter 2 in *Broca's brain* (pp. 15-21).

## **Class 2: Occam's razor: Overfitting and cross validation**

### **Decision trees and memory-based methods**

Is it a good idea to prefer simple hypotheses and why? And what does the word “simple” mean anyway? The main objective of this class is to stress that success in learning is measured by the performance on novel cases rather than fits of past data. Two widely used machine learning techniques are introduced along the way: decision trees and memory-based methods.

- Induction of decision trees (ID3)
- Preference (search) bias vs. restriction (language) bias
- Overfitting and cross validation
- The bias-variance dilemma in statistical estimation
- Model selection criteria
- Memory-based (lazy) methods do the inductive leap at test rather than during training
- $k$  nearest neighbor
- Locally weighted regression
- The curse of dimensionality
- Radial-basis function networks
- Gordon Logan's instance-based theory of automatization

#### *Required readings:*

Mitchell, T. (1997). Decision tree learning. Chapter 3 in *Machine learning* (pp. 52-80).

Mitchell, T. (1997). Instance-based learning. Ch. 8 in *Machine learning* (pp. 230-248).

#### *Optional readings:*

Atkeson et al. (1997). Locally weighted learning. *AI Review*, 11, 11-73.

Gell-Mann, M. (1994). Information and crude complexity. Chapter 3 in *The quark and the jaguar* (pp. 23-41).

Logan, G. (1988). Toward an instance theory of automatization. *Psychological Review*, 95, 492-527.

Zucchini, W. (2000). An introduction to model selection. *Journal of Mathematical Psychology*, 44, 41-61.

## **Class 3: Occam's razor in neural nets. Genetic connectionism**

### **Active learning and explanation-based learning**

This class consists of two quite unrelated parts cramped into a single day because of time constraints. The first part explicates how connectionist learning systems fit into the general scheme formulated earlier. The intimate link between evolution and learning is discussed on the basis of a simulated genetic experiment on learning rules for neural networks.

The second part of the class introduces briefly two exciting extensions of the learning paradigm -- active and theory-based approaches. Two psychological theories are discussed in the light of these new concepts.

- Neural nets and backpropagation
- Occam in the network world: weight decay and weight elimination
- Experiment in genetic connectionism
- The two nested loops of adaptation: alliance of nature and nurture
- Neurophysiological mechanisms of learning and gene expression
- Active learning: experimentation augments passive observation
- David Klahr's dual search model of scientific discovery
- Explanation-based learning: getting the most of sparse data and a prior theory of the domain
- The Theory Theory of cognitive development

*Required readings:*

Chalmers, D. (1990). The evolution of learning: An experiment in genetic connectionism.

Klahr, D. (2000). Scientific discovery as problem solving. Chapter 2 in *Exploring science* (pp. 21-39).

*Optional readings:*

Jacobs, R. & Jordan, M. (1992). Computational consequences of a bias toward short connections. *Journal of Cognitive Neuroscience*, 4, 323-336.

Mitchell, T. (1997). Combining inductive and analytical learning. Chapter 12 in *Machine learning* (pp. 334-366).

#### **Class 4: Bayesian learning**

The true logic of this world is in the calculus of probabilities, proclaimed James Clerk Maxwell. Following his dictum, we will discuss the important class of systems that represent knowledge as probability distributions and learn by explicitly or implicitly manipulating probabilities. It turns out that the learning algorithms discussed so far have straightforward Bayesian interpretation. Occam's razor reappears under the guise of the minimum description length principle. As an extra bonus, the EM algorithm is introduced in the context of density estimation with Gaussian mixtures.

- Subjective probabilities
- Bayes theorem: data-based conversion of priors into posteriors
- Maximum a posteriori (MAP) learner
- The minimum description length principle
- Application: Naive Bayes classification of documents
- Bayes optimal classifier
- John Anderson's rational analysis of cognition
- Belief networks
- Generative neural networks
- Gibbs sampling

- Gaussian mixtures for density estimation
- Expectation maximization (EM) algorithm
- Bayesian framework for philosophy of science

*Required readings:*

Feynman(1963). The meaning of it all (pp. 15-28, 64-71).  
 Mitchell, T. (1997). Bayesian learning. Chapter 6 in *Machine learning* (pp. 154-200).

*Optional readings:*

Anderson, J.R. (1991). Is human cognition adaptive? *Behavioral and Brain Sciences*, 14, 471-517.  
 Charniak, E. (1991). Bayesian networks without tears. *AI Magazine*, 12, 50-63.  
 Hinton, G. (1999). Products of Experts. *ICANN 99*.  
 Howson, C. & Urbach, P. (1993). *Scientific reasoning: The Bayesian approach*.  
 Chapter 1: Introduction (pp. 3-15).  
 Neal, R. (1996). *Bayesian learning for neural networks*. Ch 1: Introduction (pp. 1-28).

**Class 5: Reinforcement learning: Discovering actions to maximize rewards**

Why is this obsession with all these function approximators, posterior probabilities, density estimators, and other obscurities? Answer: they are useful tricks for coping with the real, dangerous, and constantly changing world out there. This last session considers the challenges faced by the agents that try to go the full distance—live in an uncertain environment, explore it, and figure out what actions yield maximal reward in the long run.

- The reinforcement learning problem
- Markov systems
- Immediate vs. future rewards; discounting
- Dynamic programming
- Learning policies for Markov decision processes
- The credit assignment problem
- Exploration vs. exploitation
- Temporal difference learning
- Reinforcement learning for neural networks
- Application: TD-Gammon
- Application: Robot control
- The actor-critic architecture

*Required reading:*

Sutton, R. & Barto, A. (1998). *Reinforcement learning: An introduction*. Chapters 1 and 2 only (pp. 3-49).

*Optional readings:*

Bertsekas, D. & Tsitsiklis, J. (1996). *Neuro-dynamic programming*. Chapter 1: Introduction (pp. 1-10).  
 Harmon, M. (1996). *Reinforcement learning: A tutorial*.  
<http://www-anw.cs.umass.edu/~mharmon/rltutorial/>

Kaelbling et al. (1996). Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4, 237-285.

Mitchell, T. (1997). Reinforcement learning. Ch 13 in *Machine learning* (pp. 367-390).

### **Small groups:**

Discussion of the issues raised in the main classes and/or brought up by the participants. Design of learning systems for tasks proposed by the participants.

The course is designed to be accessible to graduate students with non-technical background (e.g. psychology, linguistics). A minimal level of proficiency with some basic concepts of statistics (e.g. conditional probability) and computer science (e.g. gradient descent search) will be extremely helpful, however. The first afternoon will offer an “executive briefing” for students who need to be brought up to speed with such concepts.

### **Assessment:**

Students who desire credit should write a 10-page paper describing how the ideas discussed in the course relate to their own research, challenge the opinions of the instructor, propose a concrete learning system for a task of interest to the student, etc.

### **Bibliography:**

Anderson, J.R. (1991). Is human cognition adaptive? *Behavioral and Brain Sciences*, 14, 471-517.

Atkeson, C., Moore, A., & Schall, S. (1997). Locally weighted learning. *AI Review*, 11, 11-73.

[Available on-line at <ftp://ftp.cc.gatech.edu/pub/people/cga/air1.html>]

Bertsekas, D. & Tsitsiklis, J. (1996). *Neuro-dynamic programming*. Belmont, MA: Athena Scientific.

Chalmers, D. (1990). The evolution of learning: An experiment in genetic connectionism. In D. Touretzky, J. Elman, T. Sejnowski, & G. Hinton (Eds.), *Proceedings of the 1990 Connectionist Models Summer School*. San Mateo, CA: Morgan Kaufmann. [Available on line at <ftp://ftp.cogsci.indiana.edu/pub/chalmers.evolution.ps>]

Charniak, E. (1991). Bayesian networks without tears. *AI Magazine*, 12, 50-63.

Feynman(1998/1963). *The meaning of at all: Thoughts of a citizen scientist*. Reading, MA: Perseus Books.

Gell-Mann, M. (1994). *The quark and the jaguar: Adventures in the simple and the complex*. New York: W.H. Freeman and Company.

Harmon, M. (1996). Reinforcement learning: A tutorial.

<http://www-anw.cs.umass.edu/~mharmon/rltutorial/>

Hinton, G. (1999). Products of Experts. *Proceedings of the Ninth International Conference on Artificial Neural Networks* (ICANN 99, vol 1, pp. 1-6). [Available on-line at

<http://www.gatsby.ucl.ac.uk/publications/papers/06-1999.html>]

Howson, C. & Urbach, P. (1993). *Scientific reasoning: The Bayesian approach*. Chicago: Open Court.

Jacobs, R. & Jordan, M. (1992). Computational consequences of a bias toward short connections. *Journal of Cognitive Neuroscience*, 4, 323-336.

Kaelbling, L., Littman, M., & Moore, A. (1996). Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4, 237-285.

<http://www.cs.washington.edu/research/jair/abstracts/kaelbling96a.html>

Klahr, D. (2000). *Exploring science: The cognition and development of discovery processes*. Cambridge, MA: MIT Press.

Logan, G. (1988). Toward an instance theory of automatization. *Psychological Review*, 95, 492-527.

**Mitchell, T. (1997). *Machine learning*. New York: McGraw-Hill.**

<http://www.cs.cmu.edu/afs/cs.cmu.edu/user/mitchell/ftp/mlbook.html>

Neal, R. (1996). *Bayesian learning for neural networks*. New York: Springer-Verlag.

Sagan, C. (1974). Can we know the universe? Reflections on a grain of salt. In *Broca's brain* (chap. 2, pp. 15-21). New York: Ballantine Books. [Also available in M. Gardner (Ed.) (1994), *Great essays in science*. Amherst, NY: Prometheus Books.]

Sutton, R. & Barto, A. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press. [Full text available on-line at <http://www-anw.cs.umass.edu/~rich/book/the-book.html>]

Zucchini, W. (2000). An introduction to model selection. *Journal of Mathematical Psychology*, 44, 41-61. <http://quantrm2.psy.ohio-state.edu/injae/jmpsp.htm>

The following web pages also provide wonderful resources:

Machine learning resources (D. Aha): <http://www.aic.nrl.navy.mil/~aha/research/machine-learning.html>

Reinforcement learning links: <http://www-iiuf.unifr.ch/~aperezu/robotreinfo.html>

-----  
Alexander Petrov has a M.S. in computer science from Sofia University and a Ph.D. in cognitive science from New Bulgarian University. At present he is a post-doctoral research associate at the Department of Psychology at Carnegie Mellon University. More information about his vita and his research can be found on his personal web page at <http://www.andrew.cmu.edu/~apetrov/>

Syllabus last updated June 30, 2000