

## SIMPLIFYING NEURAL NETS BY DISCOVERING FLAT MINIMA

Sepp Hochreiter, Jürgen Schmidhuber

We present a new algorithm for finding low complexity networks with high generalization capability. The algorithm searches for large connected regions of so-called 'flat' minima of the error function. In the weight-space environment of a "flat" minimum, the error remains approximately constant. Using an MDL-based argument, flat minima can be shown to correspond to low expected overfitting. Although our algorithm requires the computation of second order derivatives, it has backprop's order of complexity. Experiments with feedforward and recurrent nets are described. In an application to stock market prediction, the method outperforms conventional backprop, weight decay, and "optimal brain surgeon".

\*\*\*\*\*

## A Model of the Neural Basis of the Rat's Sense of Direction

William Skaggs, James Knierim, Hemant Kudrimoti, Bruce McNaughton

In the last decade the outlines of the neural structures subserving the sense of direction have begun to emerge. Several investigations have shed light on the effects of vestibular input and visual input on the head direction representation. In this paper, a model is formulated of the neural mechanisms underlying the head direction system. The model is built out of simple ingredients, depending on nothing more complicated than connectional specificity, attractor dynamics, Hebbian learning, and sigmoidal nonlinearities, but it behaves in a sophisticated way and is consistent with most of the observed properties of real head direction cells. In addition it makes a number of predictions that ought to be testable by reasonably straightforward experiments.

\*\*\*\*\*

## A Mixture Model System for Medical and Machine Diagnosis

Magnus Stensmo, Terrence J. Sejnowski

Diagnosis of human disease or machine fault is a missing data problem since many variables are initially unknown. Additional information needs to be obtained. The joint probability distribution of the data can be used to solve this problem. We model this with mixture models whose parameters are estimated by the EM algorithm. This gives the benefit that missing data in the database itself can also be handled correctly. The request for new information to refine the diagnosis is performed using the maximum utility principle. Since the system is based on learning it is domain independent and less labor intensive than expert systems or probabilistic networks. An example using a heart disease database is presented.

\*\*\*\*\*

## Learning with Preknowledge: Clustering with Point and Graph Matching Distance Measures

Steven Gold, Anand Rangarajan, Eric Mjolsness

Prior constraints are imposed upon a learning problem in the form of distance measures. Prototypical 2-D point sets and graphs are learned by clustering with point matching and graph matching distance measures. The point matching distance measure is approximately invariant under affine transformations - translation, rotation, scale and shear - and permutations. It operates between noisy images with missing and spurious points. The graph matching distance measure operates on weighted graphs and is invariant under permutations. Learning is formulated as an optimization problem. Large objectives so formulated (hundreds of million variables) are efficiently minimized using a combination of optimization techniques - algebraic transformations, iterative projective scaling, clocked objectives, and deterministic annealing.

\*\*\*\*\*

## Learning Local Error Bars for Nonlinear Regression

David Nix, Andreas Weigend

We present a new method for obtaining local error bars for nonlinear regression, i.e., estimates of the confidence in predicted values that depend on the input. We approach this problem by applying a maximum likelihood

likelihood framework to an assumed distribution of errors. We demonstrate our method first on computer-generated data with locally varying, normally distributed target noise. We then apply it to laser data from the Santa Fe Time Series Competition where the underlying system noise is known quantization error and the error bars give local estimates of model misspecification. In both cases, the method also provides a weighted regression effect that improves generalization performance.

\*\*\*\*\*

#### Reinforcement Learning Methods for Continuous-Time Markov Decision Problems

Steven Bradtke, Michael Duff

Semi-Markov Decision Problems are continuous time generalizations of discrete time Markov Decision Problems. A number of reinforcement learning algorithms have been developed recently for the solution of Markov Decision Problems, based on the ideas of asynchronous dynamic programming and stochastic approximation. Among these are TD( $x$ ), Q-Learning, and Real-time Dynamic Programming. After reviewing semi-Markov Decision Problems and Bellman's optimality equation in that context, we propose algorithms similar to those named above, adapted to the solution of semi-Markov Decision Problems. We demonstrate these algorithms by applying them to the problem of determining the optimal control for a simple queueing system. We conclude with a discussion of circumstances under which these algorithms may be usefully applied.

\*\*\*\*\*

#### Connectionist Speaker Normalization with Generalized Resource Allocating Networks

Cesare Furlanello, Diego Giuliani, Edmondo Trentin

The paper presents a rapid speaker-normalization technique based on neural network spectral mapping. The neural network is used as a front-end of a continuous speech recognition system (speaker dependent, HMM-based) to normalize the input acoustic data from a new speaker. The spectral difference between speakers can be reduced using a limited amount of new acoustic data (40 phonetically rich sentences). Recognition error of phone units from the acoustic-phonetic continuous speech corpus APASCI is decreased with an adaptability ratio of 25%. We used local basis networks of elliptical Gaussian kernels, with recursive allocation of units and on-line optimization of parameters (GRAN model). For this application, the model included a linear term. The results compare favorably with multivariate linear mapping based on constrained orthonormal transformations.

\*\*\*\*\*

#### A Novel Reinforcement Model of Birdsong Vocalization Learning

Kenji Doya, Terrence J. Sejnowski

Songbirds learn to imitate a tutor song through auditory and motor learning. We have developed a theoretical framework for song learning that accounts for response properties of neurons that have been observed in many of the nuclei that are involved in song learning. Specifically, we suggest that the anterior forebrain pathway, which is not needed for song production in the adult but is essential for song acquisition, provides synaptic perturbations and adaptive evaluations for syllable vocalization learning. A computer model based on reinforcement learning was constructed that could replicate a real zebra finch song with 90% accuracy based on a spectrographic measure. The second generation of the bird song model replicated the tutor song with 96% accuracy.

\*\*\*\*\*

#### Bias, Variance and the Combination of Least Squares Estimators

Ronny Meir

We consider the effect of combining several least squares estimators on the expected performance of a regression problem. Computing the exact bias and variance curves as a function of the sample size we are able to quantitatively compare the effect of the combination on the bias and variance separately, and

d thus on the expected error which is the sum of the two. Our exact calculations, demonstrate that the combination of estimators is particularly useful in the case where the data set is small and noisy and the function to be learned is unrealizable. For large data sets the single estimator produces superior results. Finally, we show that by splitting the data set into several independent parts and training each estimator on a different subset, the performance can in some cases be significantly improved.

\*\*\*\*\*

#### Hierarchical Mixtures of Experts Methodology Applied to Continuous Speech Recognition

Ying Zhao, Richard Schwartz, Jason Sroka, John Makhoul

In this paper, we incorporate the Hierarchical Mixtures of Experts (HME) method of probability estimation, developed by Jordan [1], into an HMM(cid:173) based continuous speech recognition system. The resulting system can be thought of as a continuous-density HMM system, but instead of using gaussian mixtures, the HME system employs a large set of hierarchically organized but relatively small neural networks to perform the probability density estimation. The hierarchical structure is reminiscent of a decision tree except for two important differences: each "expert" or neural net performs a "soft" decision rather than a hard decision, and, unlike ordinary decision trees, the parameters of all the neural nets in the HME are automatically trainable using the EM algorithm. We report results on the ARPA 5,000-word and 40,000-word Wall Street Journal corpus using HME models.

\*\*\*\*\*

#### A Comparison of Discrete-Time Operator Models for Nonlinear System Identification

Andrew Back, Ah Tsoi

We present a unifying view of discrete-time operator models used in the context of finite word length linear signal processing. Comparisons are made between the recently presented gamma operator model, and the delta and rho operator models for performing nonlinear system identification and prediction using neural networks. A new model based on an adaptive bilinear transformation which generalizes all of the above models is presented.

\*\*\*\*\*

#### Learning Many Related Tasks at the Same Time with Backpropagation

Rich Caruana

Hinton [6] proposed that generalization in artificial neural nets should improve if nets learn to represent the domain's underlying regularities. Abu-Mustafa's hints work [1] shows that the outputs of a backprop net can be used as inputs through which domain(cid:173) specific information can be given to the net. We extend these ideas by showing that a backprop net learning many related tasks at the same time can use these tasks as inductive bias for each other and thus learn better. We identify five mechanisms by which multitask backprop improves generalization and give empirical evidence that multi task backprop generalizes better in real domains.

\*\*\*\*\*

#### Multidimensional Scaling and Data Clustering

Thomas Hofmann, Joachim Buhmann

Visualizing and structuring pairwise dissimilarity data are difficult combinatorial op(cid:173) timization problems known as multidimensional scaling or pairwise data clustering. Algorithms for embedding dissimilarity data set in a Euclidean space, for clustering these data and for actively selecting data to support the clustering process are discussed in the maximum entropy framework. Active data selection provides a strategy to discover structure in a data set efficiently with partially unknown data.

\*\*\*\*\*

#### Predicting the Risk of Complications in Coronary Artery Bypass Operations using Neural Networks

Richard P. Lippmann, Linda Kukolich, David Shahian

Experiments demonstrated that sigmoid multilayer perceptron (MLP) networks provide slightly better risk prediction than conventional logistic regression when used to predict the risk of death, stroke, and renal failure on 1257 patients who underwent coronary artery bypass operations at the Lahey Clinic. MLP networks with no hidden layer and networks with one hidden layer were trained using stochastic gradient descent with early stopping. MLP networks and logistic regression used the same input features and were evaluated using bootstrap sampling with 50 replications. ROC areas for predicting mortality using preoperative input features were 70.5% for logistic regression and 76.0% for MLP networks. Regularization provided by early stopping was an important component of improved performance. A simplified approach to generating confidence intervals for MLP risk predictions using an auxiliary "confidence MLP" was developed. The confidence MLP is trained to reproduce confidence intervals that were generated during training using the outputs of 50 MLP networks trained with different bootstrap samples.

\*\*\*\*\*

#### Spatial Representations in the Parietal Cortex May Use Basis Functions

Alexandre Pouget, Terrence J. Sejnowski

The parietal cortex is thought to represent the egocentric positions of objects in particular coordinate systems. We propose an alternative approach to spatial perception of objects in the parietal cortex from the perspective of sensorimotor transformations. The responses of single parietal neurons can be modeled as a gaussian function of retinal position multiplied by a sigmoid function of eye position, which form a set of basis functions. We show here how these basis functions can be used to generate receptive fields in either retinotopic or head-centered coordinates by simple linear transformations. This raises the possibility that the parietal cortex does not attempt to compute the positions of objects in a particular frame of reference but instead computes a general purpose representation of the retinal location and eye position from which any transformation can be synthesized by direct projection. This representation predicts that hemineglect, a neurological syndrome produced by parietal lesions, should not be confined to egocentric coordinates, but should be observed in multiple frames of reference in single patients, a prediction supported by several experiments.

\*\*\*\*\*

#### Reinforcement Learning Algorithm for Partially Observable Markov Decision Problems

Tommi Jaakkola, Satinder Singh, Michael Jordan

Increasing attention has been paid to reinforcement learning algorithms in recent years, partly due to successes in the theoretical analysis of their behavior in Markov environments. If the Markov assumption is removed, however, neither generally the algorithms nor the analyses continue to be usable. We propose and analyze a new learning algorithm to solve a certain class of non-Markov decision problems. Our algorithm applies to problems in which the environment is Markov, but the learner has restricted access to state information. The algorithm involves a Monte-Carlo policy evaluation combined with a policy improvement method that is similar to that of Markov decision problems and is guaranteed to converge to a local maximum. The algorithm operates in the space of stochastic policies, a space which can yield a policy that performs considerably better than any deterministic policy. Although the space of stochastic policies is continuous-even for a discrete action space-our algorithm is computationally tractable.

\*\*\*\*\*

#### FINANCIAL APPLICATIONS OF LEARNING FROM HINTS

Yaser Abu-Mostafa

The basic paradigm for learning in neural networks is 'learning from examples' where a training set of input-output examples is used to teach the network the

target function. Learning from hints is a generalization of learning from examples where additional information about the target function can be incorporated in the same learning process. Such information can come from common sense rules or special expertise. In financial market applications where the training data is very noisy, the use of such hints can have a decisive advantage. We demonstrate the use of hints in foreign-exchange trading of the U.S. Dollar versus the British Pound, the German Mark, the Japanese Yen, and the Swiss Franc, over a period of 32 months. We explain the general method of learning from hints and how it can be applied to other markets. The learning model for this method is not restricted to neural networks.

\*\*\*\*\*

#### An Auditory Localization and Coordinate Transform Chip

Timothy Horiuchi

The localization and orientation to various novel or interesting events in the environment is a critical sensorimotor ability in all animals, predator or prey. In mammals, the superior colliculus (SC) plays a major role in this behavior, the deeper layers exhibiting topographically mapped responses to visual, auditory, and somatosensory stimuli. Sensory information arriving from different modalities should then be represented in the same coordinate frame. Auditory cues, in particular, are thought to be computed in head-based coordinates which must then be transformed to retinal coordinates. In this paper, an analog VLSI implementation for auditory localization in the azimuthal plane is described which extends the architecture proposed for the barn owl to a primate eye movement system where further transformation is required. This transformation is intended to model the projection in primates from auditory cortical areas to the deeper layers of the primate superior colliculus. This system is interfaced with an analog VLSI-based saccadic eye movement system also being constructed in our laboratory.

\*\*\*\*\*

#### Limits on Learning Machine Accuracy Imposed by Data Quality

Corinna Cortes, L. D. Jackel, Wan-Ping Chiang

Random errors and insufficiencies in databases limit the performance of any classifier trained from and applied to the database. In this paper we propose a method to estimate the limiting performance of classifiers imposed by the database. We demonstrate this technique on the task of predicting failure in telecommunication paths.

\*\*\*\*\*

#### Interference in Learning Internal Models of Inverse Dynamics in Humans

Reza Shadmehr, Tom Brashers-Krug, Ferdinando Mussa-Ivaldi

Experiments were performed to reveal some of the computational properties of the human motor memory system. We show that as humans practice reaching movements while interacting with a novel mechanical environment, they learn an internal model of the inverse dynamics of that environment. Subjects show recall of this model at testing sessions 24 hours after the initial practice. The representation of the internal model in memory is such that there is interference when there is an attempt to learn a new inverse dynamics map immediately after an anticorrelated mapping was learned. We suggest that this interference is an indication that the same computational elements used to encode the first inverse dynamics map are being used to learn the second mapping. We predict that this leads to a forgetting of the initially learned skill.

\*\*\*\*\*

#### Optimal Movement Primitives

Terence Sanger

The theory of Optimal Unsupervised Motor Learning shows how a network can discover a reduced-order controller for an unknown nonlinear system by representing only the most significant modes. Here, I extend the theory to apply to command sequences, so that the most significant components discovered by the network correspond to motion "primitives". Combina

tions of these primitives can be used to produce a wide variety of different movements. I demonstrate applications to human handwriting decomposition and synthesis, as well as to the analysis of electrophysiological experiments on movements resulting from stimulation of the frog spinal cord.

\*\*\*\*\*

Using a Saliency Map for Active Spatial Selective Attention: Implementation & Initial Results

Shumeet Baluja, Dean A. Pomerleau

In many vision based tasks, the ability to focus attention on the important portions of a scene is crucial for good performance on the tasks. In this paper we present a simple method of achieving spatial selective attention through the use of a saliency map. The saliency map indicates which regions of the input retina are important for performing the task. The saliency map is created through predictive auto-encoding. The performance of this method is demonstrated on two simple tasks which have multiple very strong distracting features in the input retina. Architectural extensions and application directions for this model are presented.

\*\*\*\*\*

Factorial Learning and the EM Algorithm

Zoubin Ghahramani

Many real world learning problems are best characterized by an interaction of multiple independent causes or factors. Discovering such causal structure from the data is the focus of this paper. Based on Zemel and Hinton's cooperative vector quantizer (CVQ) architecture, an unsupervised learning algorithm is derived from the Expectation-Maximization (EM) framework. Due to the combinatorial nature of the data generation process, the exact E-step is computationally intractable. Two alternative methods for computing the E-step are proposed: Gibbs sampling and mean-field approximation, and some promising empirical results are presented.

\*\*\*\*\*

Pairwise Neural Network Classifiers with Probabilistic Outputs

David Price, Stefan Knerr, Léon Personnaz, Gérard Dreyfus

Multi-class classification problems can be efficiently solved by partitioning the original problem into sub-problems involving only two classes: for each pair of classes, a (potentially small) neural network is trained using only the data of these two classes. We show how to combine the outputs of the two-class neural networks in order to obtain posterior probabilities for the class decisions. The resulting probabilistic pairwise classifier is part of a handwriting recognition system which is currently applied to check reading. We present results on real world data bases and show that, from a practical point of view, these results compare favorably to other neural network approaches.

\*\*\*\*\*

Real-Time Control of a Tokamak Plasma Using Neural Networks

Chris M. Bishop, Paul S. Haynes, Mike E U Smith, Tom N. Todd, David L. Trotman, Colin G. Windsor

This paper presents results from the first use of neural networks for the real-time feedback control of high temperature plasmas in a tokamak fusion experiment. The tokamak is currently the principal experimental device for research into the magnetic confinement approach to controlled fusion. In the tokamak, hydrogen plasmas, at temperatures of up to 100 Million K, are confined by strong magnetic fields.

Accurate control of the position and shape of the plasma boundary requires real-time feedback control of the magnetic field structure on a time-scale of a few tens of microseconds. Software simulations have demonstrated that a neural network approach can give significantly better performance than the linear technique currently used on most tokamak experiments. The practical application of the neural network approach requires high-speed hardware, for which a fully parallel implementation of the multil

ayer perceptron, using a hybrid of digital and analogue technology, has been developed.

\*\*\*\*\*

#### An Actor/Critic Algorithm that is Equivalent to Q-Learning

Robert Crites, Andrew Barto

We prove the convergence of an actor/critic algorithm that is equivalent to Q-learning by construction. Its equivalence is achieved by encoding Q-values within the policy and value function of the actor and critic. The resultant actor/critic algorithm is novel in two ways: it updates the critic only when the most probable action is executed from any given state, and it rewards the actor using criteria that depend on the relative probability of the action that was executed.

\*\*\*\*\*

#### Template-Based Algorithms for Connectionist Rule Extraction

Jay Alexander, Michael C. Mozer

Casting neural network weights in symbolic terms is crucial for interpreting and explaining the behavior of a network. Additionally, in some domains, a symbolic description may lead to more robust generalization. We present a principled approach to symbolic rule extraction based on the notion of weight templates, parameterized regions of weight space corresponding to specific symbolic expressions. With an appropriate choice of representation, we show how template parameters may be efficiently identified and instantiated to yield the optimal match to a unit's actual weights. Depending on the requirements of the application domain, our method can accommodate arbitrary disjunctions and conjunctions with  $O(k)$  complexity, simple n-of-m expressions with  $O(k!)$  complexity, or a more general class of recursive n-of-m expressions with  $O(k!)$  complexity, where  $k$  is the number of inputs to a unit. Our method of rule extraction offers several benefits over alternative approaches in the literature, and simulation results on a variety of problems demonstrate its effectiveness.

\*\*\*\*\*

#### Reinforcement Learning with Soft State Aggregation

Satinder Singh, Tommi Jaakkola, Michael Jordan

It is widely accepted that the use of more compact representations than lookup tables is crucial to scaling reinforcement learning (RL) algorithms to real-world problems. Unfortunately almost all of the theory of reinforcement learning assumes lookup table representations. In this paper we address the pressing issue of combining function approximation and RL, and present 1) a function approximator based on a simple extension to state aggregation (a commonly used form of compact representation), namely soft state aggregation, 2) a theory of convergence for RL with arbitrary, but fixed, soft state aggregation, 3) a novel intuitive understanding of the effect of state aggregation on online RL, and 4) a new heuristic adaptive state aggregation algorithm that finds improved compact representations by exploiting the non-discrete nature of soft state aggregation. Preliminary empirical results are also presented.

\*\*\*\*\*

#### A Connectionist Technique for Accelerated Textual Input: Letting a Network Do the Typing

Dean Pomerleau

Each year people spend a huge amount of time typing. The text people type typically contains a tremendous amount of redundancy due to predictable word usage patterns and the text's structure. This paper describes a neural network system called AutoTypist that monitors a person's typing and predicts what will be entered next. AutoTypist displays the most likely subsequent word to the typist, who can accept it with a single keystroke, instead of typing it in its entirety. The multi-layer perceptron at the heart of AutoTypist adapts its predictions of likely subsequent text to the user's word usage pattern, and to the characteristics of the text currently being typed. Increases in typing speed of 2-3% when typing English prose and 10-20% when typing C code have been de

monstrated using the system, suggesting a potential time savings of more than 20 hours per user per year. In addition to increasing typing speed, AutoTypist reduces the number of keystrokes a user must type by a similar amount (2-3% for English, 10- 20% for computer programs). This keystroke savings has the potential to significantly reduce the frequency and severity of repeated stress injuries caused by typing, which are the most common injury suffered in today's office environment.

\*\*\*\*\*

#### Advantage Updating Applied to a Differential Game

Mance E. Harmon, Leemon Baird, A. Harry Klopff

An application of reinforcement learning to a linear-quadratic, differential game is presented. The reinforcement learning system uses a recently developed algorithm, the residual gradient form of advantage updating. The game is a Markov Decision Process (MDP) with continuous time, states, and actions, linear dynamics, and a quadratic cost function. The game consists of two players, a missile and a plane; the missile pursues the plane and the plane evades the missile. The reinforcement learning algorithm for optimal control is modified for differential games in order to find the minimax point, rather than the maximum. Simulation results are compared to the optimal solution, demonstrating that the simulated reinforcement learning system converges to the optimal answer. The performance of both the residual gradient and on-residual gradient forms of advantage updating and Q-learning are compared. The results show that advantage updating converges faster than Q-learning in all simulations. The results also show advantage updating converges regardless of the time step duration; Q-learning is unable to converge as the time step duration ~rows small.

\*\*\*\*\*

#### Pattern Playback in the 90s

Malcolm Slaney

Deciding the appropriate representation to use for modeling human auditory processing is a critical issue in auditory science. While engineers have successfully performed many single-speaker tasks with LPC and spectrogram methods, more difficult problems will need a richer representation. This paper describes a powerful auditory representation known as the correlogram and shows how this non-linear representation can be converted back into sound, with no loss of perceptually important information. The correlogram is interesting because it is a neurophysiologically plausible representation of sound. This paper shows improved methods for spectrogram inversion (conventional pattern playback), inversion of a cochlear model, and inversion of the correlogram representation.

\*\*\*\*\*

#### An Analog Neural Network Inspired by Fractal Block Coding

Fernando Pineda, Andreas Andreou

We consider the problem of decoding block coded data, using a physical dynamical system. We sketch out a decompression algorithm for fractal block codes and then show how to implement a recurrent neural network using physically simple but highly-nonlinear, analog circuit models of neurons and synapses. The nonlinear system has many fixed points, but we have at our disposal a procedure to choose the parameters in such a way that only one solution, the desired solution, is stable. As a partial proof of the concept, we present experimental data from a small system a 16-neuron analog CMOS chip fabricated in a 2 $\mu$ m analog p-well process. This chip operates in the subthreshold regime and, for each choice of parameters, converges to a unique stable state. Each state exhibits a qualitatively fractal shape.

\*\*\*\*\*

#### Phase-Space Learning

Fu-Sheng Tsung, Garrison Cottrell

Existing recurrent net learning algorithms are inadequate. We introduce the conceptual framework of viewing recurrent training as matching vector fields of dynamical systems in phase space. Phase space reconstruction



on techniques make the hidden states explicit, reducing temporal learning to a feed-forward problem. In short, we propose viewing iterated prediction [LF88] as the best way of training recurrent networks on deterministic signals. Using this framework, we can train multiple trajectories, insure their stability, and design arbitrary dynamical systems.

\*\*\*\*\*

An experimental comparison of recurrent neural networks

Bill Horne, C. Giles

Many different discrete-time recurrent neural network architectures have been proposed. However, there has been virtually no effort to compare these architectures experimentally. In this paper we review and categorize many of these architectures and compare how they perform on various classes of simple problems including grammatical inference and nonlinear system identification.

\*\*\*\*\*

Inferring Ground Truth from Subjective Labelling of Venus Images

Padhraic Smyth, Usama Fayyad, Michael Burl, Pietro Perona, Pierre Baldi

In remote sensing applications "ground-truth" data is often used as the basis for training pattern recognition algorithms to generate thematic maps or to detect objects of interest. In practical situations, experts may visually examine the images and provide a subjective noisy estimate of the truth. Calibrating the reliability and bias of expert labellers is a non-trivial problem. In this paper we discuss some of our recent work on this topic in the context of detecting small volcanoes in Magellan SAR images of Venus. Empirical results (using the Expectation-Maximization procedure) suggest that accounting for subjective noise can be quite significant in terms of quantifying both human and algorithm detection performance.

\*\*\*\*\*

Learning Prototype Models for Tangent Distance

Trevor Hastie, Patrice Simard

Simard, LeCun & Denker (1993) showed that the performance of nearest-neighbor classification schemes for handwritten character recognition can be improved by incorporating invariance to the so-called specific transformations in the underlying distance metric - called tangent distance. The resulting classifier, however, can be prohibitively slow and memory intensive due to the large amount of prototypes that need to be stored and used in the distance comparisons. In this paper we develop rich models for representing large subsets of the prototypes. These models are either used singly per class, or as basic building blocks in conjunction with the K-means clustering algorithm.

\*\*\*\*\*

Diffusion of Credit in Markovian Models

Yoshua Bengio, Paolo Frasconi

This paper studies the problem of diffusion in Markovian models, such as hidden Markov models (HMMs) and how it makes very difficult the task of learning of long-term dependencies in sequences. Using results from Markov chain theory, we show that the problem of diffusion is reduced if the transition probabilities approach 0 or 1. Under this condition, standard HMMs have very limited modeling capabilities, but input/output HMMs can still perform interesting computations.

\*\*\*\*\*

The Ni1000: High Speed Parallel VLSI for Implementing Multilayer Perceptrons

Michael Perrone, Leon Cooper

In this paper we present a new version of the standard multilayer perceptron (MLP) algorithm for the state-of-the-art in neural network VLSI implementations: the Intel Ni1000. This new version of the MLP uses a fundamental property of high dimensional spaces which allows the  $l_2$ -norm to be accurately approximated by the  $l_1$ -norm. This approach enables the

standard MLP to utilize the parallel architecture of the Ni1000 to achieve on the order of 40000, 256-dimensional classifications per second.

\*\*\*\*\*

#### Model of a Biological Neuron as a Temporal Neural Network

Sean D. Murphy, Edward W. Kairiss

A biological neuron can be viewed as a device that maps a multidimensional temporal event signal (dendritic postsynaptic activations) into a unidimensional temporal event signal (action potentials). We have designed a network, the Spatio-Temporal Event Mapping (STEM) architecture, which can learn to perform this mapping for arbitrary biological physical models of neurons. Such a network appropriately trained, called a STEM cell, can be used in place of a conventional compartmental model in simulations where only the transfer function is important, such as network simulations. The STEM cell offers advantages over compartmental models in terms of computational efficiency, analytical tractability, and as a framework for VLSI implementations of biological neurons.

\*\*\*\*\*

#### Non-linear Prediction of Acoustic Vectors Using Hierarchical Mixtures of Experts

Steve Waterhouse, Anthony Robinson

In this paper we consider speech coding as a problem of speech modelling. In particular, prediction of parameterised speech over short time segments is performed using the Hierarchical Mixture of Experts (HME) (Jordan & Jacobs 1994). The HME gives two advantages over traditional non-linear function approximators such as the Multi-Layer Perceptron (MLP); a statistical understanding of the operation of the predictor and provision of information about the performance of the predictor in the form of likelihood information and local error bars. These two issues are examined on both toy and real world problems of regression and time series prediction. In the speech coding context, we extend the principle of combining local predictions via the HME to a Vector Quantization scheme in which fixed local codebooks are combined on-line for each observation.

\*\*\*\*\*

#### JPMAX: Learning to Recognize Moving Objects as a Model-fitting Problem

Suzanna Becker

Unsupervised learning procedures have been successful at low-level feature extraction and preprocessing of raw sensor data. So far, however, they have had limited success in learning higher-order representations, e.g., of objects in visual images. A promising approach is to maximize some measure of agreement between the outputs of two groups of units which receive inputs physically separated in space, time or modality, as in (Becker and Hinton, 1992; Becker, 1993; de Sa, 1993). Using the same approach, a much simpler learning procedure is proposed here which discovers features in a single-layer network consisting of several populations of units, and can be applied to multi-layer networks trained one layer at a time. When trained with this algorithm on image sequences of moving geometric objects a two-layer network can learn to perform accurate position-invariant object classification.

\*\*\*\*\*

#### The Electrotonic Transformation: a Tool for Relating Neuronal Form to Function

Nicholas T. Carnevale, Kenneth Y. Tsai, Brenda Claiborne, Thomas Brown

The spatial distribution and time course of electrical signals in neurons have important theoretical and practical consequences. Because it is difficult to infer how neuronal form affects electrical signaling, we have developed a quantitative yet intuitive approach to the analysis of electrotonus. This approach transforms the architecture of the cell from anatomical to electrotonic space, using the logarithm of voltage attenuation as the distance metric. We describe the theory behind this approach and illustrate its use.

\*\*\*\*\*

#### Dynamic Modelling of Chaotic Time Series with Neural Networks

Jose C. Principe, Jyh-Ming Kuo

The auditory system of the barn owl contains several spatial maps. In young barn owls raised with optical prisms over their eyes, these auditory maps are shifted to stay in register with the visual map, suggesting that the visual input imposes a frame of reference on the auditory maps. However, the optic tectum, the first site of convergence of visual with auditory information, is not the site of plasticity for the shift of the auditory maps; the plasticity occurs instead in the inferior colliculus, which contains an auditory map and projects into the optic tectum. We explored a model of the owl remapping in which a global reinforcement signal whose delivery is controlled by visual foveation. A Hebbian learning rule gated by reinforcement learned to appropriately adjust auditory maps. In addition, reinforcement learning preferentially adjusted the weights in the inferior colliculus, as in the owl brain, even though the weights were allowed to change throughout the auditory system. This observation raises the possibility that the site of learning does not have to be genetically specified, but could be determined by how the learning procedure interacts with the network architecture.

\*\*\*\*\*

#### Boltzmann Chains and Hidden Markov Models

Lawrence Saul, Michael Jordan

We propose a statistical mechanical framework for the modeling of discrete time series. Maximum likelihood estimation is done via Boltzmann learning in one-dimensional networks with tied weights. We call these networks Boltzmann chains and show that they contain hidden Markov models (HMMs) as a special case. Our framework also motivates new architectures that address particular shortcomings of HMMs. We look at two such architectures: parallel chains that model feature sets with disparate time scales, and looped networks that model long-term dependencies between hidden states. For these networks, we show how to implement the Boltzmann learning rule exactly, in polynomial time, without resort to simulated or mean-field annealing. The necessary computations are done by exact decimation procedures from statistical mechanics.

\*\*\*\*\*

#### Learning with Product Units

Laurens Leerink, C. Giles, Bill Horne, Marwan Jabri

The TNM staging system has been used since the early 1960's to predict breast cancer patient outcome. In an attempt to increase prognostic accuracy, many putative prognostic factors have been identified.

Because the TNM stage model can not accommodate these new factors, the proliferation of factors in breast cancer has led to clinical confusion. What is required is a new computerized prognostic system that can test putative prognostic factors and integrate the predictive factors with the TNM variables in order to increase prognostic accuracy. Using the area under the curve of the receiver operating characteristic, we compare the accuracy of the following predictive models in terms of five year breast cancer-specific survival: pTNM staging system, principal component analysis, classification and regression trees, logistic regression, cascade correlation neural network, conjugate gradient descent neural, probabilistic neural network, and backpropagation neural network. Several statistical models are significantly more ac-

\*\*\*\*\*

#### Efficient Methods for Dealing with Missing Data in Supervised Learning

Volker Tresp, Ralph Neuneier, Subutai Ahmad

We present efficient algorithms for dealing with the problem of missing inputs (incomplete feature vectors) during training and recall. Our approach is based on the approximation of the input data distribution using Parzen windows. For recall, we obtain closed form solutions for arbitrary feedforward networks. For training, we show how the backpropagation step for an incomplete pattern can be approximated by a weighted averaged backpropagation step. The complexity of the solutions for training and recall is independent of the number of missing features. We verify our theoretical

results using one classification and one regression problem.

\*\*\*\*\*

#### Predictive Coding with Neural Nets: Application to Text Compression

Jürgen Schmidhuber, Stefan Heil

To compress text files, a neural predictor network  $P$  is used to approximate the conditional probability distribution of possible "next characters", given  $n$  previous characters.  $P$ 's outputs are fed into standard coding algorithms that generate short codes for characters with high predicted probability and long codes for highly unpredictable characters. Tested on short German newspaper articles, our method outperforms widely used Lempel-Ziv algorithms (used in UNIX functions such as "compress" and "gzip").

\*\*\*\*\*

#### Computational Structure of coordinate transformations: A generalization study

Zoubin Ghahramani, Daniel M. Wolpert, Michael Jordan

One of the fundamental properties that both neural networks and the central nervous system share is the ability to learn and generalize from examples. While this property has been studied extensively in the neural network literature it has not been thoroughly explored in human perceptual and motor learning. We have chosen a coordinate transformation system-the visuomotor map which transforms visual coordinates into motor coordinates-to study the generalization effects of learning new input-output pairs. Using a paradigm of computer controlled altered visual feedback, we have studied the generalization of the visuomotor map subsequent to both local and context-dependent remappings. A local remapping of one or two input-output pairs induced a significant global, yet decaying, change in the visuomotor map, suggesting a representation for the map composed of units with large functional receptive fields. Our study of context-dependent remappings indicated that a single point in visual space can be mapped to two different finger locations depending on a context variable-the starting point of the movement. Furthermore, as the context is varied there is a gradual shift between the two remappings, consistent with two visuomotor modules being learned and gated smoothly with the context.

\*\*\*\*\*

#### Recognizing Handwritten Digits Using Mixtures of Linear Models

Geoffrey E. Hinton, Michael Revow, Peter Dayan

We construct a mixture of locally linear generative models of a collection of pixel-based images of digits, and use them for recognition. Different models of a given digit are used to capture different styles of writing, and new images are classified by evaluating their log-likelihoods under each model. We use an EM-based algorithm in which the M-step is computationally straightforward principal components analysis (PCA). Incorporating tangent-plane information [12] about expected local deformations only requires adding tangent vectors into the sample covariance matrices for the PCA, and it demonstrably improves performance.

\*\*\*\*\*

#### A Critical Comparison of Models for Orientation and Ocular Dominance Columns in the Striate Cortex

E. Erwin, K. Obermayer, K. Schulten

More than ten of the most prominent models for the structure and for the activity dependent formation of orientation and ocular dominance columns in the striate cortex have been evaluated. We implemented those models on parallel machines, we extensively explored parameter space, and we quantitatively compared model predictions with experimental data which were recorded optically from macaque striate cortex. In our contribution we present a summary of our results to date. Briefly, we find that (i) despite apparent differences, many models are based on similar principles and, consequently, make similar predictions, (ii) certain "pattern models" as well as the developmental "correlation-based learning" models disagree with the experimental data, and (iii)

i) of the models we have investigated, "competitive Hebbian" models and the recent model of Swindale provide the best match with experimental data.

\*\*\*\*\*

#### Classifying with Gaussian Mixtures and Clusters

Nanda Kambhatla, Todd Leen

In this paper, we derive classifiers which are winner-take-all (WTA) approximations to a Bayes classifier with Gaussian mixtures for class conditional densities. The derived classifiers include clustering based algorithms like LVQ and k-Means. We propose a constrained rank Gaussian mixtures model and derive a WTA algorithm for it. Our experiments with two speech classification tasks indicate that the constrained rank model and the WTA approximations improve the performance over the unconstrained models.

\*\*\*\*\*

#### Anatomical origin and computational role of diversity in the response properties of cortical neurons

Kalanit Spector, Shimon Edelman, Rafael Malach

The maximization of diversity of neuronal response properties has been recently suggested as an organizing principle for the formation of such prominent features of the functional architecture of the brain as the cortical columns and the associated patchy projection patterns (Malach, 1994).

We show that (1) maximal diversity is attained when the ratio of dendritic and axonal arbor sizes is equal to one, as found in many cortical areas and across species (Lund et al., 1993; Malach, 1994), and (2) that maximization of diversity leads to better performance in systems of receptive fields implementing steerable/shiftable filters, and in matching spatially distributed signals, a problem that arises in many high-level visual tasks.

\*\*\*\*\*

#### Synchrony and Desynchrony in Neural Oscillator Networks

Deliang Wang, David Terman

An novel class of locally excitatory, globally inhibitory oscillator networks is proposed. The model of each oscillator corresponds to a standard relaxation oscillator with two time scales. The network exhibits a mechanism of selective gating, whereby an oscillator jumping up to its active phase rapidly recruits the oscillators stimulated by the same pattern, while preventing others from jumping up. We show analytically that with the selective gating mechanism the network rapidly achieves both synchronization within blocks of oscillators that are stimulated by connected regions and desynchronization between different blocks. Computer simulations demonstrate the network's promising ability for segmenting multiple input patterns in real time. This model lays a physical foundation for the oscillatory correlation theory of feature binding, and may provide an effective computational framework for scene segmentation and figure/ground segregation.

\*\*\*\*\*

#### A Computational Model of Prefrontal Cortex Function

Todd Braver, Jonathan D. Cohen, David Servan-Schreiber

Accumulating data from neurophysiology and neuropsychology have suggested two information processing roles for prefrontal cortex (PFC): 1) short-term active memory; and 2) inhibition. We present a new behavioral task and a computational model which were developed in parallel.

The task was developed to probe both of these prefrontal functions simultaneously, and produces a rich set of behavioral data that act as constraints on the model. The model is implemented in continuous-time, thus providing a natural framework in which to study the temporal dynamics of processing in the task. We show how the model can be used to examine the behavioral consequences of neuromodulation in PFC. Specifically, we use the model to make novel and testable predictions regarding the behavioral performance of schizophrenics, who are hypothesized to suffer from reduced dopaminergic tone in this brain area.

\*\*\*\*\*

### Combining Estimators Using Non-Constant Weighting Functions

Volker Tresp, Michiaki Taniguchi

This paper discusses the linearly weighted combination of estimators in which the weighting functions are dependent on the input. We show that the weighting functions can be derived either by evaluating the input dependent variance of each estimator or by estimating how likely it is that a given estimator has seen data in the region of the input space close to the input pattern. The latter solution is closely related to the mixture of experts approach and we show how learning rules for the mixture of experts can be derived from the theory about learning with missing features. The presented approaches are modular since the weighting functions can easily be modified (no retraining) if more estimators are added. Furthermore, it is easy to incorporate estimators which were not derived from data such as expert systems or algorithms.

\*\*\*\*\*

### Stochastic Dynamics of Three-State Neural Networks

Toru Ohira, Jack Cowan

We present here an analysis of the stochastic neurodynamics of a neural network composed of three-state neurons described by a master equation. An outer-product representation of the master equation is employed. In this representation, an extension of the analysis from two to three-state neurons is easily performed. We apply this formalism with approximation schemes to a simple three-state network and compare the results with Monte Carlo simulations.

\*\*\*\*\*

### On the Computational Utility of Consciousness

Donald Mathis, Michael C. Mozer

We propose a computational framework for understanding and modeling human consciousness. This framework integrates many existing theoretical perspectives, yet is sufficiently concrete to allow simulation experiments. We do not attempt to explain qualia (subjective experience), but instead ask what differences exist within the cognitive information processing system when a person is conscious of mentally-represented information versus when that information is unconscious. The central idea we explore is that the contents of consciousness correspond to temporally persistent states in a network of computational modules. Three simulations are described illustrating that the behavior of persistent states in the models corresponds roughly to the behavior of conscious states people experience when performing similar tasks. Our simulations show that periodic settling to persistent (i.e., conscious) states improves performance by cleaning up inaccuracies and noise, forcing decisions, and helping keep the system on track toward a solution.

\*\*\*\*\*

### Ocular Dominance and Patterned Lateral Connections in a Self-Organizing Model of the Primary Visual Cortex

Joseph Sirosh, Risto Miikkulainen

A neural network model for the self-organization of ocular dominance and lateral connections from binocular input is presented. The self-organizing process results in a network where (1) afferent weights of each neuron organize into smooth hill-shaped receptive fields primarily on one of the retinas, (2) neurons with common eye preference form connected, intertwined patches, and (3) lateral connections primarily link regions of the same eye preference. Similar self-organization of cortical structures has been observed experimentally in strabismic kittens. The model shows how patterned lateral connections in the cortex may develop based on correlated activity and explains why lateral connection patterns follow receptive field properties such as ocular dominance.

\*\*\*\*\*

### Effects of Noise on Convergence and Generalization in Recurrent Networks

Kam Jim, Bill Horne, C. Giles

We introduce and study methods of inserting synaptic noise into dynamically-driven recurrent neural networks and show that applying a controlled amount of noise during training may improve convergence and generalization. In addition, we analyze the effects of each noise parameter (additive vs. multiplicative, cumulative vs. non-cumulative, per time step vs. per string) and predict that best overall performance can be achieved by injecting additive noise at each time step. Extensive simulations on learning the dual parity grammar from temporal strings substantiate these predictions.

\*\*\*\*\*

#### An Integrated Architecture of Adaptive Neural Network Control for Dynamic Systems

Ke Liu, Robert Tokar, Brain McVey

In this study, an integrated neural network control architecture for nonlinear dynamic systems is presented. Most of the recent emphasis in the neural network control field has no error feedback as the control input, which rises the lack of adaptation problem. The integrated architecture in this paper combines feed forward control and error feedback adaptive control using neural networks. The paper reveals the different internal functionality of these two kinds of neural network controllers for certain input styles, e.g., state feedback and error feedback. With error feedback, neural network controllers learn the slopes or the gains with respect to the error feedback, producing an error driven adaptive control systems. The results demonstrate that the two kinds of control scheme can be combined to realize their individual advantages. Testing with disturbances added to the plant shows good tracking and adaptation with the integrated neural control architecture.

\*\*\*\*\*

#### Implementation of Neural Hardware with the Neural VLSI of URAN in Applications with Reduced Representations

Il Han, Ki-Chul Kim, Hwang-Soo Lee

Implement Korean

\*\*\*\*\*

#### Estimating Conditional Probability Densities for Periodic Variables

Chris M. Bishop, Claire Legleye

Most of the common techniques for estimating conditional probability densities are inappropriate for applications involving periodic variables. In this paper we introduce three novel techniques for tackling such problems, and investigate their performance using synthetic data. We then apply these techniques to the problem of extracting the distribution of wind vector directions from radar scatterometer data gathered by a remote-sensing satellite.

\*\*\*\*\*

#### Analysis of Unstandardized Contributions in Cross Connected Networks

Thomas Shultz, Yuriiko Oshima-Takane, Yoshio Takane

Understanding knowledge representations in neural nets has been a difficult problem. Principal components analysis (PCA) of contributions (products of sending activations and connection weights) has yielded valuable insights into knowledge representations, but much of this work has focused on the correlation matrix of contributions. The present work shows that analyzing the variance-covariance matrix of contributions yields more valid insights by taking account of weights.

\*\*\*\*\*

#### A Rigorous Analysis of Linsker-type Hebbian Learning

J. Feng, H. Pan, V. P. Roychowdhury

We propose a novel rigorous approach for the analysis of Linsker's unsupervised Hebbian learning network. The behavior of this model is determined by the underlying nonlinear dynamics which are parameterized by a set of parameters originating from the Hebbian rule and the arbor density of the synapses. These parameters determine the presence or absence of a specific receptive field (also referred to as a 'connection pattern') as a saturated fixed point attract

or of the model. In this paper, we perform a qualitative analysis of the underlying nonlinear dynamics over the parameter space, determine the effects of the system parameters on the emergence of various receptive fields, and predict precisely within which parameter regime the network will have the potential to develop a specially designated connection pattern. In particular, this approach exposes, for the first time, the crucial role played by the synaptic density functions, and provides a complete precise picture of the parameter space that defines the relationships among the different receptive fields. Our theoretical predictions are confirmed by numerical simulations.

\*\*\*\*\*

#### Associative Decorrelation Dynamics: A Theory of Self-Organization and Optimization in Feedback Networks

Dawei Dong

This paper outlines a dynamic theory of development and adaptation in neural networks with feedback connections. Given input ensemble, the connections change in strength according to an associative learning rule and approach a stable state where the neuronal outputs are decorrelated. We apply this theory to primary visual cortex and examine the implications of the dynamical decorrelation of the activities of orientation selective cells by the intracortical connections.

The theory gives a unified and quantitative explanation of the psychophysical experiments on orientation contrast and orientation adaptation. Using only one parameter, we achieve good agreements between the theoretical predictions and the experimental data.

\*\*\*\*\*

#### Visual Speech Recognition with Stochastic Networks

Javier Movellan

This paper presents ongoing work on a speaker independent visual speech recognition system. The work presented here builds on previous research efforts in this area and explores the potential use of simple hidden Markov models for limited vocabulary, speaker independent visual speech recognition. The task at hand is recognition of the first four English digits, a task with possible applications in car-phone images were modeled as mixtures of independent dialing. The Gaussian distributions, and the temporal dependencies were captured with standard left-to-right hidden Markov models. The results indicate that simple hidden Markov models may be used to successfully recognize relatively unprocessed image sequences. The system achieved performance levels equivalent to untrained humans when asked to recognize the first four English digits.

\*\*\*\*\*

#### Finding Structure in Reinforcement Learning

Sebastian Thrun, Anton Schwartz

Reinforcement learning addresses the problem of learning to select actions in order to maximize one's performance in unknown environments. To scale reinforcement learning to complex real-world tasks, such as typically studied in AI, one must ultimately be able to discover the structure in the world, in order to abstract away the myriad of details and to operate in more tractable problem spaces. This paper presents the SKILLS algorithm. SKILLS discovers skills, which are partially defined action policies that arise in the context of multiple, related tasks. Skills collapse whole action sequences into single operators. They are learned by minimizing the compactness of action policies, using a description length argument on their representation. Empirical results in simple grid navigation tasks illustrate the successful discovery of structure in reinforcement learning.

\*\*\*\*\*

#### Active Learning with Statistical Models

David Cohn, Zoubin Ghahramani, Michael Jordan

For many types of learners one can compute the statistically "optimal" way to select data. We review how these techniques have been used with feedforward neural networks [MacKay, 1992; Cohn, 1994]. We then show how the same pri-



nciples may be used to select data for two alternative, statistically-based learning architectures: mixtures of Gaussians and locally weighted regression. While the techniques for neural networks are expensive and approximate, the techniques for mixtures of Gaussians and locally weighted regression are both efficient and accurate.

\*\*\*\*\*

From Data Distributions to Regularization in Invariant Learning

Todd Leen

Ideally pattern recognition machines provide constant output when the inputs are transformed under a group of desired invariances. These invariances can be achieved by enhancing the training data to include examples of inputs transformed by elements of  $G$ , while leaving the corresponding targets unchanged.

Alternatively the cost function for training can include a regularization term that penalizes changes in the output when the input is transformed under the group.

\*\*\*\*\*

An Input Output HMM Architecture

Yoshua Bengio, Paolo Frasconi

We introduce a recurrent architecture having a modular structure and we formulate a training procedure based on the EM algorithm. The resulting model has similarities to hidden Markov models, but supports recurrent networks processing style and allows to exploit the supervised learning paradigm while using maximum likelihood estimation.

\*\*\*\*\*

Grouping Components of Three-Dimensional Moving Objects in Area MST of Visual Cortex

Richard Zemel, Terrence J. Sejnowski

Many cells in the dorsal part of the medial superior temporal (MST) area of visual cortex respond selectively to spiral flow patterns-specific combinations of expansion/contraction and rotation motions. Previous investigators have suggested that these cells may represent self-motion. Spiral patterns can also be generated by the relative motion of the observer and a particular object. An MST cell may then account for some portion of the complex flow field, and the set of active cells could encode the entire flow; in this manner, MST effectively segments moving objects. Such a grouping operation is essential in interpreting scenes containing several independent moving objects and observer motion. We describe a model based on the hypothesis that the selective tuning of MST cells reflects the grouping of object components undergoing coherent motion. Inputs to the model were generated from sequences of ray-traced images that simulated realistic motions, combining observer motion, eye movements, and independent object motion. The input representation was modeled after response properties of neurons in area MT, which provides the primary input to area MST. After applying an unsupervised learning algorithm, the units became tuned to patterns signaling coherent motion. The results match many of the known properties of MST cells and are consistent with recent studies indicating that these cells process 3-D object motion information.

\*\*\*\*\*

Higher Order Statistical Decorrelation without Information Loss

Gustavo Deco, Wilfried Brauer

A neural network learning paradigm based on information theory is proposed as a way to perform in an unsupervised fashion, redundancy reduction among the elements of the output layer without loss of information from the sensory input. The model developed performs nonlinear decorrelation up to higher orders of the cumulant tensors and results in probabilistically independent components of the output layer. This means that we don't need to assume Gaussian distribution neither at the input nor at the output. The theory presented is related to the unsupervised-learning theory of Barlow, which

proposes redundancy reduction as the goal of cognition. When nonlinear units are used nonlinear principal component analysis is obtained. In this case nonlinear manifolds can be reduced to minimum dimension manifolds. If such units are used the network performs a generalized principal component analysis in the sense that non-Gaussian distributions can be linearly decorrelated and higher orders of the correlation tensors are also taken into account. The basic structure of the architecture involves a general transformation that is volume conserving and therefore the entropy, yielding a map without loss of information. Minimization of the mutual information among the output neurons eliminates the redundancy between the outputs and results in statistical decorrelation of the extracted features. This is known as factorially learning.

\*\*\*\*\*

#### Sample Size Requirements for Feedforward Neural Networks

Michael Turmon, Terrence L. Fine

We estimate the number of training samples required to ensure that the performance of a neural network on its training data matches that obtained when fresh data is applied to the network. Existing estimates are higher by orders of magnitude than practice indicates. This work seeks to narrow the gap between theory and practice by transforming the problem into determining the distribution of the supremum of a random field in the space of weight vectors, which in turn is attacked by application of a recent technique called the Poisson clumping heuristic.

\*\*\*\*\*

#### Generalisation in Feedforward Networks

Adam Kowalczyk, Herman Ferrá

We discuss a model of consistent learning with an additional restriction on the probability distribution of training samples, the target concept and hypothesis class. We show that the model provides a significant improvement on the upper bounds of sample complexity, i.e. the minimal number of random training samples allowing a selection of the hypothesis with a predefined accuracy and confidence. Further, we show that the model has the potential for providing a finite sample complexity even in the case of infinite VC-dimension as well as for a sample complexity below VC-dimension. This is achieved by linking sample complexity to an "average" number of implementable dichotomies of a training sample rather than the maximal size of a shattered sample, i.e. VC-dimension.

\*\*\*\*\*

#### The Use of Dynamic Writing Information in a Connectionist On-Line Cursive Handwriting Recognition System

Stefan Manke, Michael Finke, Alex Waibel

In this paper we present NPen++, a connectionist system for writer independent, large vocabulary on-line cursive handwriting recognition. This system combines a robust input representation, which preserves the dynamic writing information, with a neural network architecture, a so called Multi-State Time Delay Neural Network (MS-TDNN), which integrates recognition and segmentation in a single framework. Our preprocessing transforms the original coordinate sequence into a (still temporal) sequence of feature vectors, which combine strictly local features, like curvature or writing direction, with a bitmap-like representation of the coordinate's proximity. The MS-TDNN architecture is well suited for handling temporal sequences as provided by this input representation. Our system is tested both on writer dependent and writer independent tasks with vocabulary sizes ranging from 400 up to 20,000 words. For example, on a 20,000 word vocabulary we achieve word recognition rates up to 88.9% (writer dependent) and 84.1% (writer independent) without using any language models.

\*\*\*\*\*

#### Direct Multi-Step Time Series Prediction Using TD( $\lambda$ )

Peter T. Kazlas, Andreas Weigend

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.

\*\*\*\*\*

#### Capacity and Information Efficiency of a Brain-like Associative Net

Bruce Graham, David Willshaw

We have determined the capacity and information efficiency of an associative net configured in a brain-like way with partial connectivity and noisy input cues. Recall theory was used to calculate the capacity when pattern recall is achieved using a winners-take-all strategy. Transforming the dendritic sum according to input activity and unit usage can greatly increase the capacity of the associative net under these conditions. For moderately sparse patterns, maximum information efficiency is achieved with very low connectivity levels (~ 10%). This corresponds to the level of connectivity commonly seen in the brain and invites speculation that the brain is connected in the most information efficient way.

\*\*\*\*\*

#### SARDNET: A Self-Organizing Feature Map for Sequences

Daniel L. James, Risto Miikkulainen

A self-organizing neural network for sequence classification called SARDNET is described and analyzed experimentally. SARDNET extends the Kohonen Feature Map architecture with activation retention and decay in order to create unique distributed response patterns for different sequences.

SARDNET yields extremely dense yet descriptive representations of sequential input in very few training iterations. The network has proven successful on mapping arbitrary sequences of binary and real numbers, as well as phonemic representations of English words. Potential applications include isolated spoken word recognition and cognitive science models of sequence processing.

\*\*\*\*\*

#### Deterministic Annealing Variant of the EM Algorithm

Naonori Ueda, Ryohei Nakano

We present a deterministic annealing variant of the EM algorithm for maximum likelihood parameter estimation problems. In our approach, the EM process is reformulated as the problem of minimizing the thermodynamic free energy by using the principle of maximum entropy and statistical mechanics analogy. Unlike simulated annealing approaches, this minimization is deterministically performed. Moreover, the derived algorithm, unlike the conventional EM algorithm, can obtain better estimates free of the initial parameter values.

\*\*\*\*\*

#### A Non-linear Information Maximisation Algorithm that Performs Blind Separation

Anthony Bell, Terrence J. Sejnowski

A new learning algorithm is derived which performs online stochastic gradient ascent in the mutual information between outputs and inputs of a network. In the absence of a priori knowledge about the 'signal' and 'noise' components of the input, propagation of information depends on calibrating network non-linearities to the detailed higher-order moments of the input density functions. By incidentally minimising mutual information between outputs, as well as maximising their individual entropies, the network 'factorises' the input into independent components. As an example application, we have achieved near-perfect separation of ten digitally mixed speech signals. Our simulations lead us to believe that our network performs better at blind separation than the Herault-Jutten network, reflecting the fact that it is derived rigorously from the mutual information objective.

\*\*\*\*\*

#### Pulstream Synapses with Non-Volatile Analogue Amorphous-Silicon Memories

A. Holmes, Alan Murray, Stephen Churcher, J. Hajto, M. Rose

This paper presents results from the first use of neural networks for

the real-time feedback control of high temperature plasmas in a tokamak fusion experiment. The tokamak is currently the principal experimental device for research into the magnetic confinement approach to controlled fusion. In the tokamak, hydrogen plasmas, at temperatures of up to 100 Million K, are confined by strong magnetic fields.

Accurate control of the position and shape of the plasma boundary requires real-time feedback control of the magnetic field structure on a time-scale of a few tens of microseconds. Software simulations have demonstrated that a neural network approach can give significantly better performance than the linear technique currently used on most tokamak experiments. The practical application of the neural network approach requires high-speed hardware, for which a fully parallel implementation of the multilayer perceptron, using a hybrid of digital and analogue technology, has been developed.

\*\*\*\*\*

#### Dynamic Cell Structures

Jörg Bruske, Gerald Sommer

Dynamic Cell Structures (DCS) represent a family of artificial neural architectures suited both for unsupervised and supervised learning. They belong to the recently [Martinetz94] introduced class of Topology Representing Networks (TRN) which build perfectly topology preserving feature maps.

DCS employ a modified Kohonen learning rule in conjunction with competitive Hebbian learning. The Kohonen type learning rule serves to adjust the synaptic weight vectors while Hebbian learning establishes a dynamic lateral connection structure between the units reflecting the topology of the feature manifold. In case of supervised learning, i.e. function approximation, each neural unit implements a Radial Basis Function, and an additional layer of linear output units adjusts according to a delta-rule. DCS is the first RBF-based approximation scheme attempting to concurrently learn and utilize a perfectly topology preserving map for improved performance. Simulations on a selection of CMU-Benchmarks indicate that the DCS idea applied to the Growing Cell Structure algorithm [Fritzke93] leads to an efficient and elegant algorithm that can beat conventional models on similar tasks.

\*\*\*\*\*

#### Single Transistor Learning Synapses

Paul Hasler, Chris Diorio, Bradley Minch, Carver Mead

We describe single-transistor silicon synapses that compute, learn, and provide non-volatile memory retention. The single transistor synapses simultaneously perform long term weight storage, compute the product of the input and the weight value, and update the weight value according to a Hebbian or a backpropagation learning rule. Memory is accomplished via charge storage on polysilicon floating gates, providing long-term retention without refresh. The synapses efficiently use the physics of silicon to perform weight updates; the weight value is increased using tunneling and the weight value decreases using hot electron injection. The small size and low power operation of single transistor synapses allows the development of dense synaptic arrays. We describe the design, fabrication, characterization, and modeling of an array of single transistor synapses. When the steady state source current is used as the representation of the weight value, both the incrementing and decrementing functions are proportional to a power of the source current. The synaptic array was fabricated in the standard 2.1µm double - poly, analog process available from MOSIS.

\*\*\*\*\*

#### Comparing the prediction accuracy of artificial neural networks and other statistical models for breast cancer survival

Harry B. Burke, David B. Rosen, Philip H. Goodman

The TNM staging system has been used since the early 1960's to predict breast cancer patient outcome. In an attempt to increase prognostic accuracy, many putative prognostic factors have been identified.

Because the TNM stage model can not accommodate these new factors, the proliferation of factors in breast cancer has led to clinical confusion. What is required is a new computerized prognostic system that can test putative prognostic factors and integrate the predictive factors with the TNM variables in order to increase prognostic accuracy. Using the area under the curve of the receiver operating characteristic, we compare the accuracy of the following predictive models in terms of five year breast cancer-specific survival: pTNM staging system, principal component analysis, classification and regression trees, logistic regression, cascade correlation neural network, conjugate gradient descent neural, probabilistic neural network, and backpropagation neural network. Several statistical models are significantly more ac-

\*\*\*\*\*

Learning direction in global motion: two classes of psychophysically-motivated models

V. Sundareswaran, Lucia Vaina

Perceptual learning is defined as fast improvement in performance and retention of the learned ability over a period of time. In a set of psychophysical experiments we demonstrated that perceptual learning occurs for the discrimination of direction in stochastic motion stimuli.

Here we model this learning using two approaches: a clustering model that learns to accommodate the motion noise, and an averaging model that learns to ignore the noise. Simulations of the models show performance similar to the psychophysical results.

\*\*\*\*\*

On-line Learning of Dichotomies

N. Barkai, H. Seung, H. Sompolinsky

The performance of on-line algorithms for learning dichotomies is studied. In on-line learning, the number of examples  $P$  is equivalent to the learning time, since each example is presented only once. The learning curve, or generalization error as a function of  $P$ , depends on the schedule at which the learning rate is lowered. For a target that is a perceptron rule, the learning curve of the perceptron algorithm can decrease as fast as  $p^{-1}$ , if the schedule is optimized. If the target is not realizable by a perceptron, the perceptron algorithm does not generally converge to the solution with lowest generalization error. For the case of unrealizability due to a simple output noise, we propose a new on-line algorithm for a perceptron yielding a learning curve that can approach the optimal generalization error as fast as  $p^{-1/2}$ . We then generalize the perceptron algorithm to any class of thresholded smooth functions learning a target from that class. For "well-behaved" input distributions, if this algorithm converges to the optimal solution, its learning curve can decrease as fast as  $p^{-1}$ .

\*\*\*\*\*

Asymptotics of Gradient-based Neural Network Training Algorithms

Sayandev Mukherjee, Terrence L. Fine

We study the asymptotic properties of the sequence of iterates of weight-vector estimates obtained by training a multilayer feed forward neural network with a basic gradient-descent method using a fixed learning constant and no batch-processing. In the one-dimensional case, an exact analysis establishes the existence of a limiting distribution that is not Gaussian in general.

For the general case and small learning constant, a linearization approximation permits the application of results from the theory of random matrices to again establish the existence of a limiting distribution. We study the first few moments of this distribution to compare and contrast the results of our analysis with those of techniques of stochastic approximation.

\*\*\*\*\*

Convergence Properties of the K-Means Algorithms

Léon Bottou, Yoshua Bengio

This paper studies the convergence properties of the well known K-Means clustering algorithm. The K-Means algorithm can be described either as a gradient

ient descent algorithm or by slightly extending the mathematics of the EM algorithm to this hard threshold case. We show that the K-Means algorithm actually minimizes the quantization error using the very fast Newton algorithm.

\*\*\*\*\*

Using Voice Transformations to Create Additional Training Talkers for Word Spotting

Eric Chang, Richard P. Lippmann

Speech recognizers provide good performance for most users but the error rate often increases dramatically for a small percentage of talkers who are "different" from those talkers used for training. One expensive solution to this problem is to gather more training data in an attempt to sample these outlier users. A second solution, explored in this paper, is to artificially enlarge the number of training talkers by transforming the speech of existing training talkers. This approach is similar to enlarging the training set for OCR digit recognition by warping the training digit images, but is more difficult because continuous speech has a much larger number of dimensions (e.g. linguistic, phonetic, style, temporal, spectral) that differ across talkers. We explored the use of simple linear spectral warping to enlarge a 48-talker training data base used for word spotting. The average detection rate overall was increased by 2.9 percentage points (from 68.3% to 71.2%) for male speakers and 2.5 percentage points (from 64.8% to 67.3%) for female speakers. This increase is small but similar to that obtained by doubling the amount of training data.

\*\*\*\*\*

Forward dynamic models in human motor control: Psychophysical evidence

Daniel M. Wolpert, Zoubin Ghahramani, Michael Jordan

Based on computational principles, with as yet no direct experimental validation, it has been proposed that the central nervous system (CNS) uses an internal model to simulate the dynamic behavior of the motor system in planning, control and learning (Sutton and Barto, 1981; Ito, 1984; Kawato et al., 1987; Jordan and Rumelhart, 1992; Miall et al., 1993). We present experimental results and simulations based on a novel approach that investigates the temporal propagation of errors in the sensorimotor integration process. Our results provide direct support for the existence of an internal model.

\*\*\*\*\*

Direction Selectivity In Primary Visual Cortex Using Massive Intracortical Connections

Humbert Suarez, Christof Koch, Rodney Douglas

Almost all models of orientation and direction selectivity in visual cortex are based on feedforward connection schemes, where geniculate input provides all excitation to both pyramidal and inhibitory neurons. The latter neurons then suppress the response of the former for non-optimal stimuli. However, anatomical studies show that up to 90 % of the excitatory synaptic input onto any cortical cell is provided by other cortical cells. The massive excitatory feedback nature of cortical circuits is embedded in the canonical microcircuit of Douglas & Martin (1991). We here investigate analytically and through biologically realistic simulations the functioning of a detailed model of this circuitry, operating in a hysteretic mode. In the model, weak geniculate input is dramatically amplified by intracortical excitation, while inhibition has a dual role: (i) to prevent the early geniculate-induced excitation in the null direction and (ii) to restrain excitation and ensure that the neurons fire only when the stimulus is in their receptive-field. Among the

\*\*\*\*\*

Bayesian Query Construction for Neural Network Models

Gerhard Paass, Jörg Kindermann

If data collection is costly, there is much to be gained by actively selecting particularly informative data points in a sequential way. In a Bayesian decision-theoretic framework we develop a query selection criterion

terion which explicitly takes into account the intended use of the model predictions. By Markov Chain Monte Carlo methods the necessary quantities can be approximated to a desired precision. As the number of data points grows, the model complexity is modified by a Bayesian model selection strategy. The proper ties of two versions of the criterion are demonstrated in numerical experiments.

\*\*\*\*\*

#### A Silicon Axon

Bradley Minch, Paul Hasler, Chris Diorio, Carver Mead

We present a silicon model of an axon which shows promise as a building block for pulse-based neural computations involving correlations of pulses across both space and time. The circuit shares a number of features with its biological counterpart including an excitation threshold, a brief refractory period after pulse completion, pulse amplitude restoration, and pulse width restoration. We provide a simple explanation of circuit operation and present data from a chip fabricated in a standard 2J1m CMOS process through the MOS Implementation Service (MOSIS). We emphasize the necessity of the restoration of the width of the pulse in time for stable propagation in axons.

\*\*\*\*\*

#### Plasticity-Mediated Competitive Learning

Nicol Schraudolph, Terrence J. Sejnowski

Differentiation between the nodes of a competitive learning network work is conventionally achieved through competition on the basis of neural activity. Simple inhibitory mechanisms are limited to sparse representations, while decorrelation and factorization schemes that support distributed representations are computationally unattractive. By letting neural plasticity mediate the competitive interaction instead, we obtain diffuse, nonadaptive alternatives for fully distributed representations. We use this technique to simplify and improve our binary information gain optimization algorithm for feature extraction (Schraudolph and Sejnowski, 1993); the same approach could be used to improve other learning algorithms.

\*\*\*\*\*

#### Active Learning for Function Approximation

Kah Sung, Partha Niyogi

We develop a principled strategy to sample a function optimally for function approximation tasks within a Bayesian framework. Using ideas from optimal experiment design, we introduce an objective function (incorporating both bias and variance) to measure the degree of approximation, and the potential utility of the data points towards optimizing this objective. We show how the general strategy can be used to derive precise algorithms to select data for two cases: learning unit step functions and polynomial functions. In particular, we investigate whether such active algorithms can learn the target with fewer examples. We obtain theoretical and empirical results to suggest that this is the case.

\*\*\*\*\*

#### Patterns of damage in neural networks: The effects of lesion area, shape and number

Eytan Ruppin, James Reggia

Current understanding of the effects of damage on neural networks is rudimentary, even though such understanding could lead to important insights concerning neurological and psychiatric disorders. Motivated by this consideration, we present a simple analytical framework for estimating the functional damage resulting from focal structural lesions to a neural network. The effects of focal lesions of varying area, shape and number on the retrieval capacities of a spatially-organized associative memory. Although our analytical results are based on some approximations, they correspond well with simulation results. This study sheds light on some important features characterizing the clinical manifestations of multi-infarct dementia, including the strong association between the number of infarct

ts and the prevalence of dementia after stroke, and the 'multiplicative' interaction that has been postulated to occur between Alzheimer's disease and multi-infarct dementia.

\*\*\*\*\*

#### A Study of Parallel Perturbative Gradient Descent

D. Lippe, Joshua Alspector

We have continued our study of a parallel perturbative learning method [Alspector et al., 1993] and implications for its implementation in analog VLSI. Our new results indicate that, in most cases, a single parallel perturbation (per pattern presentation) of the function parameters (weights in a neural network) is theoretically the best course. This is not true, however, for certain problems and may not generally be true when faced with issues of implementation such as limited precision. In these cases, multiple parallel perturbations may be best as indicated in our previous results.

\*\*\*\*\*

#### A Neural Model of Delusions and Hallucinations in Schizophrenia

Eytan Ruppin, James Reggia, David Horn

We implement and study a computational model of Stevens' [1992] theory of the pathogenesis of schizophrenia. This theory hypothesizes that the onset of schizophrenia is associated with reactive synaptic regeneration occurring in brain regions receiving degenerating temporal lobe projections. Concentrating on one such area, the frontal cortex, we model a frontal module as an associative memory neural network whose input synapses represent incoming temporal projections. We analyze how, in the face of weakened external input projections, compensatory strengthening of internal synaptic connections and increased noise levels can maintain memory capacities (which are generally preserved in schizophrenia). However, These compensatory changes adversely lead to spontaneous, biased retrieval of stored memories, which corresponds to the occurrence of schizophrenic delusions and hallucinations without any apparent external trigger, and for their tendency to concentrate on just few central themes. Our results explain why these symptoms tend to wane as schizophrenia progresses, and why delayed therapeutic intervention leads to a much slower response.

\*\*\*\*\*

#### Correlation and Interpolation Networks for Real-time Expression Analysis/Synthesis

Trevor Darrell, Irfan Essa, Alex Pentland

We describe a framework for real-time tracking of facial expressions that uses neurally-inspired correlation and interpolation methods. A distributed view-based representation is used to characterize facial state, and is computed using a replicated correlation network. The ensemble response of the set of view correlation scores is input to a network based interpolation method, which maps perceptual state to motor control states for a simulated 3-D face model. Activation levels of the motor state correspond to muscle activations in an anatomically derived model. By integrating fast and robust 2-D processing with 3-D models, we obtain a system that is able to quickly track and interpret complex facial motions in real-time.

\*\*\*\*\*

#### Neural Network Ensembles, Cross Validation, and Active Learning

Anders Krogh, Jesper Vedelsby

Learning of continuous valued functions using neural network ensembles (committees) can give improved accuracy, reliable estimation of the generalization error, and active learning. The ambiguity is defined as the variation of the output of ensemble members averaged over unlabeled data, so it quantifies the disagreement among the networks. It is discussed how to use the ambiguity in combination with cross-validation to give a reliable estimate of the ensemble generalization error, and how this type of ensemble cross-validation can sometimes improve performance. It



is shown how to estimate the optimal weights of the ensemble members using unlabeled data. By a generalization of query by committee, it is finally shown how the ambiguity can be used to select new training data to be labeled in an active learning scheme.

\*\*\*\*\*

## Extracting Rules from Artificial Neural Networks with Distributed Representations

Sebastian Thrun

Although artificial neural networks have been applied in a variety of real-world scenarios with remarkable success, they have often been criticized for exhibiting a low degree of human comprehensibility. Techniques that compile compact sets of symbolic rules out of artificial neural networks offer a promising perspective to overcome this obvious deficiency of neural network representations. This paper presents an approach to the extraction of if-then rules from artificial neural networks. Its key mechanism is validity interval analysis, which is a generic tool for extracting symbolic knowledge by propagating rule-like knowledge through backpropagation-style neural networks. Empirical studies in a robot arm domain illustrate the appropriateness of the proposed method for extracting rules from networks with real-valued and distributed representations.

\*\*\*\*\*

## A model of the hippocampus combining self-organization and associative memory function

Michael Hasselmo, Eric Schnell, Joshua Berke, Edi Barkai

A model of the hippocampus is presented which forms rapid self-organized representations of input arriving via the perforant path, performs recall of previous associations in region CA3, and performs comparison of this recall with afferent input in region CA1. This comparison drives feedback regulation of cholinergic modulation to set appropriate dynamics for learning of new representations in region CA3 and CA1. The network responds to novel patterns with increased cholinergic modulation, allowing storage of new self-organized representations, but responds to familiar patterns with a decrease in acetylcholine, allowing recall based on previous representations. This requires selectivity of the cholinergic suppression of synaptic transmission in stratum radiatum of regions CA3 and CA1, which has been demonstrated experimentally.

\*\*\*\*\*

## Glove-TalkII: Mapping Hand Gestures to Speech Using Neural Networks

Sidney Fels, Geoffrey E. Hinton

Glove-TalkII is a system which translates hand gestures to speech through an adaptive interface. Hand gestures are mapped continuously to 10 control parameters of a parallel formant speech synthesizer. The mapping allows the hand to act as an artificial vocal tract that produces speech in real time. This gives an unlimited vocabulary in addition to direct control of fundamental frequency and volume. Currently, the best version of Glove-TalkII uses several input devices (including a CyberGlove, a ContactGlove, a 3-space tracker, and a foot-pedal), a parallel formant speech synthesizer and 3 neural networks. The gesture-to-speech task is divided into vowel and consonant production by using a gating network to weight the outputs of a vowel and a consonant neural network. The gating network and the consonant network are trained with examples from the user. The vowel network implements a fixed, user-defined relationship between hand-position and vowel sound and does not require any training examples from the user. Volume, fundamental frequency and stop consonants are produced with a fixed mapping from the input devices. One subject has trained to speak intelligibly with Glove-TalkII. He speaks slowly with speech quality similar to a text-to-speech synthesizer but with far more natural-sounding pitch variations.

\*\*\*\*\*

## Learning in large linear perceptrons and why the thermodynamic limit is relevant to the real world

Peter Sollich

We present a new method for obtaining the response function  $\rho$  and its average  $G$  from which most of the properties of learning and generalization in linear perceptrons can be derived. We first rederive the known results for the 'thermodynamic limit' of infinite perceptron size  $N$  and show explicitly that  $\rho$  is self-averaging in this limit. We then discuss extensions of our method to more general learning scenarios with anisotropic teacher space priors, input distributions, and weight decay terms. Finally, we use our method to calculate the finite  $N$  corrections of order  $1/N$  to  $G$  and discuss the corresponding finite size effects on generalization and learning dynamics. An important spin-off is the observation that results obtained in the thermodynamic limit are often directly relevant to systems of fairly modest, 'real-world' sizes.

\*\*\*\*\*

Learning Saccadic Eye Movements Using Multiscale Spatial Filters

Rajesh Rao, Dana Ballard

We describe a framework for learning saccadic eye movements using a photometric representation of target points in natural scenes. The representation takes the form of a high-dimensional vector comprised of the responses of spatial filters at different orientations and scales. We first demonstrate the use of this response vector in the task of locating previously foveated points in a scene and subsequently use this property in a multisaccade strategy to derive an adaptive motor map for delivering accurate saccades.

\*\*\*\*\*

A Charge-Based CMOS Parallel Analog Vector Quantizer

Gert Cauwenberghs, Volnei Pedroni

We present an analog VLSI chip for parallel analog vector quantization. The MOSIS 2.0 J..Lm double-poly CMOS Tiny chip contains an array of  $16 \times 16$  charge-based distance estimation cells, implementing a mean absolute difference (MAD) metric operating on a 16-input analog vector field and 16 analog template vectors. The distance cell including dynamic template storage measures  $60 \times 78$  J..Lm<sup>2</sup>. Additionally, the chip features a winner-take-all (WTA) output circuit of linear complexity, with global positive feedback for fast and decisive settling of a single winner output. Experimental results on the complete  $16 \times 16$  VQ system demonstrate correct operation with 34 dB analog input dynamic range and 3 J..Lsec cycle time at 0.7 mW power dissipation.

\*\*\*\*\*

Boosting the Performance of RBF Networks with Dynamic Decay Adjustment

Michael Berthold, Jay Diamond

Radial Basis Function (RBF) Networks, also known as networks of locally-tuned processing units (see [6]) are well known for their ease of use. Most algorithms used to train these types of networks, however, require a fixed architecture, in which the number of units in the hidden layer must be determined before training starts. The RCE training algorithm, introduced by Reilly, Cooper and Elbaum (see [8]), and its probabilistic extension, the P-RCE algorithm, take advantage of a growing structure in which hidden units are only introduced when necessary. The nature of these algorithms allows training to reach stability much faster than is the case for gradient-descent based methods. Unfortunately P-RCE networks do not adjust the standard deviation of their prototypes individually, using only one global value for this parameter. This paper introduces the Dynamic Decay Adjustment (DDA) algorithm which utilizes the constructive nature of the P-RCE algorithm together with independent adaptation of each prototype's decay factor. In addition, this radial adjustment is class dependent and distinguishes between different neighbours. It is shown that networks trained with the presented algorithm perform substantially better than common RBF networks.

\*\*\*\*\*

An Alternative Model for Mixtures of Experts

Lei Xu, Michael Jordan, Geoffrey E. Hinton

We propose an alternative model for mixtures of experts which uses a different parametric form for the gating network. The modified model is trained by the EM algorithm. In comparison with earlier models-trained by either EM or gradient ascent-there is no need to select a learning stepsize. We report simulation experiments which show that the new architecture yields faster convergence. We also apply the new model to two problem domains: piecewise nonlinear function approximation and the combination of multiple previously trained classifiers.

\*\*\*\*\*

Catastrophic Interference in Human Motor Learning

Tom Brashers-Krug, Reza Shadmehr, Emanuel Todorov

Biological sensorimotor systems are not static maps that transform input (sensory information) into output (motor behavior). Evidence from many lines of research suggests that their representations are plastic, experience-dependent entities. While this plasticity is essential for flexible behavior, it presents the nervous system with difficult organizational challenges. If the sensorimotor system adapts itself to perform well under one set of circumstances, will it then perform poorly when placed in an environment with different demands (negative transfer)? Will a later experience-dependent change undo the benefits of previous learning (catastrophic interference)? We explore the first question in a separate paper in this volume (Shadmehr et al. 1995). Here we present psychophysical and computational results that explore the question of catastrophic interference in the context of a dynamic motor learning task. Under some conditions, subjects show evidence of catastrophic interference. Under other conditions, however, subjects appear to be immune to its effects. These results suggest that motor learning can undergo a process of consolidation. Modular neural networks are well suited for the demands of learning multiple input/output mappings. By incorporating the notion of fast- and slow-changing connections into a modular architecture, we were able to account for the psychophysical results.

\*\*\*\*\*

On the Computational Complexity of Networks of Spiking Neurons

Wolfgang Maass

We investigate the computational power of a formal model for networks of spiking neurons, both for the assumption of an unlimited timing precision, and for the case of a limited timing precision. We also prove upper and lower bounds for the number of examples that are needed to train such networks.

\*\*\*\*\*

New Algorithms for 2D and 3D Point Matching: Pose Estimation and Correspondence

Steven Gold, Chien-Ping Lu, Anand Rangarajan, Suguna Pappu, Eric Mjølhus

A fundamental open problem in computer vision-determining pose and correspondence between two sets of points in space-is solved with a novel, robust and easily implementable algorithm. The technique works on noisy point sets that may be of unequal sizes and may differ by non-rigid transformations. A 2D variation calculates the pose between two point sets related by an affine transformation-translation, rotation, scale and shear. A 3D to 3D variation calculates translation and rotation. An objective describing the problem is derived from Mean field theory. The objective is minimized with clocked (EM-like) dynamics. Experiments with both handwritten and synthetic data provide empirical evidence for the method.

\*\*\*\*\*

PCA-Pyramids for Image Compression

Horst Bischof, Kurt Hornik

This paper presents a new method for image compression by neural networks. First, we show that we can use neural networks in a pyramidal framework, yielding the so-called PCA pyramids. Then we present an image compression method based on the PCA pyramid, which is similar to the Laplace pyramid

d and wavelet transform. Some experimental results with real images are reported. Finally, we present a method to combine the quantization step with the learning of the PCA pyramid.

\*\*\*\*\*

#### Morphogenesis of the Lateral Geniculate Nucleus: How Singularities Affect Global Structure

Svilen Tzonev, Klaus Schulten, Joseph Malpeli

The macaque lateral geniculate nucleus (LGN) exhibits an intricate lamination pattern, which changes midway through the nucleus at a point coincident with small gaps due to the blind spot in the retina. We present a three-dimensional model of morphogenesis in which local cell interactions cause a wave of development of neuronal receptive fields to propagate through the nucleus and establish two distinct lamination patterns. We examine the interactions between the wave and the localized singularities due to the gaps, and find that the gaps induce the change in lamination pattern. We explore critical factors which determine general LGN organization.

\*\*\*\*\*

#### Reinforcement Learning Predicts the Site of Plasticity for Auditory Remapping in the Barn Owl

Alexandre Pouget, Cedric Deffayet, Terrence J. Sejnowski

The auditory system of the barn owl contains several spatial maps. In young barn owls raised with optical prisms over their eyes, these auditory maps are shifted to stay in register with the visual map, suggesting that the visual input imposes a frame of reference on the auditory maps. However, the optic tectum, the first site of convergence of visual with auditory information, is not the site of plasticity for the shift of the auditory maps; the plasticity occurs instead in the inferior colliculus, which contains an auditory map and projects into the optic tectum. We explored a model of the owl remapping in which a global reinforcement signal whose delivery is controlled by visual foveation. A Hebb learning rule gated by reinforcement learned to appropriately adjust auditory maps. In addition, reinforcement learning preferentially adjusted the weights in the inferior colliculus, as in the owl brain, even though the weights were allowed to change throughout the auditory system. This observation raises the possibility that the site of learning does not have to be genetically specified, but could be determined by how the learning procedure interacts with the network architecture.

\*\*\*\*\*

#### A Lagrangian Formulation For Optical Backpropagation Training In Kerr-Type Optical Networks

James Steck, Steven Skinner, Alvaro Cruz-Cabrera, Elizabeth Behrman

A training method based on a form of continuous spatially distributed optical error back-propagation is presented for an all optical network composed of nondiscrete neurons and weighted interconnections. The all optical network is feed-forward and is composed of thin layers of a Kerr type self focusing/defocusing nonlinear optical material. The training method is derived from a Lagrangian formulation of the constrained minimization of the network error at the output. This leads to a formulation that describes training as a calculation of the distributed error of the optical signal at the output which is then reflected back through the device to assign a spatially distributed error to the internal layers. This error is then used to modify the internal weighting values. Results from several computer simulations of the training are presented, and a simple optical table demonstration of the network is discussed.

\*\*\*\*\*

#### Instance-Based State Identification for Reinforcement Learning

R. Andrew McCallum

This paper presents instance-based state identification, an approach to reinforcement learning and hidden state that builds disambiguating amounts of short-term memory on-line, and also learns with an order of magnitude fewer training steps than several previous approaches. Inspired by a key simil

arity between learning with hidden state and learning in continuous geometrical spaces, this approach uses instance-based (or "memory-based") learning, a method that has worked well in continuous spaces.

\*\*\*\*\*

#### Nonlinear Image Interpolation using Manifold Learning

Christoph Bregler, Stephen Omohundro

The problem of interpolating between specified images in an image sequence is a simple, but important task in model-based vision. We describe an approach based on the abstract task of "manifold learning" and present results on both synthetic and real image sequences. This problem arose in the development of a combined lip-reading and speech recognition system.

\*\*\*\*\*

#### A Growing Neural Gas Network Learns Topologies

Bernd Fritzke

An incremental network model is introduced which is able to learn the important topological relations in a given set of input vectors by means of a simple Hebb-like learning rule. In contrast to previous approaches like the "neural gas" method of Martinetz and Schulten (1991, 1994), this model has no parameters which change over time and is able to continue learning, adding units and connections, until a performance criterion has been met. Applications of the model include vector quantization, clustering, and interpolation.

\*\*\*\*\*

#### Transformation Invariant Autoassociation with Application to Handwritten Character Recognition

Holger Schwenk, Maurice Milgram

When training neural networks by the classical backpropagation algorithm the whole problem to learn must be expressed by a set of inputs and desired outputs. However, we often have high-level knowledge about the learning problem. In optical character recognition (OCR), for instance, we know that the classification should be invariant under a set of transformations like rotation or translation. We propose a new modular classification system based on several autoassociative multilayer perceptrons which allows the efficient incorporation of such knowledge. Results are reported on the NIST database of upper case handwritten letters and compared to other approaches to the invariance problem.

\*\*\*\*\*

#### Learning to Play the Game of Chess

Sebastian Thrun

This paper presents NeuroChess, a program which learns to play chess from the final outcome of games. NeuroChess learns chess board evaluation functions, represented by artificial neural networks. It integrates inductive neural network learning, temporal differencing, and a variant of explanation-based learning. Performance results illustrate some of the strengths and weaknesses of this approach.

\*\*\*\*\*

#### Interior Point Implementations of Alternating Minimization Training

Michael Lemmon, Peter Szymanski

This paper presents an alternating minimization (AM) algorithm used in the training of radial basis function and linear regressor networks. The algorithm is a modification of a small-step interior point method used in solving primal linear programs. The algorithm has a convergence rate of  $O(n \log L)$  iterations where  $n$  is a measure of the network size and  $L$  is a measure of the resulting solution's accuracy. Two results are presented that specify how aggressively the two steps of the AM may be pursued to ensure convergence of each step of the alternating minimization.

\*\*\*\*\*

#### A Convolutional Neural Network Hand Tracker

Steven Nowlan, John Platt

We describe a system that can track a hand in a sequence of video frames and

d recognize hand gestures in a user-independent manner. The system locates the hand in each video frame and determines if the hand is open or closed. The tracking system is able to track the hand to within  $\pm 10$  pixels of its correct location in 99.7% of the frames from a test set containing video sequences from 18 different individuals captured in 18 different room environments. The gesture recognition network correctly determines if the hand being tracked is open or closed in 99.1 % of the frames in this test set. The system has been designed to operate in real time with existing hardware.

\*\*\*\*\*

#### Temporal Dynamics of Generalization in Neural Networks

Changfeng Wang, Santosh Venkatesh

This paper presents a rigorous characterization of how a general nonlinear learning machine generalizes during the training process when it is trained on a random sample using a gradient descent algorithm based on reduction of training error. It is shown, in particular, that best generalization performance occurs, in general, before the global minimum of the training error is achieved. The different roles played by the complexity of the machine class and the complexity of the specific machine in the class during learning are also precisely demarcated.

\*\*\*\*\*

#### Learning Stochastic Perceptrons Under k-Blocking Distributions

Mario Marchand, Saeed Haddjifaradji

We present a statistical method that PAC learns the class of stochastic perceptrons with arbitrary monotonic activation function and weights  $W_i \in \{-1, 0, +1\}$  when the probability distribution that generates the input examples is member of a family that we call k-blocking distributions. Such distributions represent an important step beyond the case where each input variable is statistically independent since the 2k-blocking family contains all the Markov distributions of order k. By stochastic perceptron we mean a perceptron which, upon presentation of input vector  $x$ , outputs 1 with probability  $f(\sum W_i x_i - B)$ . Because the same algorithm works for any monotonic (nondecreasing or nonincreasing) activation function on Boolean domain, it handles the well studied cases of sigmoids and the "usual" radial basis functions.

\*\*\*\*\*

#### Recurrent Networks: Second Order Properties and Pruning

Morten Pedersen, Lars Hansen

Second order properties of cost functions for recurrent networks are investigated. We analyze a layered fully recurrent architecture, the virtue of this architecture is that it features the conventional feedforward architecture as a special case. A detailed description of recursive computation of the full Hessian of the network cost function is provided. We discuss the possibility of invoking simplifying approximations of the Hessian and show how weight decays from the cost function and thereby greatly assist training. We present tentative pruning results, using Hassibi et al.'s Optimal Brain Surgeon, demonstrating that recurrent networks can construct an efficient internal memory.

\*\*\*\*\*

#### Unsupervised Classification of 3D Objects from 2D Views

Satoshi Suzuki, Hiroshi Ando

This paper presents an unsupervised learning scheme for categorizing 3D objects from their 2D projected images. The scheme exploits an auto-associative network's ability to encode each view of a single object into a representation that indicates its view direction. We propose two models that employ different classification mechanisms; the first model selects an auto-associative network whose recovered view best matches the input view, and the second model is based on a modular architecture whose additional network classifies the views by splitting the input space nonlinearly. We demonstrate the effectiveness of the proposed classification models through simulations using 3D wire

-frame objects.

\*\*\*\*\*

#### Hyperparameters Evidence and Generalisation for an Unrealisable Rule

Glenn Marion, David Saad

Using a statistical mechanical formalism we calculate the evidence, generalisation error and consistency measure for a linear perceptron trained and tested on a set of examples generated by a non linear teacher. The teacher is said to be unrealisable because the student can never model it without error. Our model allows us to interpolate between the known case of a linear teacher, and an unrealisable, nonlinear teacher. A comparison of the hyperparameters which maximise the evidence with those that optimise the performance measures reveals that, in the non-linear case, the evidence procedure is a misleading guide to optimising performance. Finally, we explore the extent to which the evidence procedure is unreliable and find that, despite being sub-optimal, in some circumstances it might be a useful method for fixing the hyperparameters.

\*\*\*\*\*

#### Adaptive Elastic Input Field for Recognition Improvement

Minoru Asogawa

For machines to perform classification tasks, such as speech and character recognition, appropriately handling deformed patterns is a key to achieving high performance. The authors presents a new type of classification system, an Adaptive Elastic Input Field Neural Network (AIFNN), which includes a simple pre-trained neural network and an elastic input field attached to an input layer. By using an iterative method, AIFNN can determine an optimal affine translation for an elastic input field to compensate for the original deformations. The convergence of the AIFNN algorithm is shown. AIFNN is applied for handwritten numerals recognition. Consequently, 10.83% of originally misclassified patterns are correctly categorized and total performance is improved, without modifying the neural network.

\*\*\*\*\*

#### A Model for Chemosensory Reception

Rainer Malaka, Thomas Ragg, Martin Hammer

A new model for chemosensory reception is presented. It models reactions between odor molecules and receptor proteins and the activation of second messenger by receptor proteins. The mathematical formulation of the reaction kinetics is transformed into an artificial neural network (ANN). The resulting feed-forward network provides a powerful means for parameter fitting by applying learning algorithms. The weights of the network corresponding to chemical parameters can be trained by presenting experimental data. We demonstrate the simulation capabilities of the model with experimental data from honey bee chemosensory neurons. It can be shown that our model is sufficient to rebuild the observed data and that simpler models are not able to do this task.

\*\*\*\*\*

#### A Real Time Clustering CMOS Neural Engine

Teresa Serrano-Gotarredona, Bernabé Linares-Barranco, José Huertas

We describe an analog VLSI implementation of the ARTI algorithm (Carpenter, 1987). A prototype chip has been fabricated in a standard low cost 1.5- $\mu$ m double-metal single-poly CMOS process. It has a die area of 1cm<sup>2</sup> and is mounted in a 120-pins PGA package. The chip realizes a modified version of the original ARTI architecture. Such modification has been shown to preserve all computational properties of the original algorithm (Serrano, 1994a), while being more appropriate for VLSI realizations. The chip implements an ARTI network with 100 F1 nodes and 18 F2 nodes. It can therefore cluster 100 binary pixels input patterns into up to 18 different categories. Modular expansibility of the system is possible by assembling an NxM array of chips without any extra interfacing circuitry, resulting in an F1 layer with 100xN nodes, and an F2 layer with 18xM nodes. Pattern classification is performed in less than 1.8 $\mu$ s, which means an equivalent computing power of 2.2x10<sup>9</sup>

connections and connection-updates per second. Although internally the chip is analog in nature, it interfaces to the outside world through digital signals, thus having a true asynchronous digital behavior. Experimental chip test results are available, which have been obtained through test equipments for digital chips.

\*\*\*\*\*

Learning from queries for maximum information gain in imperfectly learnable problems

Peter Sollich, David Saad

In supervised learning, learning from queries rather than from random examples can improve generalization performance significantly. We study the performance of query learning for problems where the student cannot learn the teacher perfectly, which occur frequently in practice. As a prototypical scenario of this kind, we consider a linear perceptron student learning a binary perceptron teacher. Two kinds of queries for maximum information gain, i.e., minimum entropy, are investigated: Minimum student space entropy (MSSE) queries, which are appropriate if the teacher space is unknown, and minimum teacher space entropy (MTSE) queries, which can be used if the teacher space is assumed to be known, but a student of a simpler form has deliberately been chosen. We find that for MSSE queries, the structure of the student space determines the efficacy of query learning, whereas MTSE queries lead to a higher generalization error than random examples, due to a lack of feedback about the progress of the student in the way queries are selected.

\*\*\*\*\*

Factorial Learning by Clustering Features

Joshua Tenenbaum, Emanuel V. Todorov

We introduce a novel algorithm for factorial learning, motivated by segmentation problems in computational vision, in which the underlying factors correspond to clusters of highly correlated input features. The algorithm derives from a new kind of competitive clustering model, in which the cluster generators compete to explain each feature of the data set and cooperate to explain each input example, rather than competing for examples and cooperating on features, as in traditional clustering algorithms. A natural extension of the algorithm recovers hierarchical models of data generated from multiple unknown categories, each with a different, multiple causal structure. Several simulations demonstrate the power of this approach.

\*\*\*\*\*

Generalization in Reinforcement Learning: Safely Approximating the Value Function

Justin Boyan, Andrew Moore

A straightforward approach to the curse of dimensionality in reinforcement learning and dynamic programming is to replace the lookup table with a generalizing function approximator such as a neural net. Although this has been successful in the domain of backgammon, there is no guarantee of convergence. In this paper, we show that the combination of dynamic programming and function approximation is not robust, and in even very benign cases, may produce an entirely wrong policy. We then introduce Grow-Support, a new algorithm which is safe from divergence yet can still reap the benefits of successful generalization.

\*\*\*\*\*

A Rapid Graph-based Method for Arbitrary Transformation-Invariant Pattern Classification

Alessandro Sperduti, David Stork

We present a graph-based method for rapid, accurate search through prototypes for transformation-invariant pattern classification. Our method has in theory the same recognition accuracy as other recent methods based on "tangent distance" [Simard et al., 1994], since it uses the same categorization rule. Nevertheless ours is significantly faster during classification because far fewer tangent distances need be computed. C



crucial to the success of our system are 1) a novel graph architecture in which transformation constraints and geometric relationships among prototypes are encoded during learning, and 2) an improved graph search criterion, used during classification. These architectural insights are applicable to a wide range of problem domains. Here we demonstrate that on a handwriting recognition task, a basic implementation of our system requires less than half the computation of the Euclidean sorting method.

\*\*\*\*\*

A solvable connectionist model of immediate recall of ordered lists

Neil Burgess

A model of short-term memory for serially ordered lists of verbal stimuli is proposed as an implementation of the 'articulatory loop' thought to mediate this type of memory (Baddeley, 1986). The model predicts the presence of a repeatable time-varying 'context' signal coding the timing of items' presentation in addition to a store of phonological information and a process of serial rehearsal. Items are associated with context nodes and phonemes by Hebbian connections showing both short and long term plasticity. Items are activated by phonemic input during presentation and reactivated by context and phonemic feedback during output. Serial selection of items occurs via a winner-take-all interaction amongst items, with the winner subsequently receiving decaying inhibition. An approximate analysis of error probabilities due to Gaussian noise during output is presented. The model provides an explanatory account of the probability of error as a function of serial position, list length, word length, phonemic similarity, temporal grouping, item and list familiarity, and is proposed as the starting point for a model of rehearsal and vocabulary acquisition.

\*\*\*\*\*

$H^\infty$  Optimal Training Algorithms and their Relation to Backpropagation

Babak Hassibi, Thomas Kailath

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.

\*\*\*\*\*

Using a neural net to instantiate a deformable model

Christopher Williams, Michael Revow, Geoffrey E. Hinton

Deformable models are an attractive approach to recognizing non-rigid objects which have considerable within class variability. However, there are severe search problems associated with fitting the models to data. We show that by using neural networks to provide better starting points, the search time can be significantly reduced. The method is demonstrated on a character recognition task.

\*\*\*\*\*

Grammar Learning by a Self-Organizing Network

Michiro Negishi

This paper presents the design and simulation results of a self-organizing neural network which induces a grammar from example sentences. Input sentences are generated from a simple phrase structure grammar including number agreement, verb transitivity, and recursive noun phrase construction rules. The network induces a grammar explicitly in the form of symbol categorization rules and phrase structure rules.

\*\*\*\*\*

Coarse-to-Fine Image Search Using Neural Networks

Clay Spence, John Pearson, Jim Bergen

The efficiency of image search can be greatly improved by using a coarse-to-fine search strategy with a multi-resolution image representation. However, if the resolution is so low that the objects have few distinguishing features, search becomes difficult. We show that the performance of search at such low resolutions can be improved by using context information, i.e., objects visible at low-resolution which are not the objects of

interest but are associated with them. The networks can be given explicit context information as inputs, or they can learn to detect the context objects, in which case the user does not have to be aware of their existence. We also use Integrated Feature Pyramids, which represent high-frequency information at low resolutions. The use of multi-resolution search techniques allows us to combine information about the appearance of the objects on many scales in an efficient way. A natural form of exemplar selection also arises from these techniques. We illustrate these ideas by training hierarchical systems of neural networks to find clusters of buildings in aerial photographs of farmland.

\*\*\*\*\*

ICEG Morphology Classification using an Analogue VLSI Neural Network

Richard Coggins, Marwan Jabri, Barry Flower, Stephen Pickard

An analogue VLSI neural network has been designed and tested to perform cardiac morphology classification tasks. Analogue techniques were chosen to meet the strict power and area requirements of an Implantable Cardioverter Defibrillator (ICD) system. The robustness of the neural network architecture reduces the impact of noise, drift and offsets inherent in analogue approaches. The network is a 10:6:3 multi-layer perceptron with on chip digital weight storage, a bucket brigade input to feed the Intracardiac Electrogram (ICEG) to the network and has a winner take all circuit at the output. The network was trained in loop and included a commercial ICD in the signal processing path. The system has successfully distinguished arrhythmia for different patients with better than 90% true positive and true negative detections for dangerous rhythms which cannot be detected by present ICDs. The chip was implemented in 1.2um CMOS and consumes less than 200mW maximum average power in an area of 2.2 x 2.2mm<sup>2</sup>.

\*\*\*\*\*