

Learning to Predict Visibility and Invisibility from Occlusion Events

Jonathan Marshall, Richard Alley, Robert Hubbard

Visual occlusion events constitute a major source of depth information. This paper presents a self-organizing neural network that learns to detect, represent, and predict the visibility and invisibility relationships that arise during occlusion events, after a period of exposure to motion sequences containing occlusion and disocclusion events. The network develops two parallel opponent channels or "chains" of lateral excitatory connections for every resolvable motion trajectory. One channel, the "On" chain or "visible" chain, is activated when a moving stimulus is visible. The other channel, the "Off" chain or "invisible" chain, carries a persistent, amodal representation that predicts the motion of a formerly visible stimulus that becomes invisible due to occlusion. The learning rule uses disinhibition from the On chain to trigger learning in the Off chain. The On and Off chain neurons can learn separate associations with object depth or(cid:173) dering. The results are closely related to the recent discovery (Assad & Maunsell, 1995) of neurons in macaque monkey posterior parietal cortex that respond selectively to inferred motion of invisible stimuli.

Onset-based Sound Segmentation

Leslie Smith

A technique for segmenting sounds using processing based on mam(cid:173) malian early auditory processing is presented. The technique is based on features in sound which neuron spike recording suggests are detected in the cochlear nucleus. The sound signal is band(cid:173) passed and each signal processed to enhance onsets and offsets. The onset and offset signals are compressed, then clustered both in time and across frequency channels using a network of integrate(cid:173) and-fire neurons. Onsets and offsets are signalled by spikes, and the timing of these spikes used to segment the sound.

Beating a Defender in Robotic Soccer: Memory-Based Learning of a Continuous Function

Peter Stone, Manuela Veloso

Learning how to adjust to an opponent's position is critical to the success of having intelligent agents collaborating towards the achievement of specific tasks in unfriendly environments. This pa(cid:173) per describes our work on a Memory-based technique for to choose an action based on a continuous-valued state attribute indicating the position of an opponent. We investigate the question of how an agent performs in nondeterministic variations of the training situ(cid:173) ations. Our experiments indicate that when the random variations fall within some bound of the initial training, the agent performs better with some initial training rather than from a tabula-rasa.

A Neural Network Autoassociator for Induction Motor Failure Prediction

Thomas Petsche, Angelo Marcantonio, Christian Darken, Stephen Hanson, Gary Kuhn, N. Santoso

We present results on the use of neural network based autoassociators which act as novelty or anomaly detectors to detect imminent motor failures. The autoassociator is trained to reconstruct spectra obtained from the healthy motor. In laboratory tests, we have demonstrated that the trained autoassociator has a small reconstruction error on measurements recorded from healthy motors but a larger error on those recorded from a motor with a fault. We have designed and built a motor monitoring system using an autoassociator for anomaly detection and are in the process of testing the system at three industrial and commercial sites.

The Gamma MLP for Speech Phoneme Recognition

Steve Lawrence, Ah Tsoi, Andrew Back

We define a Gamma multi-layer perceptron (MLP) as an MLP with the usual

l synaptic weights replaced by gamma filters (as proposed by de Vries and Principe (de Vries and Principe, 1992)) and associated gain terms throughout all layers. We derive gradient descent update equations and apply the model to the recognition of speech phonemes. We find that both the inclusion of gamma filters in all layers, and the inclusion of synaptic gains, improves the performance of the Gamma MLP. We compare the Gamma MLP with TDNN, Back-Tsai FIR MLP, and Back-Tsai IIR MLP architectures, and a local approximation scheme. We find that the Gamma MLP results in an substantial reduction in error rates.

Laterally Interconnected Self-Organizing Maps in Hand-Written Digit Recognition
Yoonsuck Choe, Joseph Sirosh, Risto Miikkulainen

An application of laterally interconnected self-organizing maps (LISSOM) to handwritten digit recognition is presented. The lateral connections learn the correlations of activity between units on the map. The resulting excitatory connections focus the activity into local patches and the inhibitory connections decorrelate redundant activity on the map. The map thus forms internal representations that are easy to recognize with e.g. a perceptron network. The recognition rate on a subset of NIST database 3 is 4.0% higher with LISSOM than with a regular Self-Organizing Map (SOM) as the front end, and 15.8% higher than recognition of raw input bit maps directly. These results form a promising starting point for building pattern recognition systems with a LISSOM map as a front end.

Improved Silicon Cochlea using Compatible Lateral Bipolar Transistors

André van Schaik, Eric Fragnière, Eric Vittoz

Analog electronic cochlear models need exponentially scaled filters. CMOS Compatible Lateral Bipolar Transistors (CLBTs) can create exponentially scaled currents when biased using a resistive line with a voltage difference between both ends of the line. Since these CLBTs are independent of the CMOS threshold voltage, current sources implemented with CLBTs are much better matched than current sources created with MOS transistors operated in weak inversion. Measurements from integrated test chips are shown to verify the improved matching.

On the Computational Power of Noisy Spiking Neurons

Wolfgang Maass

It has remained unknown whether one can in principle carry out reliable digital computations with networks of biologically realistic models for neurons. This article presents rigorous constructions for simulating in real-time arbitrary given boolean circuits and finite automata with arbitrarily high reliability by networks of noisy spiking neurons. In addition we show that with the help of "shunting inhibition" even networks of very unreliable spiking neurons can simulate in real-time any McCulloch-Pitts neuron (or "threshold gate"), and therefore any multilayer perceptron (or "threshold circuit") in a reliable manner. These constructions provide a possible explanation for the fact that biological neural systems can carry out quite complex computations within 100 msec. It turns out that the assumption that these constructions require about the shape of the EPSP's and the behaviour of the noise are surprisingly weak.

Parallel Optimization of Motion Controllers via Policy Iteration

Jefferson Coelho, R. Sitaraman, Roderic Grupen

This paper describes a policy iteration algorithm for optimizing the performance of a harmonic function-based controller with respect to a user-defined index. Value functions are represented as potential distribution over the problem domain, being control policies represented as gradient fields over the same domain. All intermediate policies are intrinsically safe, i.e. collisions are not promoted during the adaptation process. The algorithm has efficient implementation in parallel

SIMD architectures. One potential application - travel distance minimization - illustrates its usefulness.

Model Matching and SFMD Computation

Steven Rehfuss, Dan Hammerstrom

In systems that process sensory data there is frequently a model matching stage where class hypotheses are combined to recognize a complex entity. We introduce a new model of parallelism, the Single Function Multiple Data (SFMD) model, appropriate to this stage. SFMD functionality can be added with small hardware expense to certain existing SIMD architectures, and as an incremental addition to the programming model. Adding SFMD to an SIMD machine will not only allow faster model matching, but also increase its flexibility as a general purpose machine and its scope in performing the initial stages of sensory processing.

Learning with ensembles: How overfitting can be useful

Peter Sollich, Anders Krogh

We study the characteristics of learning with ensembles. Solving exactly the simple model of an ensemble of linear students, we find surprisingly rich behaviour. For learning in large ensembles, it is advantageous to use under-regularized students, which actually over-fit the training data. Globally optimal performance can be obtained by choosing the training set sizes of the students appropriately. For smaller ensembles, optimization of the ensemble weights can yield significant improvements in ensemble generalization performance, in particular if the individual students are subject to noise in the training process. Choosing students with a wide range of regularization parameters makes this improvement robust against changes in the unknown level of noise in the training data.

Analog VLSI Processor Implementing the Continuous Wavelet Transform

R. Edwards, Gert Cauwenberghs

We present an integrated analog processor for real-time wavelet decomposition and reconstruction of continuous temporal signals covering the audio frequency range. The processor performs complex harmonic modulation and Gaussian lowpass filtering in 16 parallel channels, each clocked at a different rate, producing a multiresolution mapping on a logarithmic frequency scale. Our implementation uses mixed-mode analog and digital circuits, oversampling techniques, and switched-capacitor filters to achieve a wide linear dynamic range while maintaining compact circuit size and low power consumption. We include experimental results on the processor and characterize its components separately from measurements on a single-channel test chip.

A Novel Channel Selection System in Cochlear Implants Using Artificial Neural Network

Marwan Jabri, Raymond Wang

State-of-the-art speech processors in cochlear implants perform channel selection using a spectral maxima strategy. This strategy can lead to confusions when high frequency features are needed to discriminate between sounds. We present in this paper a novel channel selection strategy based upon pattern recognition which allows "smart" channel selections to be made. The proposed strategy is implemented using multi-layer perceptrons trained on a multi-speaker labelled speech database. The input to the network are the energy coefficients of N energy channels. The output of the system are the indices of the M selected channels. We compare the performance of our proposed system to that of spectral maxima strategy, and show that our strategy can produce significantly better results.

Rapid Quality Estimation of Neural Network Input Representations

Kevin Cherkauer, Jude Shavlik

The choice of an input representation for a neural network can have a profound impact on its accuracy in classifying novel instances. However, neural networks are typically computationally expensive to train, making it difficult to test large numbers of alternative representations. This paper introduces fast quality measures for neural network representations, allowing one to quickly and accurately estimate which of a collection of possible representations for a problem is the best. We show that our measures for ranking representations are more accurate than a previously published measure, based on experiments with three difficult, real-world pattern recognition problems.

Generating Accurate and Diverse Members of a Neural-Network Ensemble

David Opitz, Jude Shavlik

Neural-network ensembles have been shown to be very accurate classification techniques. Previous work has shown that an effective ensemble should consist of networks that are not only highly correct, but ones that make their errors on different parts of the input space as well. Most existing techniques, however, only indirectly address the problem of creating such a set of networks. In this paper we present a technique called ADDEMUP that uses genetic algorithms to directly search for an accurate and diverse set of trained networks. ADDEMUP works by first creating an initial population, then uses genetic operators to continually create new networks, keeping the set of networks that are as accurate as possible while disagreeing with each other as much as possible. Experiments on three DNA problems show that ADDEMUP is able to generate a set of trained networks that is more accurate than several existing approaches. Experiments also show that ADDEMUP is able to effectively incorporate prior knowledge, if available, to improve the quality of its ensemble.

A Model of Spatial Representations in Parietal Cortex Explains Hemineglect

Alexandre Pouget, Terrence J. Sejnowski

We have recently developed a theory of spatial representations in which the position of an object is not encoded in a particular frame of reference but, instead, involves neurons computing basis functions of their sensory inputs. This type of representation is able to perform nonlinear sensorimotor transformations and is consistent with the response properties of parietal neurons. We now ask whether the same theory could account for the behavior of human patients with parietal lesions. These lesions induce a deficit known as hemineglect that is characterized by a lack of reaction to stimuli located in the hemispace contralateral to the lesion. A simulated lesion in a basis function representation was found to replicate three of the most important aspects of hemineglect: i) The models failed to cross the leftmost lines in line cancellation experiments, ii) the deficit affected multiple frames of reference and, iii) it could be object centered. These results strongly support the basis function hypothesis for spatial representations and provide a computational theory of hemineglect at the single cell level.

Temporal Difference Learning in Continuous Time and Space

Kenji Doya

A continuous-time, continuous-state version of the temporal difference (TD) algorithm is derived in order to facilitate the application of reinforcement learning to real-world control tasks and neurobiological modeling. An optimal nonlinear feedback control law was also derived using the derivatives of the value function. The performance of the algorithms was tested in a task of swinging up a pendulum with limited torque. Both the "critic" that specifies the paths to the upright position and the "actor" that works as a nonlinear feedback controller were successfully implemented by radial basis function (RBF) networks.

Tempering Backpropagation Networks: Not All Weights are Created Equal

Nicol Schraudolph, Terrence J. Sejnowski

Terrence J. Sejnowski

Improving Policies without Measuring Merits

Peter Dayan, Satinder Singh

Performing policy iteration in dynamic programming should only require knowledge of relative rather than absolute measures of the utility of actions (Werbos, 1991) - what Baird (1993) calls the advantages of actions at states. Nevertheless, most existing methods in dynamic programming (including Baird's) compute some form of absolute utility function. For smooth problems, advantages satisfy two differential consistency conditions (including the requirement that they be free of curl), and we show that enforcing these can lead to appropriate policy improvement solely in terms of advantages.

Selective Attention for Handwritten Digit Recognition

Ethem Alpaydin

Completely parallel object recognition is NP-complete. Achieving a recognizer with feasible complexity requires a compromise between parallel and sequential processing where a system selectively focuses on parts of a given image, one after another. Successive fixations are generated to sample the image and these samples are processed and abstracted to generate a temporal context in which results are integrated over time. A computational model based on a partially recurrent feedforward network is proposed and made credible by testing on the real-world problem of recognition of handwritten digits with encouraging results.

Learning Model Bias

Jonathan Baxter

In this paper the problem of learning appropriate domain-specific bias is addressed. It is shown that this can be achieved by learning many related tasks from the same domain, and a theorem is given bounding the number of tasks that must be learnt. A corollary of the theorem is that if the tasks are known to possess a common internal representation or preprocessing then the number of examples required per task for good generalization when learning n tasks simultaneously scales like $O(a + \sqrt{n})$, where $O(a)$ is a bound on the minimum number of examples required to learn a single task, and $O(a + b)$ is a bound on the number of examples required to learn each task independently. An experiment providing strong qualitative support for the theoretical results is reported.

Estimating the Bayes Risk from Sample Data

Robert Snapp, Tong Xu

A new nearest-neighbor method is described for estimating the Bayes risk of a multiclass pattern classification problem from sample data (e.g., a classified training set). Although it is assumed that the classification problem can be accurately described by sufficiently smooth class-conditional distributions, neither these distributions, nor the corresponding prior probabilities of the classes are required. Thus this method can be applied to practical problems where the underlying probabilities are not known. This method is illustrated using two different pattern recognition problems.

A Model of Auditory Streaming

Susan McCabe, Michael Denham

An essential feature of intelligent sensory processing is the ability to focus on the part of the signal of interest against a background of distracting signals, and to be able to direct this focus at will. In this paper the problem of auditory scene segmentation is considered and a model of the early stages of the process is proposed. The behaviour of the model is shown to be in agreement with

th a number of well known psychophysical results. The principal contribution of this model lies in demonstrating how streaming might result from interactions between the tonotopic patterns of activity of input signals and traces of previous activity which feedback and influence the way in which subsequent signals are processed.

Context-Dependent Classes in a Hybrid Recurrent Network-HMM Speech Recognition System

Dan Kershaw, Anthony Robinson, Mike Hochberg

A method for incorporating context-dependent phone classes in a connectionist-HMM hybrid speech recognition system is introduced. A modular approach is adopted, where single-layer networks discriminate between different context classes given the phone class and the acoustic data. The context networks are combined with a context-independent (CI) network to generate context-dependent (CD) phone probability estimates. Experiments show an average reduction in word error rate of 16% and 13% from the CI system on ARPA 5,000 word and SQALE 20,000 word tasks respectively. Due to improved modelling, the decoding speed of the CD system is more than twice as fast as the CI system.

Exploiting Tractable Substructures in Intractable Networks

Lawrence Saul, Michael Jordan

We develop a refined mean field approximation for inference and learning in probabilistic neural networks. Our mean field theory, unlike most, does not assume that the units behave as independent degrees of freedom; instead, it exploits in a principled way the existence of large substructures that are computationally tractable. To illustrate the advantages of this framework, we show how to incorporate weak higher order interactions into a first-order hidden Markov model, treating the corrections (but not the first order structure) within mean field theory.

Statistical Theory of Overtraining - Is Cross-Validation Asymptotically Effective?

Shun-ichi Amari, Noboru Murata, Klaus-Robert Müller, Michael Finke, Howard Yang

A statistical theory for overtraining is proposed. The analysis treats realizable stochastic neural networks, trained with Kullback-Leibler loss in the asymptotic case. It is shown that the asymptotic gain in the generalization error is small if we perform early stopping, even if we have access to the optimal stopping time. Considering cross-validation stopping we answer the question: In what ratio the examples should be divided into training and testing sets in order to obtain the optimum performance. In the non-asymptotic region cross-validated early stopping always decreases the generalization error. Our large scale simulations done on a CM5 are in nice agreement with our analytical findings.

Primitive Manipulation Learning with Connectionism

Yoky Matsuoka

Infants' manipulative exploratory behavior within the environment is a vehicle of cognitive stimulation [McCall 1974]. During this time, infants practice and perfect sensorimotor patterns that become behavioral modules which will be serialized and imbedded in more complex actions. This paper explores the development of such primitive learning systems using an embodied lightweight hand which will be used for a humanoid being developed at the MIT Artificial Intelligence Laboratory [Brooks and Stein 1993]. Primitive grasping procedures are learned from sensory inputs using a connectionist reinforcement algorithm while two submodules preprocess sensory data to recognize the hardness of objects and detect shear using competitive learning and back-propagation algorithm strategies, respectively. This system is not only consistent and quick during the initial learning stage, but also adaptable to new situations after training is complete.

leted.

Silicon Models for Auditory Scene Analysis

John Lazzaro, John Wawrzynek

We are developing special-purpose, low-power analog-to-digital converters for speech and music applications, that feature analog circuit models of biological audition to process the audio signal before conversion. This paper describes our most recent converter design, and a working system that uses several copies of the chip to compute multiple representations of sound from an analog input. This multi-representation system demonstrates the plausibility of inexpensively implementing an auditory scene analysis approach to sound processing.

Human Face Detection in Visual Scenes

Henry Rowley, Shumeet Baluja, Takeo Kanade

We present a neural network-based face detection system. A retinally connected neural network examines small windows of an image, and decides whether each window contains a face. The system arbitrates between multiple networks to improve performance over a single network. We use a bootstrap algorithm for training, which adds false detections into the training set as training progresses. This eliminates the difficult task of manually selecting non-face training examples, which must be chosen to span the entire space of non-face images. Comparisons with another state-of-the-art face detection system are presented; our system has better performance in terms of detection and false-positive rates.

SPERT-II: A Vector Microprocessor System and its Application to Large Problems in Backpropagation Training

John Wawrzynek, Krste Asanovic, Brian Kingsbury, James Beck, David Johnson, Nelson Morgan

We report on our development of a high-performance system for neural network and other signal processing applications. We have designed and implemented a vector microprocessor and packed it as an attached processor for a conventional workstation. We present performance comparisons with commercial workstations on neural network backpropagation training. The SPERT-II system demonstrates significant speedups over extensively hand-optimized code running on the workstations.

Pruning with generalization based weight saliencies: λ OBD, λ OBS

Morten Pedersen, Lars Hansen, Jan Larsen

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Stochastic Hillclimbing as a Baseline Method for Evaluating Genetic Algorithms

Ari Juels, Martin Wattenberg

We investigate the effectiveness of stochastic hillclimbing as a baseline for evaluating the performance of genetic algorithms (GAs) as combinatorial function optimizers. In particular, we address two problems to which GAs have been applied in the literature: Koza's 11-multiplexer problem and the jobshop problem. We demonstrate that simple stochastic hillclimbing methods are able to achieve results comparable or superior to those obtained by the GAs designed to address these two problems. We further illustrate, in the case of the jobshop problem, how insights obtained in the formulation of a stochastic hillclimbing algorithm can lead to improvements in the encoding used by a GA.

Using the Future to "Sort Out" the Present: Rankprop and Multitask Learning for Medical Risk Evaluation

Rich Caruana, Shumeet Baluja, Tom Mitchell

A patient visits the doctor; the doctor reviews the patient's history, asks questions, makes basic measurements (blood pressure, . . .), and prescribes tests or treatment. The prescribed course of action is based on an assessment of patient risk-patients at higher risk are given more and faster attention. It is also sequential- it is too expensive to immediately order all tests which might later be of value. This paper presents two methods that together improve the accuracy of backprop nets on a pneumonia risk assessment problem by 10-50%. Rankprop improves on backpropagation with sum of squares error in ranking patients by risk. Multitask learning takes advantage of future lab tests available in the training set, but not available in practice when predictions must be made. Both methods are broadly applicable.

Exponentially many local minima for single neurons

Peter Auer, Mark Herbster, Manfred K. K. Warmuth

We show that for a single neuron with the logistic function as the transfer function the number of local minima of the error function based on the square loss can grow exponentially in the dimension.

An Information-theoretic Learning Algorithm for Neural Network Classification

David J. Miller, Ajit Rao, Kenneth Rose, Allen Gersho

A new learning algorithm is developed for the design of statistical classifiers minimizing the rate of misclassification. The method, which is based on ideas from information theory and analogies to statistical physics, assigns data to classes in probability. The distributions are chosen to minimize the expected classification error while simultaneously enforcing the classifier's structure and a level of "randomness" measured by Shannon's entropy. Achievement of the classifier structure is quantified by an associated cost. The constrained optimization problem is equivalent to the minimization of a Helmholtz free energy, and the resulting optimization method is a basic extension of the deterministic annealing algorithm that explicitly enforces structural constraints on assignments while reducing the entropy and expected cost with temperature. In the limit of low temperature, the error rate is minimized directly and a hard classifier with the requisite structure is obtained. This learning algorithm can be used to design a variety of classifier structures.

The approach is compared with standard methods for radial basis function design and is demonstrated to substantially outperform other design methods on several benchmark examples, while retaining design complexity comparable to, or only moderately greater than that of strict descent-based methods.

Improving Elevator Performance Using Reinforcement Learning

Robert Crites, Andrew Barto

This paper describes the application of reinforcement learning (RL) to the difficult real world problem of elevator dispatching. The elevator domain poses a combination of challenges not seen in most RL research to date. Elevator systems operate in continuous state spaces and in continuous time as discrete event dynamic systems. Their states are not fully observable and they are nonstationary due to changing passenger arrival rates. In addition, we use a team of RL agents, each of which is responsible for controlling one elevator car. The team receives a global reinforcement signal which appears noisy to each agent due to the effects of the actions of the other agents, the random nature of the arrivals and the incomplete observation of the state. In spite of these complications, we show results that in simulation surpass the best of the heuristic elevator control algorithms of which we are aware. These results demonstrate the power of RL on a very large scale stochastic dynamic optimization problem of practical utility.

Optimal Asset Allocation using Adaptive Dynamic Programming

Ralph Neuneier

In recent years, the interest of investors has shifted to computerized asset allocation (portfolio management) to exploit the growing dynamics of the capital markets. In this paper, asset allocation is formalized as a Markovian Decision Problem which can be optimized by applying dynamic programming or reinforcement learning based algorithms. Using an artificial exchange rate, the asset allocation strategy optimized with reinforcement learning (Q-Learning) is shown to be equivalent to a policy computed by dynamic programming. The approach is then tested on the task to invest liquid capital in the German stock market. Here, neural networks are used as value function approximators. The resulting asset allocation strategy is superior to a heuristic benchmark policy. This is a further example which demonstrates the applicability of neural network based reinforcement learning to a problem setting with a high dimensional state space.

Recursive Estimation of Dynamic Modular RBF Networks

Visakan Kadirkamanathan, Maha Kadirkamanathan

In this paper, recursive estimation algorithms for dynamic modular networks are developed. The models are based on Gaussian RBF networks and the gating network is considered in two stages: At first, it is simply a time-varying scalar and in the second, it is based on the state, as in the mixture of local experts scheme. The resulting algorithm uses Kalman filter estimation for the model estimation and the gating probability estimation. Both, 'hard' and 'soft' competition based estimation schemes are developed where in the former, the most probable network is adapted and in the latter all networks are adapted by appropriate weighting of the data.

A Neural Network Classifier for the I100 OCR Chip

John Platt, Timothy Allen

This paper describes a neural network classifier for the I1000 chip, which optically reads the E13B font characters at the bottom of checks. The first layer of the neural network is a hardware linear classifier which recognizes the characters in this font. A second software neural layer is implemented on an inexpensive microprocessor to clean up the results of the first layer. The hardware linear classifier is mathematically specified using constraints and an optimization principle. The weights of the classifier are found using the active set method, similar to Vapnik's separating hyperplane algorithm. In 7.5 minutes on a SPARC 2 time, the method solves for 1523 Lagrange multipliers, which is equivalent to training on a data set of approximately 128,000 examples. The resulting network performs quite well: when tested on a test set of 1500 real checks, it has a 99.995% character accuracy rate.

Modeling Saccadic Targeting in Visual Search

Rajesh Rao, Gregory Zelinsky, Mary Hayhoe, Dana Ballard

Visual cognition depends critically on the ability to make rapid eye movements known as saccades that orient the fovea over targets of interest in a visual scene. Saccades are known to be ballistic: the pattern of muscle activation for foveating a prespecified target location is computed prior to the movement and visual feedback is precluded. Despite these distinctive properties, there has been no general model of the saccadic targeting strategy employed by the human visual system during visual search in natural scenes. This paper proposes a model for saccadic targeting that uses iconic scene representations derived from oriented spatial filters at multiple scales. Visual search proceeds in a coarse-to-fine fashion with the largest scale filter responses being compared first. The model was empirically tested by comparing its performance with actual eye movement data from human subjects in a natural visual search task; preliminary results

lts indicate substantial agreement between eye movements predicted by the model and those recorded from human subjects.

Plasticity of Center-Surround Opponent Receptive Fields in Real and Artificial Neural Systems of Vision

S. Yasui, T. Furukawa, M. Yamada, T. Saito

Despite the phylogenic and structural differences, the visual systems of different species, whether vertebrate or invertebrate, share certain functional properties. The center-surround opponent receptive field (CSRF) mechanism represents one such example. Here, analogous CSRFs are shown to be formed in an artificial neural network which learns to localize contours (edges) of the luminance difference. Furthermore, when the input pattern is corrupted by a background noise, the CSRFs of the hidden units becomes shallower and broader with decrease of the signal-to-noise ratio (SNR). The same kind of SNR-dependent plasticity is present in the CSRF of real visual neurons; in bipolar cells of the carp retina as is shown here experimentally, as well as in large monopolar cells of the fly compound eye as was described by others. Also, analogous SNR-dependent plasticity is shown to be present in the biphasic flash responses (BPFR) of these artificial and biological visual systems. Thus, the spatial (CSRF) and temporal (BPFR) filtering properties which a wide variety of creatures see the world appear to be optimized for detectability of changes in space and time.

Active Gesture Recognition using Learned Visual Attention

Trevor Darrell, Alex Pentland

We have developed a foveated gesture recognition system that runs in an untrained office environment with an active camera. Using vision routines previously implemented for an interactive environment, we determine the spatial location of salient body parts of a user and guide an active camera to obtain images of gestures or expressions. A hidden-state reinforcement learning paradigm is used to implement visual attention. The attention module selects targets to foveate based on the goal of successful recognition, and uses a new multiple-model Q-learning formulation. Given a set of target and distractor gestures, our system can learn where to foveate to maximally discriminate a particular gesture.

On Neural Networks with Minimal Weights

Vasken Bohossian, Jehoshua Bruck

Linear threshold elements are the basic building blocks of artificial neural networks. A linear threshold element computes a function that is a sign of a weighted sum of the input variables. The weights are arbitrary integers; actually, they can be very big integers-exponential in the number of the input variables. However, in practice, it is difficult to implement big weights. In the present literature a distinction is made between the two extreme cases: linear threshold functions with polynomial-size weights as opposed to those with exponential-size weights. The main contribution of this paper is to fill up the gap by further refining that separation. Namely, we prove that the class of linear threshold functions with polynomial-size weights can be divided into subclasses according to the degree of the polynomial. In fact, we prove a more general result: there exists a minimal weight linear threshold function for any arbitrary number of inputs and any weight size. To prove those results we have developed a novel technique for constructing linear threshold functions with minimal weights.

Neural Networks with Quadratic VC Dimension

Pascal Koiran, Eduardo Sontag

This paper shows that neural networks which use continuous activation functions have VC dimension at least as large as the square of the number

of weights w . This result settles a long-standing open question, namely whether the well-known $O(w \log w)$ bound, known for hard-threshold nets, also held for more general sigmoidal nets. Implications for the number of samples needed for valid generalization are discussed.

Factorial Hidden Markov Models

Zoubin Ghahramani, Michael Jordan

We present a framework for learning in hidden Markov models with distributed state representations. Within this framework, we derive a learning algorithm based on the Expectation-Maximization (EM) procedure for maximum likelihood estimation. Analogous to the standard Baum-Welch update rules, the M-step of our algorithm is exact and can be solved analytically. However, due to the combinatorial nature of the hidden state representation, the exact E-step is intractable. A simple and tractable mean field approximation is derived. Empirical results on a set of problems suggest that both the mean field approximation and Gibbs sampling are viable alternatives to the computationally expensive exact algorithm.

Extracting Tree-Structured Representations of Trained Networks

Mark Craven, Jude Shavlik

A significant limitation of neural networks is that the representations they learn are usually incomprehensible to humans. We present a novel algorithm, TREPAN, for extracting comprehensible, symbolic representations from trained neural networks. Our algorithm uses queries to induce a decision tree that approximates the concept represented by a given network. Our experiments demonstrate that TREPAN is able to produce decision trees that maintain a high level of fidelity to their respective networks while being comprehensible and accurate. Unlike previous work in this area, our algorithm is general in its applicability and scales well to large networks and problems with high-dimensional input spaces.

Improving Committee Diagnosis with Resampling Techniques

Bambang Parmanto, Paul Munro, Howard Doyle

Central to the performance improvement of a committee relative to individual networks is the error correlation between networks in the committee. We investigated methods of achieving error independence between the networks by training the networks with different resampling sets from the original training set. The methods were tested on the sinwave artificial task and the real-world problems of hepatoma (liver cancer) and breast cancer diagnoses.

The Capacity of a Bump

Gary Flake

Recently, several researchers have reported encouraging experimental results when using Gaussian or bump-like activation functions in multilayer perceptrons. Networks of this type usually require fewer hidden layers and units and often learn much faster than typical sigmoidal networks. To explain these results we consider a hyper-ridge network, which is a simple perceptron with no hidden units and a ridge activation function. If we are interested in partitioning points in d dimensions into two classes then in the limit as d approaches infinity the capacity of a hyper-ridge and a perceptron is identical. However, we show that for $p \sim d$, which is the usual case in practice, the ratio of hyper-ridge to perceptron dichotomies approaches $pl_2(d + 1)$.

Boosting Decision Trees

Harris Drucker, Corinna Cortes

We introduce a constructive, incremental learning system for regression problems that models data by means of locally linear experts. In contrast to other approaches, the experts are trained independently and do not compete for data during learning. Only when a prediction for a query is required do

the experts cooperate by blending their individual predictions.

Each expert is trained by minimizing a penalized local cross validation error using second order methods. In this way, an expert is able to find a local distance metric by adjusting the size and shape of the receptive field in which its predictions are valid, and also to detect relevant input features by adjusting its bias on the importance of individual input dimensions. We derive asymptotic results for our method. In a variety of simulations the properties of the algorithm are demonstrated with respect to interference, learning speed, prediction accuracy, feature detection, and task oriented incremental learning.

Some results on convergent unlearning algorithm

Serguei Semenov, Irina Shuvalova

In this paper we consider probabilities of different asymptotics of convergent unlearning algorithm for the Hopfield-type neural network work (Plakhov & Semenov, 1994) treating the case of unbiased random patterns. We show also that failed unlearning results in total memory breakdown.

Forward-backward retraining of recurrent neural networks

Andrew W. Senior, Anthony Robinson

This paper describes the training of a recurrent neural network as the letter posterior probability estimator for a hidden Markov model, off-line handwriting recognition system. The network estimates posterior distributions for each of a series of frames representing sections of a handwritten word. The supervised training algorithm, backpropagation through time, requires target outputs to be provided for each frame. Three methods for deriving these targets are presented. A novel method based upon the forward-backward algorithm is found to result in the recognizer with the lowest error rate.

KODAK IMAGELINK™ OCR Alphanumeric Handprint Module

Alexander Shustorovich, Christopher W. Thrasher

This paper describes the Kodak Imageliok™ OCR alphanumeric handprint module. There are two neural network algorithms at its core: the first network is trained to find individual characters in an alphanumeric field, while the second one performs the classification. Both networks were trained on Gabor projections of the original pixel images, which resulted in higher recognition rates and greater noise immunity. Compared to its purely numeric counterpart (Shustorovich and Thrasher, 1995), this version of the system has a significant application specific postprocessing module. The system has been implemented in specialized parallel hardware, which allows it to run at 80 char/sec/board. It has been installed at the Driver and Vehicle Licensing Agency (DVLA) in the United Kingdom. and its overall success rate exceeds 96% (character level without rejects) which translates into 85% field rate. If approximately 20% of the fields are rejected. the system achieves 99.8% character and 99.5% field success rate.

Predictive Q-Routing: A Memory-based Reinforcement Learning Approach to Adaptive Traffic Control

Samuel Choi, Dit-Yan Yeung

In this paper, we propose a memory-based Q-learning algorithm called predictive Q-routing (PQ-routing) for adaptive traffic control. We attempt to address two problems encountered in Q-routing (Boyan & Littman, 1994), namely, the inability to fine-tune routing policies under low network load and the inability to learn new optimal policies under decreasing load conditions. Unlike other memory-based reinforcement learning algorithms in which memory is used to keep past experiences to increase learning speed, PQ-routing keeps the best experiences

learned and reuses them by predicting the traffic trend. The effectiveness of PQ-routing has been verified under various network topologies and traffic conditions. Simulation results show that PQ-routing is superior to Q-routing in terms of both learning speed and adaptability.

High-Performance Job-Shop Scheduling With A Time-Delay TD(λ) Network

Wei Zhang, Thomas Dietterich

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Family Discovery

Stephen Omohundro

"Family discovery" is the task of learning the dimension and structure of a parameterized family of stochastic models. It is especially appropriate when the training examples are partitioned into "episodes" of samples drawn from a single parameter value. We present three family discovery algorithms based on surface learning and show that they significantly improve performance over two alternatives on a parameterized classification task.

Modeling Interactions of the Rat's Place and Head Direction Systems

A. Redish, David Touretzky

We have developed a computational theory of rodent navigation that includes analogs of the place cell system, the head direction system, and path integration. In this paper we present simulation results showing how interactions between the place and head direction systems can account for recent observations about hippocampal place cell responses to doubling and/or rotation of cue cards in a cylindrical arena (Sharp et al., 1990).

Empirical Entropy Manipulation for Real-World Problems

Paul Viola, Nicol Schraudolph, Terrence J. Sejnowski

No finite sample is sufficient to determine the density, and therefore the entropy, of a signal directly. Some assumption about either the functional form of the density or about its smoothness is necessary. Both amount to a prior over the space of possible density functions. By far the most common approach is to assume that the density has a parametric form.

A Neural Network Model of 3-D Lightness Perception

Luiz Pessoa, William Ross

A neural network model of 3-D lightness perception is presented which builds upon the FACADE Theory Boundary Contour System/Feature Contour System of Grossberg and colleagues. Early ratio encoding by retinal ganglion neurons as well as psychophysical results on constancy across different backgrounds (background constancy) are used to provide functional constraints to the theory and suggest a contrast negation hypothesis which states that ratio measures between coplanar regions are given more weight in the determination of lightness of the respective regions.

Simulations of the model address data on lightness perception, including the coplanar ratio hypothesis, the Benary cross, and White's illusion.

Adaptive Retina with Center-Surround Receptive Field

Shih-Chii Liu, Kwabena Boahen

Both vertebrate and invertebrate retinas are highly efficient in extracting contrast independent of the background intensity over five or more decades. This efficiency has been rendered possible by the adaptation of the DC operating point to the background intensity while maintaining high gain transient responses. The center-surround properties of the retina allows the system to extract information at the edges in the image

e. This silicon retina models the adaptation properties of the receptors and the antagonistic center-surround properties of the laminar cells of the invertebrate retina and the outer-plexiform layer of the vertebrate retina. We also illustrate the spatio-temporal responses of the silicon retina on moving bars. The chip has 59x64 pixels on a 6.9x6.8mm² die and it is fabricated in 2 J- μ m n-well technology.

Does the Wake-sleep Algorithm Produce Good Density Estimators?

Brendan J. Frey, Geoffrey E. Hinton, Peter Dayan

The wake-sleep algorithm (Hinton, Dayan, Frey and Neal 1995) is a relatively efficient method of fitting a multilayer stochastic generative model to high-dimensional data. In addition to the top-down connections in the generative model, it makes use of bottom-up connections for approximating the probability distribution over the hidden units given the data, and it trains these bottom-up connections using a simple delta rule. We use a variety of synthetic and real data sets to compare the performance of the wake-sleep algorithm with Monte Carlo and mean field methods for fitting the same generative model and also compare it with other models that are less powerful but easier to fit.

EM Optimization of Latent-Variable Density Models

Christopher Bishop, Markus Svensén, Christopher Williams

There is currently considerable interest in developing general non-linear density models based on latent, or hidden, variables. Such models have the ability to discover the presence of a relatively small number of underlying 'causes' which, acting in combination, give rise to the apparent complexity of the observed data set. Unfortunately, to train such models generally requires large computational effort. In this paper we introduce a novel latent variable algorithm which retains the general non-linear capabilities of previous models but which uses a training procedure based on the EM algorithm. We demonstrate the performance of the model on a toy problem and on data from flow diagnostics for a multi-phase oil pipeline.

Modern Analytic Techniques to Solve the Dynamics of Recurrent Neural Networks

A.C.C. Coolen, S. Laughton, D. Sherrington

We describe the use of modern analytical techniques in solving the dynamics of symmetric and nonsymmetric recurrent neural networks near saturation. These explicitly take into account the correlations between the post-synaptic potentials, and thereby allow for a reliable prediction of transients.

A Realizable Learning Task which Exhibits Overfitting

Siegfried BöS

In this paper we examine a perceptron learning task. The task is realizable since it is provided by another perceptron with identical architecture. Both perceptrons have nonlinear sigmoid output functions. The gain of the output function determines the level of nonlinearity of the learning task. It is observed that a high level of nonlinearity leads to overfitting. We give an explanation for this rather surprising observation and develop a method to avoid the overfitting. This method has two possible interpretations, one is learning with noise, the other cross-validated early stopping.

Simulation of a Thalamocortical Circuit for Computing Directional Heading in the Rat

Hugh Blair

Several regions of the rat brain contain neurons known as head-direction cells, which encode the animal's directional heading during spatial navigation. This paper presents a biophysical model of head-direction cells

ll activity, which suggests that a thalamocortical circuit might compute the rat's head direction by integrating the angular velocity of the head over time. The model was implemented using the neural simulator NEURON, and makes testable predictions about the structure and function of the rat head-direction circuit.

Clustering data through an analogy to the Potts model

Marcelo Blatt, Shai Wiseman, Eytan Domany

A new approach for clustering is proposed. This method is based on an analogy to a physical model; the ferromagnetic Potts model at thermal equilibrium is used as an analog computer for this hard optimization problem. We do not assume any structure of the underlying distribution of the data. Phase space of the Potts model is divided into three regions; ferromagnetic, super-paramagnetic and paramagnetic phases. The region of interest is that corresponding to the super-paramagnetic one, where domains of aligned spins appear. The range of temperatures where these structures are stable is indicated by a non-vanishing magnetic susceptibility. We use a very efficient Monte Carlo algorithm to measure the susceptibility and the spin spin correlation function. The values of the spin spin correlation function, at the super-paramagnetic phase, serve to identify the partition of the data points into clusters.

Learning Fine Motion by Markov Mixtures of Experts

Marina Meila, Michael Jordan

Compliant control is a standard method for performing fine manipulation tasks, like grasping and assembly, but it requires estimation of the state of contact (s.o.c.) between the robot arm and the objects involved. Here we present a method to learn a model of the movement from measured data. The method requires little or no prior knowledge and the resulting model explicitly estimates the s.o.c. The current s.o.c. is viewed as the hidden state variable of a discrete HMM. The control dependent transition probabilities between states are modeled as parameterized functions of the measurement. We show that their parameters can be estimated from measurements at the same time as the parameters of the movement in each s.o.c. The learning algorithm is a variant of the EM procedure. The E step is computed exactly; solving the M step exactly is not possible in general. Here, gradient ascent is used to produce an increase in likelihood.

Stable Linear Approximations to Dynamic Programming for Stochastic Control Problems with Local Transitions

Benjamin Van Roy, John Tsitsiklis

We consider the solution to large stochastic control problems by means of methods that rely on compact representations and a variant of the value iteration algorithm to compute approximate cost-to-go functions. While such methods are known to be unstable in general, we identify a new class of problems for which convergence, as well as graceful error bounds, are guaranteed. This class involves linear parameterizations of the cost-to-go function together with an assumption that the dynamic programming operator is a contraction with respect to the Euclidean norm when applied to functions in the parameterized class. We provide a special case where this assumption is satisfied, which relies on the locality of transitions in a state space. Other cases will be discussed in a full length version of this paper.

Generalisation of A Class of Continuous Neural Networks

John Shawe-Taylor, Jieyu Zhao

We propose a way of using boolean circuits to perform real valued computation in a way that naturally extends their boolean functionality. The functionality of multiple fan in threshold gates in this model

is shown to mimic that of a hardware implementation of continuous Neural Networks. A Vapnik-Chervonenkis dimension and sample size analysis for the systems is performed giving best known sample sizes for a real valued Neural Network. Experimental results confirm the conclusion that the sample sizes required for the networks are significantly smaller than for sigmoidal networks.

Visual gesture-based robot guidance with a modular neural system

Enno Littmann, Andrea Drees, Helge Ritter

We report on the development of the modular neural system "SEE EAGLE" for the visual guidance of robot pick-and-place actions. Several neural networks are integrated to a single system that usually recognizes human hand pointing gestures from stereo pairs of color video images. The output of the hand recognition stage is processed by a set of color-sensitive neural networks to determine the cartesian location of the target object that is referenced by the pointing gesture. Finally, this information is used to guide a robot to grab the target object and put it at another location that can be specified by a second pointing gesture. The accuracy of the current system allows to identify the location of the referenced target object to an accuracy of 1 cm in a workspace area of 50x50 cm. In our current environment, this is sufficient to pick and place arbitrarily positioned target objects within the workspace. The system consists of neural networks that perform the tasks of image segmentation, estimation of hand location, estimation of 3D-pointing direction, object recognition, and necessary coordinate transforms. Drawing heavily on the use of learning algorithms, the functions of all network modules were created from data examples only.

Symplectic Nonlinear Component Analysis

Lucas Parra

Statistically independent features can be extracted by finding a factorial representation of a signal distribution. Principal Component Analysis (PCA) accomplishes this for linear correlated and Gaussian distributed signals. Independent Component Analysis (ICA), formalized by Comon (1994), extracts features in the case of linear statistical dependent but not necessarily Gaussian distributed signals. Nonlinear Component Analysis finally should find a factorial representation for nonlinear statistical dependent distributed signals. This paper proposes for this task a novel feed-forward, information conserving, nonlinear map - the explicit symplectic transformations. It also solves the problem of non-Gaussian output distributions by considering single coordinate higher order statistics.

Independent Component Analysis of Electroencephalographic Data

Scott Makeig, Anthony Bell, Tzyy-Ping Jung, Terrence J. Sejnowski

Because of the distance between the skull and brain and their different resistivities, electroencephalographic (EEG) data collected from any point on the human scalp includes activity generated within a large brain area. This spatial smearing of EEG data by volume conduction does not involve significant time delays, however, suggesting that the Independent Component Analysis (ICA) algorithm of Bell and Sejnowski [1] is suitable for performing blind source separation on EEG data. The ICA algorithm separates the problem of source identification from that of source localization. First results of applying the ICA algorithm to EEG and event-related potential (ERP) data collected during a sustained auditory detection task show: (1) ICA training is insensitive to different random seeds. (2) ICA may be used to segregate obvious artifactual EEG components (line and muscle noise, eye movements) from other sources. (3) ICA is capable of isolating overlapping EEG phenomena, including alpha and theta bursts and spatially-separable ERP components, to separate ICA channels. (4) Nonstationarity

ies in EEG and behavioral state can be tracked using ICA via changes in the amount of residual correlation between ICA-filtered output channels.

Competence Acquisition in an Autonomous Mobile Robot using Hardware Neural Techniques

Geoffrey Jackson, Alan Murray

In this paper we examine the practical use of hardware neural networks in an autonomous mobile robot. We have developed a hardware neural system based around a custom VLSI chip, EPICILON III, designed specifically for embedded hardware neural applications. We present here a demonstration application of an autonomous mobile robot that highlights the flexibility of this system. This robot gains basic mobility competence in very few training epochs using an "instinct-rule" training methodology.

Gaussian Processes for Regression

Christopher Williams, Carl Rasmussen

The Bayesian analysis of neural networks is difficult because a simple prior over weights implies a complex prior distribution over functions. In this paper we investigate the use of Gaussian process priors over functions, which permit the predictive Bayesian analysis for fixed values of hyperparameters to be carried out exactly using matrix operations. Two methods, using optimization and averaging (via Hybrid Monte Carlo) over hyperparameters have been tested on a number of challenging problems and have produced excellent results.

Reorganisation of Somatosensory Cortex after Tactile Training

Rasmus Petersen, John Taylor

Topographic maps in primary areas of mammalian cerebral cortex reorganise as a result of behavioural training. The nature of this reorganisation seems consistent with the behaviour of competitive neural networks, as has been demonstrated in the past by computer simulation. We model tactile training on the hand representation in primate somatosensory cortex, using the Neural Field Theory of Amari and his colleagues. Expressions for changes in both receptive field size and magnification factor are derived, which are consistent with owl monkey experiments and make a prediction which goes beyond them.

Prediction of Beta Sheets in Proteins

Anders Krogh, Soren Riis

Most current methods for prediction of protein secondary structure use a small window of the protein sequence to predict the structure of the central amino acid. We describe a new method for prediction of the non-local structure called β -sheet, which consists of two or more β -strands that are connected by hydrogen bonds. Since β -strands are often widely separated in the protein chain, a network with two windows is introduced. After training on a set of proteins the network predicts the sheets well, but there are many false positives. By using a global energy function the β -sheet prediction is combined with a local prediction of the three secondary structures α -helix, β -strand and coil. The energy function is minimized using simulated annealing to give a final prediction.

Active Learning in Multilayer Perceptrons

Kenji Fukumizu

We propose an active learning method with hidden-unit reduction, which is devised specially for multilayer perceptrons (MLP). First, we review our active learning method, and point out that many Fisher-information-based methods applied to MLP have a critical problem: the information matrix may be singular. To solve this problem, we derive the singularity condition of

f an information matrix, and propose an active learning technique that is applicable to MLP. Its effectiveness is verified through experiments.

A model of transparent motion and non-transparent motion aftereffects

Alexander Grunewald

A model of human motion perception is presented. The model contains two stages of direction selective units. The first stage contains broadly tuned units, while the second stage contains units that are narrowly tuned. The model accounts for the motion aftereffect through adapting units at the first stage and inhibitory interactions at the second stage. The model explains how two populations of dots moving in slightly different directions are perceived as a single population moving in the direction of the vector sum, and how two populations moving in strongly different directions are perceived as transparent motion. The model also explains why the motion aftereffect in both cases appears as non-transparent motion.

Improved Gaussian Mixture Density Estimates Using Bayesian Penalty Terms and Network Averaging

Dirk Ormoneit, Volker Tresp

Volker Tresp

A Practical Monte Carlo Implementation of Bayesian Learning

Carl Rasmussen

A practical method for Bayesian training of feed-forward neural networks using sophisticated Monte Carlo methods is presented and evaluated. In reasonably small amounts of computer time this approach outperforms other state-of-the-art methods on 5 data limited tasks from real world domains.

A Smoothing Regularizer for Recurrent Neural Networks

Lizhong Wu, John Moody

We derive a smoothing regularizer for recurrent network models by requiring robustness in prediction performance to perturbations of the training data. The regularizer can be viewed as a generalization of the first order Tikhonov stabilizer to dynamic models. The closed-form expression of the regularizer covers both time-lagged and simultaneous recurrent nets, with feed forward nets and one layer linear nets as special cases. We have successfully tested this regularizer in a number of case studies and found that it performs better than standard quadratic weight decay.

REMAP: Recursive Estimation and Maximization of A Posteriori Probabilities - Application to Transition-Based Connectionist Speech Recognition

Yochai Konig, Hervé Bourlard, Nelson Morgan

In this paper, we introduce REMAP, an approach for the training and estimation of posterior probabilities using a recursive algorithm that is reminiscent of the EM-based Forward-Backward (Liporace 1982) algorithm for the estimation of sequence likelihoods. Although very general, the method is developed in the context of a statistical model for transition-based speech recognition using Artificial Neural Networks (ANN) to generate probabilities for Hidden Markov Models (HMMs). In the new approach, we use local conditional posterior probabilities of transitions to estimate global posterior probabilities of word sequences. Although we still use ANNs to estimate posterior probabilities, the network is trained with targets that are themselves estimates of local posterior probabilities. An initial experimental result shows a significant decrease in error-rate in comparison to a baseline system.

Harmony Networks Do Not Work

René Gourley

Harmony networks have been proposed as a means by which connectionist models can perform symbolic computation. Indeed, proponents claim that a harmony network can be built that constructs parse trees for strings in a context free language. This paper shows that harmony networks do not work in the following sense: they construct many outputs that are not valid parse trees.

Learning Sparse Perceptrons

Jeffrey Jackson, Mark Craven

We introduce a new algorithm designed to learn sparse perceptrons over input representations which include high-order features. Our algorithm, which is based on a hypothesis-boosting method, is able to PAC-learn a relatively natural class of target concepts. Moreover, the algorithm appears to work well in practice: on a set of three problem domains, the algorithm produces classifiers that utilize small numbers of features yet exhibit good generalization performance. Perhaps most importantly, our algorithm generates concept descriptions that are easy for humans to understand.

The Role of Activity in Synaptic Competition at the Neuromuscular Junction

Samuel Joseph, David Willshaw

An extended version of the dual constraint model of motor endplate morphogenesis is presented that includes activity dependent and independent competition. It is supported by a wide range of recent neurophysiological evidence that indicates a strong relationship between synaptic efficacy and survival. The computational model is justified at the molecular level and its predictions match the developmental and regenerative behaviour of real synapses.

Absence of Cycles in Symmetric Neural Networks

Xin Wang, Arun Jagota, Fernanda Botelho, Max Garzon

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Explorations with the Dynamic Wave Model

Thomas Rebotier, Jeffrey Elman

Following Shrager and Johnson (1995) we study growth of local function complexity in a network swept by two overlapping waves: one of pruning, and the other of Hebbian reinforcement of connections. Results indicate a significant spatial gradient in the appearance of both linearly separable and non linearly separable functions of the two inputs of the network; the n.l.s. cells are much sparser and their slope of appearance is sensitive to parameters in a highly non-linear way.

Generalization in Reinforcement Learning: Successful Examples Using Sparse Coarse Coding

Richard S. Sutton

On large problems, reinforcement learning systems must use parameterized function approximators such as neural networks in order to generalize between similar situations and actions. In these cases there are no strong theoretical results on the accuracy of convergence, and computational results have been mixed. In particular, Boyan and Moore reported at last year's meeting a series of negative results in attempting to apply dynamic programming together with function approximation to simple control problems with continuous state spaces. In this paper, we present positive results for all the control tasks they attempted, and for one that is significantly larger. The most important differences are that we used sparse-coarse-coded function approximators (CMACs) whereas they used mostly global function approximators, and that we learned online whereas they learned offline. Bo

yan and Moore and others have suggested that the problems they encountered could be solved by using actual outcomes ("rollouts"), as in classical Monte Carlo methods, and as in the TD(.) algorithm when $\gamma = 1$. However, in our experiments this always resulted in substantially poorer performance. We conclude that reinforcement learning can work robustly in conjunction with function approximators, and that there is little justification at present for avoiding the case of general γ .

A Multiscale Attentional Framework for Relaxation Neural Networks

Dimitris Tsioutsias, Eric Mjolsness

We investigate the optimization of neural networks governed by general objective functions. Practical formulations of such objectives are notoriously difficult to solve; a common problem is the poor local extrema that result by any of the applied methods. In this paper, a novel framework is introduced for the solution of large scale optimization problems. It assumes little about the objective function and can be applied to general nonlinear, non-convex functions; objectives in thousands of variables are thus efficiently minimized by a combination of techniques - deterministic annealing, multiscale optimization, attention mechanisms and trust region optimization methods.

VLSI Model of Primate Visual Smooth Pursuit

Ralph Etienne-Cummings, Jan Van der Spiegel, Paul Mueller

A one dimensional model of primate smooth pursuit mechanism has been implemented in 2 μ m CMOS VLSI. The model consolidates Robinson's negative feedback model with Wyatt and Pola's positive feedback scheme, to produce a smooth pursuit system which zeroes the velocity of a target on the retina. Furthermore, the system uses the current eye motion as a predictor for future target motion. Analysis, stability and biological correspondence of the system are discussed. For implementation at the focal plane, a local correlation based visual motion detection technique is used. Velocity measurements, ranging over 4 orders of magnitude with $< 15\%$ variation, provides the input to the smooth pursuit system. The system performed successful velocity tracking for high contrast scenes. Circuit design and performance of the complete smooth pursuit system is presented.

Reinforcement Learning by Probability Matching

Philip N. Sabes, Michael Jordan

We present a new algorithm for associative reinforcement learning. The algorithm is based upon the idea of matching a network's output probability with a probability distribution derived from the environment's reward signal. This Probability Matching algorithm is shown to perform faster and be less susceptible to local minima than previously existing algorithms. We use Probability Matching to train mixture of experts networks, an architecture for which other reinforcement learning rules fail to converge reliably on even simple problems. This architecture is particularly well suited for our algorithm as it can compute arbitrarily complex functions yet calculation of the output probability is simple.

Generalized Learning Vector Quantization

Atsushi Sato, Keiji Yamada

We propose a new learning method, "Generalized Learning Vector Quantization (GLVQ)," in which reference vectors are updated based on the steepest descent method in order to minimize the cost function. The cost function is determined so that the obtained learning rule satisfies the convergence condition. We prove that Kohonen's rule as used in LVQ does not satisfy the convergence condition and thus degrades recognition ability. Experimental results for printed Chinese character recognition reveal that GLVQ is superior to LVQ in recognition ability.

Bayesian Methods for Mixtures of Experts

Steve Waterhouse, David MacKay, Anthony Robinson

We present a Bayesian framework for inferring the parameters of a mixture of experts model based on ensemble learning by variational free energy minimisation. The Bayesian approach avoids the over-fitting and noise level under-estimation problems of traditional maximum likelihood inference.

We demonstrate these methods on artificial problems and sunspot time series prediction.

High-Speed Airborne Particle Monitoring Using Artificial Neural Networks

Alistair Ferguson, Theo Sabisch, Paul Kaye, Laurence Dixon, Hamid Bolouri

Current environmental monitoring systems assume particles to be spherical, and do not attempt to classify them. A laser-based system developed at the University of Hertfordshire aims at classifying airborne particles through the generation of two-dimensional scattering profiles. The performances of template matching, and two types of neural network (HyperNet and semi-linear units) are compared for image classification. The neural network approach is shown to be capable of comparable recognition performance, while offering a number of advantages over template matching.

Adaptive Mixture of Probabilistic Transducers

Yoram Singer

We introduce and analyze a mixture model for supervised learning of probabilistic transducers. We devise an online learning algorithm that efficiently infers the structure and estimates the parameters of each model in the mixture. Theoretical analysis and comparative simulations indicate that the learning algorithm tracks the best model from an arbitrarily large (possibly infinite) pool of models. We also present an application of the model for inducing a noun phrase recognizer.

SEEMORE: A View-Based Approach to 3-D Object Recognition Using Multiple Visual Cues

Bartlett Mel

A neurally-inspired visual object recognition system is described called SEEMORE, whose goal is to identify common objects from a large known set-independent of 3-D viewing angle, distance, and non-rigid distortion. SEEMORE's database consists of 100 objects that are rigid (shovel), non-rigid (telephone cord), articulated (book), statistical (shrubbery), and complex (photographs of scenes). Recognition results were obtained using a set of 102 color and shape feature channels within a simple feedforward network architecture. In response to a test set of 600 novel test views (6 of each object) presented individually in color video images, SEEMORE identified the object correctly 97% of the time (chance is 1%) using a nearest neighbor classifier. Similar levels of performance were obtained for the subset of 15 non-rigid objects. Generalization behavior reveals emergence of striking natural category structure not explicit in the input feature dimensions.

Using Pairs of Data-Points to Define Splits for Decision Trees

Geoffrey E. Hinton, Michael Revow

Conventional binary classification trees such as CART either split the data using axis-aligned hyperplanes or they perform a computationally expensive search in the continuous space of hyperplanes with unrestricted orientations. We show that the limitations of the former can be overcome without resorting to the latter. For every pair of training data-points, there is one hyperplane that is orthogonal to the line joining the data-points and bisects this line. Such hyperplanes are plausible candidates for splits. In a comparison on a suite of 12 datasets we found that this method of generating candidate splits outperformed the standard methods, particularly when the training sets were small.

Dynamics of On-Line Gradient Descent Learning for Multilayer Neural Networks

David Saad, Sara Solla

We consider the problem of on-line gradient descent learning for general two-layer neural networks. An analytic solution is presented and used to investigate the role of the learning rate in controlling the evolution and convergence of the learning process.

Geometry of Early Stopping in Linear Networks

Robert Dodier

A theory of early stopping as applied to linear models is presented. The backpropagation learning algorithm is modeled as gradient descent in continuous time. Given a training set and a validation set, all weight vectors found by early stopping must lie on a certain quadric surface, usually an ellipsoid. Given a training set and a candidate early stopping weight vector, all validation sets have least-squares weights lying on a certain plane. This latter fact can be exploited to estimate the probability of stopping at any given point along the trajectory from the initial weight vector to the least-squares weights derived from the training set, and to estimate the probability that training goes on indefinitely. The prospects for extending this theory to nonlinear models are discussed.

A Unified Learning Scheme: Bayesian-Kullback Ying-Yang Machine

Lei Xu

A Bayesian-Kullback learning scheme, called Ying-Yang Machine, is proposed based on the two complementary but equivalent Bayesian representations for joint density and their Kullback divergence. Not only the scheme unifies existing major supervised and unsupervised learnings, including the classical maximum likelihood or least square learning, the maximum information preservation, the EM algorithm and information geometry, the recent popular Helmholtz machine, as well as other learning methods with new variants and new results; but also the scheme provides a number of new learning models.

Using Feedforward Neural Networks to Monitor Alertness from Changes in EEG Correlation and Coherence

Scott Makeig, Tzyy-Ping Jung, Terrence J. Sejnowski

We report here that changes in the normalized electroencephalographic (EEG) cross-spectrum can be used in conjunction with feedforward neural networks to monitor changes in alertness of operators continuously and in near-real time. Previously, we have shown that EEG spectral amplitudes covary with changes in alertness as indexed by changes in behavioral error rate on an auditory detection task [6,4]. Here, we report for the first time that increases in the frequency of detection errors in this task are also accompanied by patterns of increased and decreased spectral coherence in several frequency bands and EEG channel pairs. Relationships between EEG coherence and performance vary between subjects, but within subjects, their topographic and spectral profiles appear stable from session to session. Changes in alertness also covary with changes in correlations among EEG waveforms recorded at different scalp sites, and neural networks can also estimate alertness from correlation changes in spontaneous and unobtrusively recorded EEG signals.

A Dynamical Model of Context Dependencies for the Vestibulo-Ocular Reflex

Olivier Coenen, Terrence J. Sejnowski

The vestibulo-ocular reflex (VOR) stabilizes images on the retina during rapid head motions. The gain of the VOR (the ratio of eye to head rotation velocity) is typically around -1 when the eyes are focused on a distant target. However, to stabilize images accurately, the VOR gain must vary with context

(eye position, eye vergence and head translation). We first describe a kinematic model of the VOR which relies solely on sensory information available from the semicircular canals (head rotation), the otoliths (head translation), and neural correlates of eye position and vergence angle. We then propose a dynamical model and compare it to the eye velocity responses measured in monkeys. The dynamical model reproduces the observed amplitude and time course of the modulation of the VOR and suggests one way to combine the required neural signals within the cerebellum and the brain stem. It also makes predictions for the responses of neurons to multiple inputs (head rotation and translation, eye position, etc.) in the oculomotor system.

Constructive Algorithms for Hierarchical Mixtures of Experts

Steve Waterhouse, Anthony Robinson

We present two additions to the hierarchical mixture of experts (HME) architecture. By applying a likelihood splitting criteria to each expert in the HME we "grow" the tree adaptively during training. Secondly, by considering only the most probable path through the tree we may "prune" branches away, either temporarily, or permanently if they become redundant. We demonstrate results for the growing and path pruning algorithms which show significant speed ups and more efficient use of parameters over the standard fixed structure in discriminating between two interlocking spirals and classifying 8-bit parity patterns.

Discovering Structure in Continuous Variables Using Bayesian Networks

Reimar Hofmann, Volker Tresp

We study Bayesian networks for continuous variables using linear conditional density estimators. We demonstrate that full structures can be extracted from a data set in a self-organized way and we present sampling techniques for belief update based on Markov blanket conditional density models.

Temporal coding in the sub-millisecond range: Model of barn owl auditory pathway

Richard Kempster, Wulfram Gerstner, J. van Hemmen, Hermann Wagner

Binaural coincidence detection is essential for the localization of external sounds and requires auditory signal processing with high temporal precision. We present an integrate-and-fire model of spike processing in the auditory pathway of the barn owl. It is shown that a temporal precision in the microsecond range can be achieved with neuronal time constants which are at least one magnitude longer. An important feature of our model is an unsupervised Hebbian learning rule which leads to a temporal fine tuning of the neuronal connections.

Stable Dynamic Parameter Adaption

Stefan Rüger

A stability criterion for dynamic parameter adaptation is given. In the case of the learning rate of backpropagation, a class of stable algorithms is presented and studied, including a convergence proof.

Optimizing Cortical Mappings

Geoffrey Goodhill, Steven Finch, Terrence J. Sejnowski

"Topographic" mappings occur frequently in the brain. A popular approach to understanding the structure of such mappings is to map points representing input features in a space of a few dimensions to points in a 2 dimensional space using some self-organizing algorithm. We argue that a more general approach may be useful, where similarities between features are not constrained to be geometric distances, and the objective function for topographic matching is chosen explicitly rather than being specified implicitly by the self-organizing algorithm. We investigate analytically an example of this more general approach applied to the structure of interdigitated mappings, such as the pattern of ocular dominance columns in primary

y visual cortex.

Experiments with Neural Networks for Real Time Implementation of Control

Peter Campbell, Michael Dale, Herman Ferrá, Adam Kowalczyk

This paper describes a neural network based controller for allocating capacity in a telecommunications network. This system was proposed in order to overcome a "real time" response constraint. Two basic architectures are evaluated: 1) a feedforward network-heuristic and; 2) a feedforward network-recurrent network. These architectures are compared against a linear programming (LP) optimiser as a benchmark. This LP optimiser was also used as a teacher to label the data samples for the feedforward neural network training algorithm. It is found that the systems are able to provide a traffic throughput of 99% and 95%, respectively, of the throughput obtained by the linear programming solution. Once trained, the neural network based solutions are found in a fraction of the time required by the LP optimiser.

Cholinergic suppression of transmission may allow combined associative memory function and self-organization in the neocortex

Michael Hasselmo, Milos Cekic

Selective suppression of transmission at feedback synapses during learning is proposed as a mechanism for combining associative feed(cid:173)back with self-organization of feed forward synapses. Experimental data demonstrates cholinergic suppression of synaptic transmission in layer I (feedback synapses), and a lack of suppression in layer IV (feed(cid:173)forward synapses). A network with this feature uses local rules to learn mappings which are not linearly separable. During learning, sensory stimuli and desired response are simultaneously presented as input. Feedforward connections form self-organized representations of input, while suppressed feedback connections learn the transpose of feedfor(cid:173)ward connectivity. During recall, suppression is removed, sensory input activates the self-organized representation, and activity generates the learned response.

Dynamics of Attention as Near Saddle-Node Bifurcation Behavior

Hiroyuki Nakahara, Kenji Doya

In consideration of attention as a means for goal-directed behav(cid:173)ior in non-stationary environments, we argue that the dynamics of attention should satisfy two opposing demands: long-term main(cid:173)tenance and quick transition. These two characteristics are con(cid:173)tradictory within the linear domain. We propose the near saddle(cid:173)node bifurcation behavior of a sigmoidal unit with self-connection as a candidate of dynamical mechanism that satisfies both of these demands. We further show in simulations of the 'bug-eat-food' tasks that the near saddle-node bifurcation behavior of recurrent networks can emerge as a functional property for survival in non(cid:173)stationary environments.

Softassign versus Softmax: Benchmarks in Combinatorial Optimization

Steven Gold, Anand Rangarajan

A new technique, termed soft assign, is applied for the first time to two classic combinatorial optimization problems, the travel(cid:173)ing salesman problem and graph partitioning. Soft assign, which has emerged from the recurrent neural network/statistical physics framework, enforces two-way (assignment) constraints without the use of penalty terms in the energy functions. The soft assign can also be generalized from two-way winner-take-all constraints to multiple membership constraints which are required for graph par(cid:173)titioning. The soft assign technique is compared to the softmax (Potts glass). Within the statistical physics framework, softmax and a penalty term has been a widely used method for enforcing the two-way constraints common within many combinatorial optimiza(cid:173)tion problems. The benchmarks present evidence that soft assign has clear advantages in accuracy, speed, parallelizability and algo(cid:173)arithmic simplicity

ty over softmax and a penalty term in optimization problems with two-way constraints.

From Isolation to Cooperation: An Alternative View of a System of Experts

Stefan Schaal, Christopher Atkeson

We introduce a constructive, incremental learning system for regression problems that models data by means of locally linear experts. In contrast to other approaches, the experts are trained independently and do not compete for data during learning. Only when a prediction for a query is required do the experts cooperate by blending their individual predictions.

Each expert is trained by minimizing a penalized local cross validation error using second order methods. In this way, an expert is able to find a local distance metric by adjusting the size and shape of the receptive field in which its predictions are valid, and also to detect relevant input features by adjusting its bias on the importance of individual input dimensions. We derive asymptotic results for our method. In a variety of simulations the properties of the algorithm are demonstrated with respect to interference, learning speed, prediction accuracy, feature detection, and task oriented incremental learning.

Fast Learning by Bounding Likelihoods in Sigmoid Type Belief Networks

Tommi Jaakkola, Lawrence Saul, Michael Jordan

Sigmoid type belief networks, a class of probabilistic neural networks, provide a natural framework for compactly representing probabilistic information in a variety of unsupervised and supervised learning problems. Often the parameters used in these networks need to be learned from examples. Unfortunately, estimating the parameters via exact probabilistic calculations (i.e., the EM-algorithm) is intractable even for networks with fairly small numbers of hidden units. We propose to avoid the infeasibility of the E step by bounding likelihoods instead of computing them exactly. We introduce extended and complementary representations for these networks and show that the estimation of the network parameters can be made fast (reduced to quadratic optimization) by performing the estimation in either of the alternative domains. The complementary networks can be used for continuous density estimation as well.

Neural Control for Nonlinear Dynamic Systems

Ssu-Hsin Yu, Anuradha Annaswamy

A neural network based approach is presented for controlling two distinct types of nonlinear systems. The first corresponds to nonlinear systems with parametric uncertainties where the parameters occur nonlinearly. The second corresponds to systems for which stabilizing control structures can not be determined. The proposed neural controllers are shown to result in closed-loop system stability under certain conditions.

Sample Complexity for Learning Recurrent Perceptron Mappings

Bhaskar DasGupta, Eduardo Sontag

Recurrent perceptron classifiers generalize the classical perceptron model. They take into account those correlations and dependences among input coordinates which arise from linear digital filtering. This paper provides tight bounds on sample complexity associated to the fitting of such models to experimental data.

Is Learning The n-th Thing Any Easier Than Learning The First?

Sebastian Thrun

This paper investigates learning in a lifelong context. Lifelong learning addresses situations in which a learner faces a whole stream of learning tasks. Such scenarios provide the opportunity to transfer knowledge across multiple learning tasks, in order to generalize more accurately from less

s training data. In this paper, several different approaches to lifelong learning are described, and applied in an object recognition domain. It is shown that across the board, lifelong learning approaches generalize consistently more accurately from less training data, by their ability to transfer knowledge across learning tasks.

Parallel analog VLSI architectures for computation of heading direction and time-to-contact

Giacomo Indiveri, Jörg Kramer, Christof Koch

We describe two parallel analog VLSI architectures that integrate optical flow data obtained from arrays of elementary velocity sensors to estimate heading direction and time-to-contact. For heading direction computation, we performed simulations to evaluate the most important qualitative properties of the optical flow field and determine the best functional operators for the implementation of the architecture. For time-to-contact we exploited the divergence theorem to integrate data from all velocity sensors present in the architecture and average out possible errors.

A Dynamical Systems Approach for a Learnable Autonomous Robot

Jun Tani, Naohiro Fukumura

This paper discusses how a robot can learn goal-directed navigation tasks using local sensory inputs. The emphasis is that such learning tasks could be formulated as an embedding problem of dynamical systems: desired trajectories in a task space should be embedded into an adequate sensory-based internal state space so that a unique mapping from the internal state space to the motor command could be established. The paper shows that a recurrent neural network suffices in self-organizing such an adequate internal state space from the temporal sensory input. In our experiments, using a real robot with a laser range sensor, the robot navigated robustly by achieving dynamical coherence with the environment. It was also shown that such coherence becomes structurally stable as the global attractor is self-organized in the coupling of the internal and the environmental dynamics.

Optimization Principles for the Neural Code

Michael DeWeese

Recent experiments show that the neural codes at work in a wide range of creatures share some common features. At first sight, these observations seem unrelated. However, we show that these features arise naturally in a linear filtered threshold crossing (LFTC) model when we set the threshold to maximize the transmitted information. This maximization process requires neural adaptation to not only the DC signal level, as in conventional light and dark adaptation, but also to the statistical structure of the signal and noise distributions. We also present a new approach for calculating the mutual information between a neuron's output spike train and any aspect of its input signal which does not require reconstruction of the input signal. This formulation is valid provided the correlations in the spike train are small, and we provide a procedure for checking this assumption. This paper is based on joint work (DeWeese [1], 1995). Preliminary results from the LFTC model appeared in a previous proceedings (DeWeese [2], 1995), and the conclusions we reached at that time have been reaffirmed by further analysis of the model.

A Framework for Non-rigid Matching and Correspondence

Suguna Pappu, Steven Gold, Anand Rangarajan

Matching feature point sets lies at the core of many approaches to object recognition. We present a framework for non-rigid matching that begins with a skeleton module, affine point matching, and then integrates multiple features to improve correspondence and develops an object representation based on spatial regions to model local transformations. The algorithm for feature matching

iteratively updates the transformation parameters and the correspondence solution, each in turn. The affine mapping is solved in closed form, which permits its use for data of any dimension. The correspondence is set via a method for two-way constraint satisfaction, called softassign, which has recently emerged from the neural network/statistical physics realm. The complexity of the non-rigid matching algorithm with multiple features is the same as that of the affine point matching algorithm. Results for synthetic and real world data are provided for point sets in 2D and 3D, and for 2D data with multiple types of features and parts.

Hierarchical Recurrent Neural Networks for Long-Term Dependencies

Salah Hihi, Yoshua Bengio

We have already shown that extracting long-term dependencies from sequential data is difficult, both for deterministic dynamical systems such as recurrent networks, and probabilistic models such as hidden Markov models (HMMs) or input/output hidden Markov models (IOHMMs). In practice, to avoid this problem, researchers have used domain specific a-priori knowledge to give meaning to the hidden or state variables representing past context. In this paper, we propose to use a more general type of a-priori knowledge, namely that the temporal dependencies are structured hierarchically. This implies that long-term dependencies are represented by variables with a long time scale. This principle is applied to a recurrent network which includes delays and multiple time scales. Experiments confirm the advantages of such structures. A similar approach is proposed for HMMs and IOHMMs.

Finite State Automata that Recurrent Cascade-Correlation Cannot Represent

Stefan Kremer

This paper relates the computational power of Fahlman's Recurrent Cascade Correlation (RCC) architecture to that of finite state automata (FSA). While some recurrent networks are FSA equivalent, RCC is not. The paper presents a theoretical analysis of the RCC architecture in the form of a proof describing a large class of FSA which cannot be realized by RCC.

Memory-based Stochastic Optimization

Andrew Moore, Jeff Schneider

In this paper we introduce new algorithms for optimizing noisy plants in which each experiment is very expensive. The algorithms build a global non-linear model of the expected output at the same time as using Bayesian linear regression analysis of locally weighted polynomial models. The local models answer queries about confidence, noise, gradient and Hessians, and use them to make automated decisions similar to those made by a practitioner of Response Surface Methodology. The global and local models are combined naturally as a locally weighted regression. We examine the question of whether the global model can really help optimization, and we extend it to the case of time-varying functions. We compare the new algorithms with a highly tuned higher-order stochastic optimization algorithm on randomly-generated functions and a simulated manufacturing task. We note significant improvements in total regret, time to converge, and final solution quality.

Handwritten Word Recognition using Contextual Hybrid Radial Basis Function Network/Hidden Markov Models

Bernard Lemarié, Michel Gilloux, Manuel Leroux

A hybrid and contextual radial basis function network/hidden Markov model off-line handwritten word recognition system is presented. The task assigned to the radial basis function networks is the estimation of emission probabilities associated to Markov states. The model is contextual because the estimation of emission probabilities takes into account the left context of the current image segment as represented by its predecessor in the

he sequence. The new system does not outperform the previous system without context but acts differently.

Universal Approximation and Learning of Trajectories Using Oscillators

Pierre Baldi, Kurt Hornik

Natural and artificial neural circuits must be capable of traversing specific state space trajectories. A natural approach to this problem is to learn the relevant trajectories from examples. Unfortunately, gradient descent learning of complex trajectories in amorphous networks is unsuccessful. We suggest a possible approach where trajectories are realized by combining simple oscillators, in various modular ways. We contrast two regimes of fast and slow oscillations. In all cases, we show that banks of oscillators with bounded frequencies have universal approximation properties. Open questions are also discussed briefly.

Correlated Neuronal Response: Time Scales and Mechanisms

Wyeth Bair, Ehud Zohary, Christof Koch

We have analyzed the relationship between correlated spike count and the peak in the cross-correlation of spike trains for pairs of simultaneously recorded neurons from a previous study of area MT in the macaque monkey (Zohary et al., 1994). We conclude that common input, responsible for creating peaks on the order of ten milliseconds wide in the spike train cross-correlograms (CCGs), is also responsible for creating the correlation in spike count observed at the two second time scale of the trial. We argue that both common excitation and inhibition may play significant roles in establishing this correlation.

Worst-case Loss Bounds for Single Neurons

David Helmbold, Jyrki Kivinen, Manfred K. K. Warmuth

We analyze and compare the well-known Gradient Descent algorithm and a new algorithm, called the Exponentiated Gradient algorithm, for training a single neuron with an arbitrary transfer function. Both algorithms are easily generalized to larger neural networks, and the generalization of Gradient Descent is the standard back-propagation algorithm. In this paper we prove worst-case loss bounds for both algorithms in the single neuron case. Since local minima make it difficult to prove worst-case bounds for gradient-based algorithms, we must use a loss function that prevents the formation of spurious local minima. We define such a matching loss function for any strictly increasing differentiable transfer function and prove worst-case loss bound for any such transfer function and its corresponding matching loss.

For example, the matching loss for the identity function is the square loss and the matching loss for the logistic sigmoid is the entropic loss. The different structure of the bounds for the two algorithms indicates that the new algorithm out-performs Gradient Descent when the inputs contain a large number of irrelevant components.

Stock Selection via Nonlinear Multi-Factor Models

Asriel Levin

This paper discusses the use of multilayer feed forward neural networks for predicting a stock's excess return based on its exposure to various technical and fundamental factors. To demonstrate the effectiveness of the approach a hedged portfolio which consists of equally capitalized long and short positions is constructed and its historical returns are benchmarked against T-bill returns and the S&P500 index.

The Geometry of Eye Rotations and Listing's Law

Amir Handzel, Tamar Flash

We analyse the geometry of eye rotations, and in particular saccades,

using basic Lie group theory and differential geometry. Various parameterizations of rotations are related through a unifying mathematical treatment, and transformations between co-ordinate systems are computed using the Campbell-Baker-Hausdorff formula. Next, we describe Listing's law by means of the Lie algebra $so(3)$. This enables us to demonstrate a direct connection to Donders' law, by showing that eye orientations are restricted to the quotient space $SO(3)/SO(2)$. The latter is equivalent to the sphere S^2 , which is exactly the space of gaze directions. Our analysis provides a mathematical framework for studying the oculomotor system and could also be extended to investigate the geometry of multi-joint arm movements.

When is an Integrate-and-fire Neuron like a Poisson Neuron?

Charles Stevens, Anthony Zador

In the Poisson neuron model, the output is a rate-modulated Poisson process (Snyder and Miller, 1991); the time varying rate parameter $\lambda(t)$ is an instantaneous function $G[\cdot]$ of the stimulus, $s(t) = G[s(t)]$. In a Poisson neuron, then, $\lambda(t)$ gives the instantaneous firing rate—the instantaneous probability of firing at any instant t —and the output is a stochastic function of the input. In part because of its great simplicity, this model is widely used (usually with the addition of a refractory period), especially in *in vivo* single unit electrophysiological studies, where $s(t)$ is usually taken to be the value of some sensory stimulus. In the integrate-and-fire neuron model, by contrast, the output is a filtered and thresholded function of the input: the input is passed through a low-pass filter (determined by the membrane time constant τ) and integrated until the membrane potential $v(t)$ reaches threshold θ , at which point $v(t)$ is reset to its initial value. By contrast with the Poisson model, in the integrate-and-fire model the output is a deterministic function of the input. Although the integrate-and-fire model is a caricature of real neural dynamics, it captures many of the qualitative features, and is often used as a starting point for conceptualizing the biophysical behavior of single neurons. Here we show how a slightly modified Poisson model can be derived from the integrate-and-fire model with noisy inputs ($\lambda(t) = s(t) + \eta(t)$). In the modified model, the transfer function $G[\cdot]$ is a sigmoid (erf) whose shape is determined by the noise variance σ^2/τ . Understanding the equivalence between the dominant *in vivo* and *in vitro* simple neuron models may help forge links between the two levels.

Examples of learning curves from a modified VC-formalism

Adam Kowalczyk, Jacek Szymanski, Peter Bartlett, Robert C. Williamson

We examine the issue of evaluation of model specific parameters in a modified VC-formalism. Two examples are analyzed: the 2-dimensional homogeneous perceptron and the 1-dimensional higher order neuron. Both models are solved theoretically, and their learning curves are compared against true learning curves. It is shown that the formalism has the potential to generate a variety of learning curves, including ones displaying 'phase transitions.'

Using Unlabeled Data for Supervised Learning

Geoffrey Towell

Many classification problems have the property that the only costly part of obtaining examples is the class label. This paper suggests a simple method for using distribution information contained in unlabeled examples to augment labeled examples in a supervised training framework. Empirical tests show that the technique described in this paper can significantly improve the accuracy of a supervised learner when the learner is well below its asymptotic accuracy level.

Implementation Issues in the Fourier Transform Algorithm

Yishay Mansour, Sigal Sahar

The Fourier transform of boolean functions has come to play an important role in proving many important learnability results. We aim to demonstrate that the Fourier transform techniques are also a useful and practical algorithm in addition to being a powerful theoretical tool. We describe the more prominent changes we have introduced to the algorithm, ones that were crucial and without which the performance of the algorithm would severely deteriorate. One of the benefits we present is the confidence level for each prediction which measures the likelihood the prediction is correct.

A Bound on the Error of Cross Validation Using the Approximation and Estimation Rates, with Consequences for the Training-Test Split

Michael Kearns

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Classifying Facial Action

Marian Bartlett, Paul Viola, Terrence J. Sejnowski, Beatrice Golomb, Jan Larsen, Joseph Hager, Paul Ekman

The Facial Action Coding System, (FACS), devised by Ekman and Friesen (1978), provides an objective means for measuring the facial muscle contractions involved in a facial expression. In this paper, we approach automated facial expression analysis by detecting and classifying facial actions. We generated a database of over 1100 image sequences of 24 subjects performing over 150 distinct facial actions or action combinations. We compare three different approaches to classifying the facial actions in these images: Holistic spatial analysis based on principal components of graylevel images; explicit measurement of local image features such as wrinkles; and template matching with motion flow fields. On a dataset containing six individual actions and 20 subjects, these methods had 89%, 57%, and 85% performances respectively for generalization to novel subjects. When combined, performance improved to 92%.

Strong Unimodality and Exact Learning of Constant Depth μ -Perceptron Networks

Mario Marchand, Saeed Hadjifaradji

Requests for name changes in the electronic proceedings will be accepted with no questions asked. However name changes may cause bibliographic tracking issues.

Authors are asked to consider this carefully and discuss it with their co-authors prior to requesting a name change in the electronic proceedings.

Gradient and Hamiltonian Dynamics Applied to Learning in Neural Networks

James Howse, Chaouki Abdallah, Gregory Heileman

The process of machine learning can be considered in two stages: model selection and parameter estimation. In this paper a technique is presented for constructing dynamical systems with desired qualitative properties. The approach is based on the fact that an n -dimensional nonlinear dynamical system can be decomposed into one gradient and $(n - 1)$ Hamiltonian systems. Thus, the model selection stage consists of choosing the gradient and Hamiltonian portions appropriately so that a certain behavior is obtainable.

To estimate the parameters, a stably convergent learning rule is presented. This algorithm has been proven to converge to the desired system trajectory for all initial conditions and system inputs. This technique can be used to design neural network models which are guaranteed to solve the trajectory learning problem.

A New Learning Algorithm for Blind Signal Separation

Shun-ichi Amari, Andrzej Cichocki, Howard Yang

A new on-line learning algorithm which minimizes a statistical dependency among outputs is derived for blind separation of mixed signals. The dependency is measured by the average mutual information (MI) of the outputs. The source signals and the mixing matrix are unknown except for the number of the sources. The Gram-Charlier expansion instead of the Edgeworth expansion is used in evaluating the MI. The natural gradient approach is used to minimize the MI. A novel activation function is proposed for the on-line learning algorithm which has an equivariant property and is easily implemented on a neural network like model. The validity of the new learning algorithm are verified by computer simulations.

Adaptive Back-Propagation in On-Line Learning of Multilayer Networks

Ansgar West, David Saad

An adaptive back-propagation algorithm is studied and compared with gradient descent (standard back-propagation) for on-line learning in two-layer neural networks with an arbitrary number of hidden units. Within a statistical mechanics framework, both numerical studies and a rigorous analysis show that the adaptive back-propagation method results in faster training by breaking the symmetry between hidden units more efficiently and by providing faster convergence to optimal generalization than gradient descent.

Neuron-MOS Temporal Winner Search Hardware for Fully-Parallel Data Processing

Tadashi Shibata, Tsutomu Nakai, Tatsuo Morimoto, Ryu Kaihara, Takeo Yamashita, Tadahiro Ohmi

A unique architecture of winner search hardware has been developed using a novel neuron-like high functionality device called Neuron MOS transistor (or vMOS in short) [1,2] as a key circuit element. The circuits developed in this work can find the location of the maximum (or minimum) signal among a number of input data on the continuous-time basis, thus enabling real-time winner tracking as well as fully-parallel sorting of multiple input data. We have developed two circuit schemes. One is an ensemble of self-loop-selecting vMOS ring oscillators finding the winner as an oscillating node. The other is an ensemble of vMOS variable threshold inverters receiving a common ramp-voltage for competitive excitation where data sorting is conducted through consecutive winner search actions. Test circuits were fabricated by a double-polysilicon CMOS process and their operation has been experimentally verified.

Human Reading and the Curse of Dimensionality

Gale Martin

Whereas optical character recognition (OCR) systems learn to classify single characters; people learn to classify long character strings in parallel, within a single fixation. This difference is surprising because high dimensionality is associated with poor classification learning. This paper suggests that the human reading system avoids these problems because the number of to-be-classified images is reduced by consistent and optimal eye fixation positions, and by character sequence regularities.

Unsupervised Pixel-prediction

William Softky

When a sensory system constructs a model of the environment from its input, it might need to verify the model's accuracy. One method of verification is multivariate time-series prediction: a good model could predict the near-future activity of its inputs, much as a good scientific theory predicts future data. Such a predicting model would require copious top-down connections to compare the predictions with the input. That feedback could improve the model's per

formance in two ways: by biasing internal activity to ward expected patterns, and by generating specific error signals if the predictions fail. A proof-of-concept model-an event-driven, computationally efficient layered network, incorporating "cortical" features like all-excitatory synapses and local inhibition- was constructed to make near-future predictions of a simple, moving stimulus. After unsupervised learning, the network contained units not only tuned to obvious features of the stimulus like contour orientation and motion, but also to contour discontinuity ("end-stopping") and illusory contours.

Control of Selective Visual Attention: Modeling the "Where" Pathway

Ernst Niebur, Christof Koch

Intermediate and higher vision processes require selection of a subset of the available sensory information before further processing. Usually, this selection is implemented in the form of a spatially circumscribed region of the visual field, the so-called "focus of attention" which scans the visual scene dependent on the input and on the attentional state of the subject. We here present a model for the control of the focus of attention in primates, based on a saliency map. This mechanism is not only expected to model the functionality of biological vision but also to be essential for the understanding of complex scenes in machine vision.

Quadratic-Type Lyapunov Functions for Competitive Neural Networks with Different Time-Scales

Anke Meyer-Bäse

The dynamics of complex neural networks modelling the self-organization process in cortical maps must include the aspects of long and short-term memory. The behaviour of the network is such characterized by an equation of neural activity as a fast phenomenon and an equation of synaptic modification as a slow part of the neural system. We present a quadratic-type Lyapunov function for the flow of a competitive neural system with fast and slow dynamic variables. We also show the consequences of the stability analysis on the neural net parameters.

Learning long-term dependencies is not as difficult with NARX networks

Tsungnan Lin, Bill Horne, Peter Tiffo, C. Giles

It has recently been shown that gradient descent learning algorithms for recurrent neural networks can perform poorly on tasks that involve long-term dependencies. In this paper we explore this problem for a class of architectures called NARX networks, which have powerful representational capabilities. Previous work reported that gradient descent learning is more effective in NARX networks than in recurrent networks with "hidden states". We show that although NARX networks do not circumvent the problem of long-term dependencies, they can greatly improve performance on such problems. We present some experimental results that show that NARX networks can often retain information for two to three times as long as conventional recurrent networks.

Learning the Structure of Similarity

Joshua Tenenbaum

The additive clustering (ADCLUS) model (Shepard & Arabie, 1979) treats the similarity of two stimuli as a weighted additive measure of their common features. Inspired by recent work in unsupervised learning with multiple cause models, we propose anew, statistically well-motivated algorithm for discovering the structure of natural stimulus classes using the ADCLUS model, which promises substantial gains in conceptual simplicity, practical efficiency, and solution quality over earlier efforts. We also present preliminary results with artificial data and two classic similarity data sets.

Investment Learning with Hierarchical PSOMs

Jörg Walter, Helge Ritter

We propose a hierarchical scheme for rapid learning of context dependent "skills" that is based on the recently introduced "Parameterized Self-Organizing Map" ("PSOM"). The underlying idea is to first invest some learning effort to specialize the system into a rapid learner for a more restricted range of contexts.

How Perception Guides Production in Birdsong Learning

Christopher Fry

A computational model of song learning in the song sparrow (*Melospiza melodia*) learns to categorize the different syllables of a song sparrow song and uses this categorization to train itself to reproduce song. The model fills a crucial gap in the computational explanation of birdsong learning by exploring the organization of perception in songbirds. It shows how competitive learning may lead to the organization of a specific nucleus in the bird brain, replicates the song production results of a previous model (Doya and Sejnowski, 1995), and demonstrates how perceptual learning can guide production through reinforcement learning.

Stable Fitted Reinforcement Learning

Geoffrey J. Gordon

We describe the reinforcement learning problem, motivate algorithms which seek an approximation to the Q function, and present new convergence results for two such algorithms.

Information through a Spiking Neuron

Charles Stevens, Anthony Zador

While it is generally agreed that neurons transmit information about their synaptic inputs through spike trains, the code by which this information is transmitted is not well understood. An upper bound on the information encoded is obtained by hypothesizing that the precise timing of each spike conveys information. Here we develop a general approach to quantifying the information carried by spike trains under this hypothesis, and apply it to the leaky integrate-and-fire (IF) model of neuronal dynamics. We formulate the problem in terms of the probability distribution $p(T)$ of interspike intervals (ISIs), assuming that spikes are detected with arbitrary but finite temporal resolution. In the absence of added noise, all the variability in the ISIs could encode information, and the information rate is simply the entropy of the ISI distribution, $H(T) = (-p(T) \log_2 p(T))$, times the spike rate. $H(T)$ thus provides an exact expression for the information rate. The methods developed here can be used to determine experimentally the information carried by spike trains, even when the lower bound of the information rate provided by the stimulus reconstruction method is not tight. In a preliminary series of experiments, we have used these methods to estimate information rates of hippocampal neurons in slice in response to somatic current injection. These pilot experiments suggest information rates as high as 6.3 bits/spike.

Discriminant Adaptive Nearest Neighbor Classification and Regression

Trevor Hastie, Robert Tibshirani

Nearest neighbor classification expects the class conditional probabilities to be locally constant, and suffers from bias in high dimensions. We propose a locally adaptive form of nearest neighbor classification to try to finesse this curse of dimensionality. We use a local linear discriminant analysis to estimate an effective metric for computing neighborhoods. We determine the local decision boundaries from centroid information, and then shrink neighborhoods in directions orthogonal to these local decision boundaries, and elongate

them parallel to the boundaries. Thereafter, any neighborhood-based classifier can be employed, using the modified neighborhoods. We also propose a method for global dimension reduction, that combines local dimension information. We indicate how these techniques can be extended to the regression problem.

A Predictive Switching Model of Cerebellar Movement Control

Andrew Barto, James Houk

We present a hypothesis about how the cerebellum could participate in regulating movement in the presence of significant feedback delays without resorting to a forward model of the motor plant. We show how a simplified cerebellar model can learn to control end point positioning of a nonlinear spring-mass system with realistic delays in both afferent and efferent pathways. The model's operation involves prediction, but instead of predicting sensory input, it directly regulates movement by reacting in an anticipatory fashion to input patterns that include delayed sensory feedback.

Recurrent Neural Networks for Missing or Asynchronous Data

Yoshua Bengio, Francois Gingras

In this paper we propose recurrent neural networks with feedback into the input units for handling two types of data analysis problems. On the one hand, this scheme can be used for static data when some of the input variables are missing. On the other hand, it can also be used for sequential data, when some of the input variables are missing or are available at different frequencies. Unlike in the case of probabilistic models (e.g. Gaussian) of the missing variables, the network does not attempt to model the distribution of the missing variables given the observed variables. Instead it is a more "discriminant" approach that fills in the missing variables for the sole purpose of minimizing a learning criterion (e.g., to minimize an output error).
