

Regular and Irregular Gallager-type Error-Correcting Codes

Yoshiyuki Kabashima, Tatsuto Murayama, David Saad, Renato Vicente

The performance of regular and irregular Gallager-type error-correcting code is investigated via methods of statistical physics. The transmitted codeword comprises products of the original message bits selected by two randomly-constructed sparse matrices; the number of non-zero row/column elements in these matrices constitutes a family of codes.

We show that Shannon's channel capacity may be saturated in equilibrium for many of the regular codes while slightly lower performance is obtained for others which may be of higher practical relevance. Decoding aspects are considered by employing the TAP approach which is identical to the commonly used belief-propagation-based decoding. We show that irregular codes may saturate Shannon's capacity but with improved dynamical properties.

An MEG Study of Response Latency and Variability in the Human Visual System During a Visual-Motor Integration Task

Akaysha Tang, Barak Pearlmutter, Tim Hely, Michael Zibulevsky, Michael Weisend

Human reaction times during sensory-motor tasks vary considerably. To begin to understand how this variability arises, we examined neuronal population response time variability at early versus late visual processing stages. The conventional view is that precise temporal information is gradually lost as information is passed through a layered network of mean-rate "units." We tested humans whether neuronal populations at different processing stages behave like mean-rate "units".

A blind source separation algorithm was applied to MEG signals from sensory-motor integration tasks. Response time latency and variability for multiple visual sources were estimated by detecting single-trial stimulus-locked events for each source. In two subjects tested on four visual reaction time tasks, we reliably identified sources belonging to early and late visual processing stages. The standard deviation of response latency was smaller for early rather than late processing stages. This supports the hypothesis that human population response time variability increases from early to late visual processing stages.

The Entropy Regularization Information Criterion

Alex Smola, John Shawe-Taylor, Bernhard Schölkopf, Robert C. Williamson

Effective methods of capacity control via uniform convergence bounds for function expansions have been largely limited to Support Vector machines, where good bounds are obtainable by the entropy number approach. We extend these methods to systems with expansions in terms of arbitrary (parameterized) basis functions and a wide range of regularization methods covering the whole range of general linear additive models. This is achieved by a data dependent analysis of the eigenvalues of the corresponding design matrix.

Invariant Feature Extraction and Classification in Kernel Spaces

Sebastian Mika, Gunnar Rätsch, Jason Weston, Bernhard Schölkopf, Alex Smola, Klaus-Robert Müller

In hyperspectral imagery one pixel typically consists of a mixture of the reflectance spectra of several materials, where the mixture coefficients correspond to the abundances of the constituting materials. We assume linear combinations of reflectance spectra with some additive sensor noise and derive a probabilistic MAP framework for analyzing hyperspectral data. As the material reflectance characteristics are not known a priori, we face the problem of unsupervised linear unmixing. The incorporation of different prior information (e.g. positivity and normalization of the abundances) naturally leads to a family of interesting algorithms, for example in the noise-free case yielding an algorithm that can be understood as constrained independent component analysis (ICA).

). Simulations underline the usefulness of our theory.

Correctness of Belief Propagation in Gaussian Graphical Models of Arbitrary Topology

Yair Weiss, William Freeman

Local "belief propagation" rules of the sort proposed by Pearl [15] are guaranteed to converge to the correct posterior probabilities in singly connected graphical models. Recently, a number of researchers have empirically demonstrated good performance of "loopy belief propagation" using these same rules on graphs with loops. Perhaps the most dramatic instance is the near Shannon-limit performance of "Turbo codes", whose decoding algorithm is equivalent to loopy belief propagation. Except for the case of graphs with a single loop, there has been little theoretical understanding of the performance of loopy propagation. Here we analyze belief propagation in networks with arbitrary topologies when the nodes in the graph describe jointly Gaussian random variables. We give an analytical formula relating the true posterior probabilities with those calculated using loopy propagation. We give sufficient conditions for convergence and show that when belief propagation converges it gives the correct posterior means for all graph topologies, not just networks with a single loop. The related "max-product" belief propagation algorithm finds the maximum posterior probability estimate for singly connected networks. We show that, even for non-Gaussian probability distributions, the convergence points of the max-product algorithm in loopy networks are maxima over a particular large local neighborhood of the posterior probability. These results help clarify the empirical performance results and motivate using the powerful belief propagation algorithm in a broader class of networks.

Population Decoding Based on an Unfaithful Model

Si Wu, Hiroyuki Nakahara, Noboru Murata, Shun-ichi Amari

We study a population decoding paradigm in which the maximum likelihood inference is based on an unfaithful decoding model (UMLI). This is usually the case for neural population decoding because the encoding process of the brain is not exactly known, or because a simplified decoding model is preferred for saving computational cost. We consider an unfaithful decoding model which neglects the pair-wise correlation between neuronal activities, and prove that UMLI is asymptotically efficient when the neuronal correlation is uniform or of limited-range. The performance of UMLI is compared with that of the maximum likelihood inference based on a faithful model and that of the center of mass decoding method. It turns out that UMLI has advantages of decreasing the computational complexity remarkably and maintaining a high-level decoding accuracy at the same time. The effect of correlation on the decoding accuracy is also discussed.

Broadband Direction-Of-Arrival Estimation Based on Second Order Statistics

Justinian Rosca, Joseph Ruanaidh, Alexander Jourjine, Scott Rickard

N wideband sources recorded using N closely spaced receivers can feasibly be separated based only on second order statistics when using a physical model of the mixing process. In this case we show that the parameter estimation problem can be essentially reduced to considering directions of arrival and attenuations of each signal. The paper presents two demixing methods operating in the time and frequency domain and experimentally shows that it is always possible to demix signals arriving at different angles. Moreover, one can use spatial cues to solve the channel selection problem and a post-processing Wiener filter to ameliorate the artifacts caused by demixing.

Emergence of Topography and Complex Cell Properties from Natural Images using Extensions of ICA

Aapo Hyvärinen, Patrik Hoyer

Independent component analysis of natural images leads to emergence

of simple cell properties, i.e. linear filters that resemble wavelets or Gabor functions. In this paper, we extend ICA to explain further properties of VI cells. First, we decompose natural images into independent subspaces instead of scalar components. This model leads to emergence of phase and shift invariant features, similar to those in VI complex cells. Second, we define a topography between the linear components obtained by ICA. The topographic distance between two components is defined by their higher-order correlations, so that two components are close to each other in the topography if they are strongly dependent on each other. This leads to simultaneous emergence of both topography and invariances similar to complex cell properties.

Reinforcement Learning Using Approximate Belief States

Andres Rodriguez, Ronald Parr, Daphne Koller

The problem of developing good policies for partially observable Markov decision problems (POMDPs) remains one of the most challenging areas of research in stochastic planning. One line of research in this area involves the use of reinforcement learning with belief states, probability distributions over the underlying model states. This is a promising method for small problems, but its application is limited by the intractability of computing or representing a full belief state for large problems. Recent work shows that, in many settings, we can maintain an approximate belief state, which is fairly close to the true belief state. In particular, great success has been shown with approximate belief states that marginalize out correlations between state variables. In this paper, we investigate two methods of full belief state reinforcement learning and one novel method for reinforcement learning using factored approximate belief states. We compare the performance of these algorithms on several well-known problem from the literature. Our results demonstrate the importance of approximate belief state representations for large problems.

Building Predictive Models from Fractal Representations of Symbolic Sequences

Peter Tiffo, Georg Dorffner

We propose a novel approach for building finite memory predictive models similar in spirit to variable memory length Markov models (VLMs). The models are constructed by first transforming the n-block structure of the training sequence into a spatial structure of points in a unit hypercube, such that the longer is the common suffix shared by any two n-blocks, the closer lie their point representations. Such a transformation embodies a Markov assumption - n-blocks with long common suffixes are likely to produce similar continuations. Finding a set of prediction contexts is formulated as a resource allocation problem solved by vector quantizing the spatial n-block representation. We compare our model with both the classical and variable memory length Markov models on three data sets with different memory and stochastic components. Our models have a superior performance, yet, their construction is fully automatic, which is shown to be problematic in the case of VLMs.

Neural Computation with Winner-Take-All as the Only Nonlinear Operation

Wolfgang Maass

Everybody "knows" that neural networks need more than a single layer of nonlinear units to compute interesting functions. We show that this is false if one employs winner-take-all as nonlinear unit:

Support Vector Method for Multivariate Density Estimation

Vladimir Vapnik, Sayan Mukherjee

A new method for multivariate density estimation is developed based on the Support Vector Method (SVM) solution of inverse ill-posed problems. The solution has the form of a mixture of densities. This method with Gaussian kernels compared favorably to both Parzen's method and the Gaussian Mixture Model method. For synthetic data we achieve more accurate

rate estimates for densities of 2, 6, 12, and 40 dimensions.

Leveraged Vector Machines

Yoram Singer

We describe an iterative algorithm for building vector machines used in classification tasks. The algorithm builds on ideas from support vector machines, boosting, and generalized additive models. The algorithm can be used with various continuously differentiable functions that bound the discrete (0-1) classification loss and is very simple to implement. We test the proposed algorithm with two different loss functions on synthetic and natural data. We also describe a norm-penalized version of the algorithm for the exponential loss function used in AdaBoost. The performance of the algorithm on natural data is comparable to support vector machines while typically its running time is shorter than of SVM.

Learning Factored Representations for Partially Observable Markov Decision Processes

Brian Sallans

The problem of reinforcement learning in a non-Markov environment is explored using a dynamic Bayesian network, where conditional independence assumptions between random variables are compactly represented by network parameters.

The parameters are learned on-line, and approximations are used to perform inference and to compute the optimal value function. The relative effects of inference and value function approximations on the quality of the final policy are investigated, by learning to solve a moderately difficult driving task. The two value function approximations, linear and quadratic, were found to perform similarly, but the quadratic model was more sensitive to initialization. Both performed below the level of human performance on the task. The dynamic Bayesian network performed comparably to a model using a localist hidden state representation, while requiring exponentially fewer parameters.

Variational Inference for Bayesian Mixtures of Factor Analysers

Zoubin Ghahramani, Matthew Beal

We present an algorithm that infers the model structure of a mixture of factor analysers using an efficient and deterministic variational approximation to full Bayesian integration over model parameters. This procedure can automatically determine the optimal number of components and the local dimensionality of each component (i.e. the number of factors in each factor analyser). Alternatively it can be used to infer posterior distributions over number of components and dimensionalities. Since all parameters are integrated out the method is not prone to overfitting. Using a stochastic procedure for adding components it is possible to perform the variational optimisation incrementally and to avoid local maxima. Results show that the method works very well in practice and correctly infers the number and dimensionality of nontrivial synthetic examples. By importance sampling from the variational approximation we show how to obtain unbiased estimates of the true evidence, the exact predictive density, and the KL divergence between the variational posterior and the true posterior, not only in this model but for variational approximations in general.

Topographic Transformation as a Discrete Latent Variable

Nebojsa Jojic, Brendan J. Frey

Invariance to topographic transformations such as translation and shearing in an image has been successfully incorporated into feedforward mechanisms, e.g., "convolutional neural networks", "tangent propagation".

We describe a way to add transformation invariance to a generative density model by approximating the nonlinear transformation manifold by a discrete set of transformations. An EM algorithm for the original model can

be extended to the new model by computing expectations over the set of transformations. We show how to add a discrete transformation variable to Gaussian mixture modeling, factor analysis and mixtures of factor analysis. We give results on filtering microscopy images, face and facial pose clustering, and handwritten digit modeling and recognition.

Channel Noise in Excitable Neural Membranes

Amit Manwani, Peter Steinmetz, Christof Koch

Stochastic fluctuations of voltage-gated ion channels generate current and voltage noise in neuronal membranes. This noise may be a critical determinant of the efficacy of information processing within neural systems. Using Monte-Carlo simulations, we carry out a systematic investigation of the relationship between channel kinetics and the resulting membrane voltage noise using a stochastic Markov version of the Mainen-Sejnowski model of dendritic excitability in cortical neurons. Our simulations show that kinetic parameters which lead to an increase in membrane excitability (increasing channel densities, decreasing temperature) also lead to an increase in the magnitude of the sub-threshold voltage noise. Noise also increases as the membrane is depolarized from rest towards threshold. This suggests that channel fluctuations may interfere with a neuron's ability to function as an integrator of its synaptic inputs and may limit the reliability and precision of neural information processing.

Efficient Approaches to Gaussian Process Classification

Lehel Csató, Ernest Fokoué, Manfred Opper, Bernhard Schottky, Ole Winther

We present three simple approximations for the calculation of the posterior mean in Gaussian Process classification. The first two methods are related to mean field ideas known in Statistical Physics. The third approach is based on Bayesian online approach which was motivated by recent results in the Statistical Mechanics of Neural Networks. We present simulation results showing: 1. that the mean field Bayesian evidence may be used for hyperparameter tuning and 2. that the online approach may achieve a low training error fast.

Optimal Sizes of Dendritic and Axonal Arbors

Dmitri Chklovskii

I consider a topographic projection between two neuronal layers with different densities of neurons. Given the number of output neurons connected to each input neuron (divergence or fan-out) and the number of input neurons synapsing on each output neuron (convergence or fan-in) I determine the widths of axonal and dendritic arbors which minimize the total volume of axons and dendrites. My analytical results can be summarized qualitatively in the following rule: neurons of the sparser layer should have arbors wider than those of the denser layer. This agrees with the anatomical data from retinal and cerebellar neurons whose morphology and connectivity are known. The rule may be used to infer connectivity of neurons from their morphology.

v-Arc: Ensemble Learning in the Presence of Outliers

Gunnar Rätsch, Bernhard Schölkopf, Alex Smola, Klaus-Robert Müller, Takashi Onoda, Sebastian Mika

AdaBoost and other ensemble methods have successfully been applied to a number of classification tasks, seemingly defying problems of overfitting. AdaBoost performs gradient descent in an error function with respect to the margin, asymptotically concentrating on the patterns which are hardest to learn. For very noisy problems, however, this can be disadvantageous. Indeed, theoretical analysis has shown that the margin distribution, as opposed to just the minimal margin, plays a crucial role in understanding this phenomenon. Loosely speaking, some outliers should be tolerated if this has the benefit of substantially increasing

the margin on the remaining points. We propose a new boosting algorithm which allows for the possibility of a pre-specified fraction of points to lie in the margin area Or even on the wrong side of the decision boundary.

Monte Carlo POMDPs

Sebastian Thrun

We present a Monte Carlo algorithm for learning to act in partially observable Markov decision processes (POMDPs) with real-valued state and action spaces. Our approach uses importance sampling for representing beliefs, and Monte Carlo approximation for belief propagation. A reinforcement learning algorithm, value iteration, is employed to learn value functions over belief states. Finally, a sample-based version of nearest neighbor is used to generalize across states. Initial empirical results suggest that our approach works well in practical applications.

A Recurrent Model of the Interaction Between Prefrontal and Inferotemporal Cortex in Delay Tasks

Alfonso Renart, Néstor Parga, Edmund Rolls

A very simple model of two reciprocally connected attractor neural networks works is studied analytically in situations similar to those encountered in delay match-to-sample tasks with intervening stimuli and in tasks of memory guided attention. The model qualitatively reproduces many of the experimental data on these types of tasks and provides a framework for the understanding of the experimental observations in the context of the attractor neural network scenario.

Information Factorization in Connectionist Models of Perception

Javier Movellan, James McClelland

We examine a psychophysical law that describes the influence of stimulus and context on perception. According to this law choice probability ratios factorize into components independently controlled by stimulus and context. It has been argued that this pattern of results is incompatible with feedback models of perception. In this paper we examine this claim using neural network models defined via stochastic differential equations. We show that the law is related to a condition named channel separability and has little to do with the existence of feedback connections. In essence, channels are separable if they converge into the response units without direct lateral connections to other channels and if their sensors are not directly contaminated by external inputs to the other channels. Implications of the analysis for cognitive and computational neuroscience are discussed.

Hierarchical Image Probability (HIP) Models

Clay Spence, Lucas Parra

We formulate a model for probability distributions on image spaces. We show that any distribution of images can be factored exactly into conditional distributions of feature vectors at one resolution (pyramid level) conditioned on the image information at lower resolutions. We would like to factor this over positions in the pyramid levels to make it tractable, but such factoring may miss long-range dependencies. To fix this, we introduce hidden class labels at each pixel in the pyramid. The result is a hierarchical mixture of conditional probabilities, similar to a hidden Markov model on a tree. The model parameters can be found with maximum likelihood estimation using the EM algorithm. We have obtained encouraging preliminary results on the problems of detecting various objects in SAR images and target recognition in optical aerial images.

Reinforcement Learning for Spoken Dialogue Systems

Satinder Singh, Michael Kearns, Diane Litman, Marilyn Walker

Recently, a number of authors have proposed treating dialogue systems as Markov

decision processes (MDPs). However, the practical application of MDP algorithms to dialogue systems faces a number of severe technical challenges. We have built a general software tool (RLDS, for Reinforcement Learning for Dialogue Systems) based on the MDP framework, and have applied it to dialogue corpora gathered from two dialogue systems built at AT&T Labs. Our experiments demonstrate that RLDS holds promise as a tool for "browsing" and understanding correlations in complex, temporally dependent dialogue corpora.

Distributed Synchrony of Spiking Neurons in a Hebbian Cell Assembly

David Horn, Nir Levy, Isaac Meilijson, Eytan Ruppin

We investigate the behavior of a Hebbian cell assembly of spiking neurons formed via a temporal synaptic learning curve. This learning function is based on recent experimental findings. It includes potentiation for short time delays between pre- and post-synaptic neuronal spiking, and depression for spiking events occurring in the reverse order. The coupling between the dynamics of the synaptic learning and of the neuronal activation leads to interesting results. We find that the cell assembly can fire asynchronously, but may also function in complete synchrony, or in distributed synchrony. The latter implies spontaneous division of the Hebbian cell assembly into groups of cells that fire in a cyclic manner. We investigate the behavior of distributed synchrony both by simulations and by analytic calculations of the resulting synaptic distributions.

Image Representations for Facial Expression Coding

Marian Bartlett, Gianluca Donato, Javier Movellan, Joseph Hager, Paul Ekman, Terrence J. Sejnowski

The Facial Action Coding System (FACS) (9) is an objective method for quantifying facial movement in terms of component actions. This system is widely used in behavioral investigations of emotion, cognitive processes, and social interaction. The coding is presently performed by highly trained human experts. This paper explores and compares techniques for automatically recognizing facial actions in sequences of images. These methods include unsupervised learning techniques for finding basis images such as principal component analysis, independent component analysis and local feature analysis, and supervised learning techniques such as Fisher's linear discriminants. These data-driven bases are compared to Gabor wavelets, in which the basis images are predefined. Best performances were obtained using the Gabor wavelet representation and the independent component representation, both of which achieved 96% accuracy for classifying 12 facial actions. The ICA representation employs 2 orders of magnitude fewer basis images than the Gabor representation and takes 90% less CPU time to compute for new images. The results provide converging support for using local basis images, high spatial frequencies, and statistical independence for classifying facial actions.

Algorithms for Independent Components Analysis and Higher Order Statistics

Daniel Lee, Uri Rokni, Haim Sompolinsky

A latent variable generative model with finite noise is used to describe several different algorithms for Independent Components Analysis (ICA). In particular, the Fixed Point ICA algorithm is shown to be equivalent to the Expectation-Maximization algorithm for maximum likelihood under certain constraints, allowing the conditions for global convergence to be elucidated. The algorithms can also be explained by their generic behavior near a singular point where the size of the optimal generative bases vanishes. An expansion of the likelihood about this singular point indicates the role of higher order correlations in determining the features discovered by ICA. The application and convergence of these algorithms are demonstrated on a simple illustrative example.

A SNoW-Based Face Detector

Ming-Hsuan Yang, Dan Roth, Narendra Ahuja

A novel learning approach for human face detection using a network of linear units is presented. The SNoW learning architecture is a sparse network of linear functions over a pre-defined or incrementally learned feature space and is specifically tailored for learning in the presence of a very large number of features. A wide range of face images in different poses, with different expressions and under different lighting conditions are used as a training set to capture the variations of human faces. Experimental results on commonly used benchmark data sets of a wide range of face images show that the SNoW-based approach outperforms methods that use neural networks, Bayesian methods, support vector machines and others. Furthermore, learning and evaluation using the SNoW-based method are significantly more efficient than with other methods.

A Winner-Take-All Circuit with Controllable Soft Max Property

Shih-Chii Liu

I describe a silicon network consisting of a group of excitatory neurons and a global inhibitory neuron. The output of the inhibitory neuron is normalized with respect to the input strengths. This output models the normalization property of the wide-field direction selective cells in the fly visual system. This normalizing property is also useful in any system where we wish the output signal to code only the strength of the inputs, and not be dependent on the number of inputs. The circuitry in each neuron is equivalent to that in Lazzaro's winner-take-all (WTA) circuit with one additional transistor and a voltage reference. Just as in Lazzaro's circuit, the outputs of the excitatory neurons code the neuron with the largest input. The difference here is that multiple winners can be chosen.

By varying the voltage reference of the neuron, the network can transition between a soft-max behavior and a hard WTA behavior. I show results from a fabricated chip of 20 neurons in a 1.2J.Lm CMOS technology.

Bayesian Model Selection for Support Vector Machines, Gaussian Processes and Other Kernel Classifiers

Matthias Seeger

We present a variational Bayesian method for model selection over families of kernels classifiers like Support Vector machines or Gaussian processes. The algorithm needs no user interaction and is able to adapt a large number of kernel parameters to given data without having to sacrifice training cases for validation. This opens the possibility to use sophisticated families of kernels in situations where the small "standard kernel" classes are clearly inappropriate. We relate the method to other work done on Gaussian processes and clarify the relation between Support Vector machines and certain Gaussian process models.

Kirchoff Law Markov Fields for Analog Circuit Design

Richard Golden

Three contributions to developing an algorithm for assisting engineers in designing analog circuits are provided in this paper. First, a method for representing highly nonlinear and non-continuous analog circuits using Kirchoff current law potential functions within the context of a Markov field is described. Second, a relatively efficient algorithm for optimizing the Markov field objective function is briefly described and the convergence proof is briefly sketched. And third, empirical results illustrating the strengths and limitations of the approach are provided within the context of a JFET transistor design problem. The proposed algorithm generated a set of circuit components for the JFET circuit model that accurately generated the desired characteristic curves.

Can VI Mechanisms Account for Figure-Ground and Medial Axis Effects?

Zhaoping Li

When a visual image consists of a figure against a background, V1 cells are physiologically observed to give higher responses to image regions corresponding to the figure relative to their responses to the background. The medial axis of the figure also induces relatively higher responses compared to responses to other locations in the figure (except for the boundary between the figure and the background). Since the receptive fields of V1 cells are very small compared with the global scale of the figure-ground and medial axis effects, it has been suggested that these effects may be caused by feedback from higher visual areas. I show how these effects can be accounted for by V1 mechanisms when the size of the figure is small or is of a certain scale. They are a manifestation of the processes of pre-attentive segmentation which detect and highlight the boundaries between homogeneous image regions.

Policy Gradient Methods for Reinforcement Learning with Function Approximation

Richard S. Sutton, David McAllester, Satinder Singh, Yishay Mansour

Function approximation is essential to reinforcement learning, but the standard approach of approximating a value function and determining a policy from it has so far proven theoretically intractable. In this paper we explore an alternative approach in which the policy is explicitly represented by its own function approximator, independent of the value function, and is updated according to the gradient of expected reward with respect to the policy parameters. Williams's REINFORCE method and actor-critic methods are examples of this approach. Our main new result is to show that the gradient can be written in a form suitable for estimation from experience aided by an approximate action-value or advantage function. Using this result, we prove for the first time that a version of policy iteration with arbitrary differentiable function approximation is convergent to a locally optimal policy.

Lower Bounds on the Complexity of Approximating Continuous Functions by Sigmoidal Neural Networks

Michael Schmitt

We calculate lower bounds on the size of sigmoidal neural networks that approximate continuous functions. In particular, we show that for the approximation of polynomials the network size has to grow as $O((\log k)^{1/4})$ where k is the degree of the polynomials. This bound is valid for any input dimension, i.e. independently of the number of variables. The result is obtained by introducing a new method employing upper bounds on the Vapnik-Chervonenkis dimension for proving lower bounds on the size of networks that approximate continuous functions.

Evolving Learnable Languages

Bradley Tonkes, Alan Blair, Janet Wiles

Recent theories suggest that language acquisition is assisted by the evolution of languages towards forms that are easily learnable. In this paper, we evolve combinatorial languages which can be learned by a recurrent neural network quickly and from relatively few examples. Additionally, we evolve languages for generalization in different "worlds", and for generalization from specific examples. We find that languages can be evolved to facilitate different forms of impressive generalization for a minimally biased, general purpose learner. The results provide empirical support for the theory that the language itself, as well as the language environment of a learner, plays a substantial role in learning: that there is far more to language acquisition than the language acquisition device.

Large Margin DAGs for Multiclass Classification

John Platt, Nello Cristianini, John Shawe-Taylor

We present a new learning architecture: the Decision Directed Acyclic G

raph (DDAG), which is used to combine many two-class classifiers into a multiclass classifier. For an N -class problem, the DDAG contains $N(N-1)/2$ classifiers, one for each pair of classes. We present a VC analysis of the case when the node classifiers are hyperplanes; the resulting bound on the test error depends on N and on the margin achieved at the nodes, but not on the dimension of the space. This motivates an algorithm, DAGSVM, which operates in a kernel-induced feature space and uses two-class maximal margin hyperplanes at each decision-node of the DDAG. The DAGSVM is substantially faster to train and evaluate than either the standard algorithm or Max Wins, while maintaining comparable accuracy to both of these algorithms.

Approximate Planning in Large POMDPs via Reusable Trajectories

Michael Kearns, Yishay Mansour, Andrew Ng

We consider the problem of reliably choosing a near-best strategy from a restricted class of strategies Π in a partially observable Markov decision process (POMDP). We assume we are given the ability to simulate the POMDP, and study what might be called the sample complexity - that is, the amount of data one must generate in the POMDP in order to choose a good strategy.

We prove upper bounds on the sample complexity showing that, even for infinitely large and arbitrarily complex POMDPs, the amount of data needed can be finite, and depends only linearly on the complexity of the restricted strategy class Π , and exponentially on the horizon time. This latter dependence can be eased in a variety of ways, including the application of gradient and local search algorithms. Our measure of complexity generalizes the classical supervised learning notion of VC dimension to the settings of reinforcement learning and planning.

Maximum Entropy Discrimination

Tommi Jaakkola, Marina Meila, Tony Jebara

We present a general framework for discriminative estimation based on the maximum entropy principle and its extensions. All calculations involve distributions over structures and/or parameters rather than specific settings and reduce to relative entropy projections. This holds even when the data is not separable within the chosen parametric class, in the context of anomaly detection rather than classification, or when the labels in the training set are uncertain or incomplete. Support vector machines are naturally subsumed under this class and we provide several extensions. We are also able to estimate exactly and efficiently discriminative distributions over tree structures of class-conditional models within this framework. Preliminary experimental results are indicative of the potential in these techniques.

The Relaxed Online Maximum Margin Algorithm

Yi Li, Philip Long

We describe a new incremental algorithm for training linear threshold functions: the Relaxed Online Maximum Margin Algorithm, or ROMMA. ROMMA can be viewed as an approximation to the algorithm that repeatedly chooses the hyperplane that classifies previously seen examples correctly with the maximum margin. It is known that such a maximum-margin hypothesis can be computed by minimizing the length of the weight vector subject to a number of linear constraints. ROMMA works by maintaining a relatively simple relaxation of these constraints that can be efficiently updated. We prove a mistake bound for ROMMA that is the same as that proved for the perceptron algorithm. Our analysis implies that the more computationally intensive maximum-margin algorithm also satisfies this mistake bound; this is the first worst-case performance guarantee for this algorithm. We describe some experiments using ROMMA and a variant that updates its hypothesis more aggressively as batch algorithms to recognize handwritten digits. The computational complexity and simplicity of these algorithms is similar to t

hat of per(cid:173) ceptron algorithm , but their generalization is much better.

We describe a sense in which the performance of ROMMA converges to that of SVM in the limit if bias isn't considered.

Bayesian Modelling of fMRI time Series

Pedro Højén-Sørensen, Lars Hansen, Carl Rasmussen

We present a Hidden Markov Model (HMM) for inferring the hidden psychological state (or neural activity) during single trial fMRI activation experiments with blocked task paradigms. Inference is based on Bayesian methodology, using a combination of analytical and a variety of Markov Chain Monte Carlo (MCMC) sampling techniques. The advantage of this method is that detection of short time learning effects between repeated trials is possible since inference is based only on single trial experiments.

Bayesian Averaging is Well-Tempered

Lars Hansen

Bayesian predictions are stochastic just like predictions of any other inference scheme that generalize from a finite sample. While a simple variational argument shows that Bayes averaging is generalization optimal given that the prior matches the teacher parameter distribution the situation is less clear if the teacher distribution is unknown. I define a class of averaging procedures, the tempered likelihoods, including both Bayes averaging with a uniform prior and maximum likelihood estimation as special cases. I show that Bayes is generalization optimal in this family for any teacher distribution for two learning problems that are analytically tractable: learning the mean of a Gaussian and asymptotics of smooth learners.

Policy Search via Density Estimation

Andrew Ng, Ronald Parr, Daphne Koller

We propose a new approach to the problem of searching a space of stochastic controllers for a Markov decision process (MDP) or a partially observable Markov decision process (POMDP). Following several other authors, our approach is based on searching in parameterized families of policies (for example, via gradient descent) to optimize solution quality. However, rather than trying to estimate the values and derivatives of a policy directly, we do so indirectly using estimates for the probability densities that the policy induces on states at the different points in time. This enables our algorithms to exploit the many techniques for efficient and robust approximate density propagation in stochastic systems. We show how our techniques can be applied both to deterministic propagation schemes (where the MDP's dynamics are given explicitly in compact form,) and to stochastic propagation schemes (where we have access only to a generative model, or simulator, of the MDP). We present empirical results for both of these variants on complex problems.

Low Power Wireless Communication via Reinforcement Learning

Timothy Brown

This paper examines the application of reinforcement learning to a wireless communication problem. The problem requires that channel utility be maximized while simultaneously minimizing battery usage. We present a solution to this multi-criteria problem that is able to significantly reduce power consumption. The solution uses a variable discount factor to capture the effects of battery usage.

Learning to Parse Images

Geoffrey E. Hinton, Zoubin Ghahramani, Yee Whye Teh

We describe a class of probabilistic models that we call credibility networks. Using parse trees as internal representations of images, credibility networks are able to perform segmentation and recognition simultaneously

ously, removing the need for ad hoc segmentation heuristics. Promising results in the problem of segmenting handwritten digits were obtained.

Robust Recognition of Noisy and Superimposed Patterns via Selective Attention Soo-Young Lee, Michael C. Mozer

In many classification tasks, recognition accuracy is low because input patterns are corrupted by noise or are spatially or temporally overlapping. We propose an approach to overcoming these limitations based on a model of human selective attention. The model, an early selection filter guided by top-down attentional control, entertains each candidate output class in sequence and adjusts attentional gain coefficients in order to produce a strong response for that class. The chosen class is then the one that obtains the strongest response with the least modulation of attention. We present simulation results on classification of corrupted and superimposed handwritten digit patterns, showing a significant improvement in recognition rates. The algorithm has also been applied in the domain of speech recognition, with comparable results.

Bayesian Network Induction via Local Neighborhoods Dimitris Margaritis, Sebastian Thrun

In recent years, Bayesian networks have become highly successful tools for diagnosis, prognosis, analysis, and decision making in real-world domains. We present an efficient algorithm for learning Bayesian networks from data. Our approach constructs Bayesian networks by first identifying each node's Markov blankets, then connecting nodes in a maximally consistent way. In contrast to the majority of work, which typically uses hill-climbing approaches that may produce dense and causally incorrect nets, our approach yields much more compact causal networks by heeding independencies in the data. Compact causal networks facilitate fast inference and are also easier to understand. We prove that under mild assumptions, our approach requires time polynomial in the size of the data and the number of nodes. A randomized variant, also presented here, yields comparable results at much higher speeds.

Spiking Boltzmann Machines

Geoffrey E. Hinton, Andrew Brown

We first show how to represent sharp posterior probability distributions using real valued coefficients on broadly-tuned basis functions. Then we show how the precise times of spikes can be used to convey the real-valued coefficients on the basis functions quickly and accurately. Finally we describe a simple simulation in which spiking neurons learn to model an image sequence by fitting a dynamic generative model.

Actor-Critic Algorithms

Vijay Konda, John Tsitsiklis

We propose and analyze a class of actor-critic algorithms for simulation-based optimization of a Markov decision process over a parameterized family of randomized stationary policies. These are two-time-scale algorithms in which the critic uses TD learning with a linear approximation architecture and the actor is updated in an approximate gradient direction based on information provided by the critic. We show that the features for the critic should span a subspace prescribed by the choice of parameterization of the actor. We conclude by discussing convergence properties and some open problems.

Training Data Selection for Optimal Generalization in Trigonometric Polynomial Networks

Masashi Sugiyama, Hidemitsu Ogawa

In this paper, we consider the problem of active learning in trigonometric polynomial networks and give a necessary and sufficient condition

on of sample points to provide the optimal generalization capability. By analyzing the condition from the functional analytic point of view, we clarify the mechanism of achieving the optimal generalization capability. We also show that a set of training examples satisfying the condition does not only provide the optimal generalization but also reduces the computational complexity and memory required for the calculation of learning results. Finally, examples of sample points satisfying the condition are given and computer simulations are performed to demonstrate the effectiveness of the proposed active learning method.

Image Recognition in Context: Application to Microscopic Urinalysis

Xubo Song, Joseph Sill, Yaser Abu-Mostafa, Harvey Kasdan

We propose a new and efficient technique for incorporating contextual information into object classification. Most of the current techniques face the problem of exponential computation cost. In this paper, we propose a new general framework that incorporates partial context at a linear cost. This technique is applied to microscopic urinalysis image recognition, resulting in a significant improvement of recognition rate over the context free approach. This gain would have been impossible using conventional context incorporation techniques.

Bayesian Map Learning in Dynamic Environments

Kevin P. Murphy

We consider the problem of learning a grid-based map using a robot with noisy sensors and actuators. We compare two approaches: online EM, where the map is treated as a fixed parameter, and Bayesian inference, where the map is a (matrix-valued) random variable. We show that even on a very simple example, online EM can get stuck in local minima, which causes the robot to get "lost" and the resulting map to be useless. By contrast, the Bayesian approach, by maintaining multiple hypotheses, is much more robust. We then introduce a method for approximating the Bayesian solution, called Rao-Blackwellised particle filtering. We show that this approximation, when coupled with an active learning strategy, is fast but accurate.

Better Generative Models for Sequential Data Problems: Bidirectional Recurrent Mixture Density Networks

Mike Schuster

This paper describes bidirectional recurrent mixture density networks, which can model multi-modal distributions of the type $P(X_t | y^f)$ and $P(X_t | x^I, x^2, \dots, x_{t-1}, y^f)$ without any explicit assumptions about the use of context. These expressions occur frequently in pattern recognition problems with sequential data, for example in speech recognition. Experiments show that the proposed generative models give a higher likelihood on test data compared to a traditional modeling approach, indicating that they can summarize the statistical properties of the data better.

Statistical Dynamics of Batch Learning

Song Li, K. Y. Michael Wong

An important issue in neural computing concerns the description of learning dynamics with macroscopic dynamical variables. Recent progress on on-line learning only addresses the often unrealistic case of an infinite training set. We introduce a new framework to model batch learning of restricted sets of examples, widely applicable to any learning cost function, and fully taking into account the temporal correlations introduced by the recycling of the examples. For illustration we analyze the effects of weight decay and early stopping during the learning of teacher-generated examples.

Scale Mixtures of Gaussians and the Statistics of Natural Images

Martin J. Wainwright, Eero Simoncelli

The statistics of photographic images, when represented using multiscale (wavelet) bases, exhibit two striking types of non-Gaussian behavior. First, the marginal densities of the coefficients have extended heavy tails. Second, the joint densities exhibit variance dependencies not captured by second-order models. We examine properties of the class of Gaussian scale mixtures, and show that these densities can accurately characterize both the marginal and joint distributions of natural image wavelet coefficients. This class of model suggests a Markov structure, in which wavelet coefficients are linked by hidden scaling variables corresponding to local image structure. We derive an estimator for these hidden variables, and show that a nonlinear "normalization" procedure can be used to Gaussianize the coefficients.

Independent Factor Analysis with Temporally Structured Sources

Hagai Attias

We present a new technique for time series analysis based on dynamic probabilistic networks. In this approach, the observed data are modeled in terms of unobserved, mutually independent factors, as in the recently introduced technique of Independent Factor Analysis (IFA). However, unlike in IFA, the factors are not i.i.d.; each factor has its own temporal statistical characteristics. We derive a family of EM algorithms that learn the structure of the underlying factors and their relation to the data.

These algorithms perform source separation and noise reduction in an integrated manner, and demonstrate superior performance compared to IFA.

Managing Uncertainty in Cue Combination

Zhiyong Yang, Richard Zemel

We develop a hierarchical generative model to study cue combination. The model maps a global shape parameter to local cue-specific parameters, which in turn generate an intensity image. Inferring shape from images is achieved by inverting this model. Inference produces a probability distribution at each level; using distributions rather than a single value of underlying variables at each stage preserves information about the validity of each local cue for the given image. This allows the model, unlike standard combination models, to adaptively weight each cue based on general cue reliability and specific image context. We describe the results of a cue combination psychophysics experiment we conducted that allows a direct comparison with the model. The model provides a good fit to our data and a natural account for some interesting aspects of cue combination.

Potential Boosters?

Nigel Duffy, David Helmbold

Recent interpretations of the Adaboost algorithm view it as performing a gradient descent on a potential function. Simply changing the potential function allows one to create new algorithms related to AdaBoost. However, these new algorithms are generally not known to have the formal boosting property. This paper examines the question of which potential functions lead to new algorithms that are boosters. The two main results are general sets of conditions on the potential; one set implies that the resulting algorithm is a booster, while the other implies that the algorithm is not. These conditions are applied to previously studied potential functions, such as those used by LogitBoost and Doolittle.

Resonance in a Stochastic Neuron Model with Delayed Interaction

Toru Ohira, Yuzuru Sato, Jack Cowan

We study here a simple stochastic single neuron model with delayed self-feedback capable of generating spike trains. Simulations show that its spike trains exhibit resonant behavior between "noise" and "delay". In order to gain insight into this resonance, we simplify the model and study a stochastic bin

ary element whose transition probability depends on its state at a fixed interval in the past. With this simplified model we can analytically compute interspike interval histograms, and show how the resonance between noise and delay arises. The resonance is also observed when such elements are coupled through delayed interaction.

Wiring Optimization in the Brain

Dmitri Chklovskii, Charles Stevens

The complexity of cortical circuits may be characterized by the number of synapses per neuron. We study the dependence of complexity on the fraction of the cortical volume that is made up of "wire" (that is, of axons and dendrites), and find that complexity is maximized when wire takes up about 60% of the cortical volume. This prediction is in good agreement with experimental observations. A consequence of our arguments is that any rearrangement of neurons that takes more wire would sacrifice computational power.

Learning from User Feedback in Image Retrieval Systems

Nuno Vasconcelos, Andrew Lippman

We formulate the problem of retrieving images from visual databases as a problem of Bayesian inference. This leads to natural and effective solutions for two of the most challenging issues in the design of a retrieval system: providing support for region-based queries without requiring prior image segmentation, and accounting for user-feedback during a retrieval session. We present a new learning algorithm that relies on belief propagation to account for both positive and negative examples of the user's interests.

Online Independent Component Analysis with Local Learning Rate Adaptation

Nicol Schraudolph, Xavier Giannakopoulos

Stochastic meta-descent (SMD) is a new technique for online adaptation of local learning rates in arbitrary twice-differentiable systems. Like matrix momentum it uses full second-order information while retaining $O(n)$ computational complexity by exploiting the efficient computation of Hessian-vector products. Here we apply SMD to independent component analysis, and employ the resulting algorithm for the blind separation of time-varying mixtures. By matching individual learning rates to the rate of change in each source signal's mixture coefficients, our technique is capable of simultaneously tracking sources that move at very different, a priori unknown speeds.

A Variational Bayesian Framework for Graphical Models

Hagai Attias

This paper presents a novel practical framework for Bayesian model averaging and model selection in probabilistic graphical models. Our approach approximates full posterior distributions over model parameters and structures, as well as latent variables, in an analytical manner. These posteriors fall out of a free-form optimization procedure, which naturally incorporates conjugate priors. Unlike in large sample approximations, the posteriors are generally non-Gaussian and no Hessian needs to be computed. Predictive quantities are obtained analytically. The resulting algorithm generalizes the standard Expectation Maximization algorithm, and its convergence is guaranteed. We demonstrate that this approach can be applied to a large class of models in several domains, including mixture models and source separation.

Algebraic Analysis for Non-regular Learning Machines

Sumio Watanabe

Hierarchical learning machines are non-regular and non-identifiable statistical models, whose true parameter sets are analytic sets with singularities. Using algebraic analysis, we rigorously prove that the stochastic complexity

of a non-identifiable learning machine $(m_1 - 1) \log \log n + \text{const.}$, is asymptotically equal to $>1 \log n$ - where n is the number of training samples. Moreover we show that the rational number >1 and the integer m_1 can be algorithmically calculated using resolution of singularities in algebraic geometry. Also we obtain inequalities $0 < >1 \sim d/2$ and $1 \sim m_1 \sim d$, where d is the number of parameters.

Model Selection in Clustering by Uniform Convergence Bounds

Joachim Buhmann, Marcus Held

Unsupervised learning algorithms are designed to extract structure from data samples. Reliable and robust inference requires a guarantee that extracted structures are typical for the data source, i.e., similar structures have to be inferred from a second sample set of the same data source. The overfitting phenomenon in maximum entropy based annealing algorithms is exemplarily studied for a class of histogram clustering models. Bernstein's inequality for large deviations is used to determine the maximally achievable approximation quality parameterized by a minimal temperature. Monte Carlo simulations support the proposed model selection criterion by finite temperature annealing.

Unmixing Hyperspectral Data

Lucas Parra, Clay Spence, Paul Sajda, Andreas Ziehe, Klaus-Robert Müller

In hyperspectral imagery one pixel typically consists of a mixture of the reflectance spectra of several materials, where the mixture coefficients correspond to the abundances of the constituting materials. We assume linear combinations of reflectance spectra with some additive sensor noise and derive a probabilistic MAP framework for analyzing hyperspectral data. As the material reflectance characteristics are not known a priori, we face the problem of unsupervised linear unmixing. The incorporation of different prior information (e.g. positivity and normalization of the abundances) naturally leads to a family of interesting algorithms, for example in the noise-free case yielding an algorithm that can be understood as constrained independent component analysis (ICA). Simulations underline the usefulness of our theory.

Some Theoretical Results Concerning the Convergence of Compositions of Regularized Linear Functions

Tong Zhang

Recently, sample complexity bounds have been derived for problems involving linear functions such as neural networks and support vector machines. In this paper, we extend some theoretical results in this area by deriving dimensional independent covering number bounds for regularized linear functions under certain regularization conditions. We show that such bounds lead to a class of new methods for training linear classifiers with similar theoretical advantages of the support vector machine. Furthermore, we also present a theoretical analysis for these new methods from the asymptotic statistical point of view. This technique provides better description for large sample behaviors of these algorithms.

Inference for the Generalization Error

Claude Nadeau, Yoshua Bengio

In order to compare learning algorithms, experimental results reported in the machine learning literature often use statistical tests of significance. Unfortunately, most of these tests do not take into account the variability due to the choice of training set. We perform a theoretical investigation of the variance of the cross-validation estimate of the generalization error that takes into account the variability due to the choice of training sets. This allows us to propose two new ways to estimate this variance. We show, via simulations, that these new statistics perform well relative to the statistics considered by Dietterich (Dietterich, 199

8).

LTD Facilitates Learning in a Noisy Environment

Paul Munro, Gerardina Hernández

Long-term potentiation (LTP) has long been held as a biological substrate for a ssociative learning. Recently, evidence has emerged that long-term depression (LTD) results when the presynaptic cell fires after the postsynaptic cell. The computational utility of LTD is explored here. Synaptic modification kernels for both LTP and LTD have been proposed by other laboratories based studies of one postsynaptic unit. Here, the interaction between time-dependent LTP and LTD is studied in small networks.

An Oculo-Motor System with Multi-Chip Neuromorphic Analog VLSI Control

Oliver Landolt, Steve Gyger

A system emulating the functionality of a moving eye-hence the name oculo-motor system-has been built and successfully tested. It is made of an optical device for shifting the field of view of an image sensor by up to 45 ° in any direction, four neuromorphic analog VLSI circuits implementing an oculo-motor control loop, and some off-the-shelf electronics. The custom integrated circuits communicate with each other primarily by non-arbitrated address-event buses. The system implements the behaviors of saliency-based saccadic exploration, and smooth pursuit of light spots. The duration of saccades ranges from 45 ms to 100 ms, which is comparable to human eye performance. Smooth pursuit operates on light sources moving at up to 50 °/s in the visual field.

Understanding Stepwise Generalization of Support Vector Machines: a Toy Model

Sebastian Risau-Gusman, Mirta Gordon

In this article we study the effects of introducing structure in the input distribution of the data to be learnt by a simple perceptron. We determine the learning curves within the framework of Statistical Mechanics. Stepwise generalization occurs as a function of the number of examples when the distribution of patterns is highly anisotropic. Although extremely simple, the model seems to capture the relevant features of a class of Support Vector Machines which was recently shown to present this behavior.

A Geometric Interpretation of v-SVM Classifiers

David Crisp, Christopher J. C. Burges

We show that the recently proposed variant of the Support Vector machine (SVM) algorithm, known as v-SVM, can be interpreted as a maximal separation between subsets of the convex hulls of the data, which we call soft convex hulls. The soft convex hulls are controlled by choice of the parameter v . If the intersection of the convex hulls is empty, the hyperplane is positioned halfway between them such that the distance between convex hulls, measured along the normal, is maximized; and if it is not, the hyperplane's normal is similarly determined by the soft convex hulls, but its position (perpendicular distance from the origin) is adjusted to minimize the error sum. The proposed geometric interpretation of v-SVM also leads to necessary and sufficient conditions for the existence of a choice of v for which the v-SVM solution is nontrivial.

Robust Learning of Chaotic Attractors

Rembrandt Bakker, Jaap Schouten, Marc-Olivier Coppins, Floris Takens, C. Giles, Cor van den Bleek

A fundamental problem with the modeling of chaotic time series data is that minimizing short-term prediction errors does not guarantee a match between the reconstructed attractors of model and experiments. We introduce a modeling paradigm that simultaneously learns to short-term predict and to locate the outlines of the attractor by a new way of nonlinear principal component analysis. Closed-

loop predictions are constrained to stay within these outlines, to prevent divergence from the attractor. Learning is exceptionally fast: parameter estimation for the 1000 sample laser data from the 1991 Santa Fe time series competition took less than a minute on a 166 MHz Pentium PC.

Greedy Importance Sampling

Dale Schuurmans

I present a simple variation of importance sampling that explicitly searches for important regions in the target distribution. I prove that the technique yields unbiased estimates, and show empirically it can reduce the variance of standard Monte Carlo estimators. This is achieved by concentrating samples in more significant regions of the sample space.

Recurrent Cortical Competition: Strengthen or Weaken?

Péter Adorján, Lars Schwabe, Christian Piepenbrock, Klaus Obermayer

We investigate the short term dynamics of the recurrent competition and neural activity in the primary visual cortex in terms of information processing and in the context of orientation selectivity. We propose that after stimulus onset, the strength of the recurrent excitation decreases due to fast synaptic depression. As a consequence, the network shifts from an initially highly nonlinear to a more linear operating regime. Sharp orientation tuning is established in the first highly competitive phase. In the second and less competitive phase, precise signaling of multiple orientation selectivities and long range modulation, e.g., by intra- and inter-areal connections becomes possible (surround effects). Thus the network first extracts the salient features from the stimulus, and then starts to process the details. We show that this signal processing strategy is optimal if the neurons have limited bandwidth and their objective is to transmit the maximum amount of information in any time interval beginning with the stimulus onset.

Constrained Hidden Markov Models

Sam Roweis

By thinking of each state in a hidden Markov model as corresponding to some spatial region of a fictitious topology space it is possible to naturally define neighbouring states as those which are connected in that space. The transition matrix can then be constrained to allow transitions only between neighbours; this means that all valid state sequences correspond to connected paths in the topology space. I show how such constrained HMMs can learn to discover underlying structure in complex sequences of high dimensional data, and apply them to the problem of recovering mouth movements from acoustics in continuous speech.

Approximate Inference Algorithms for Two-Layer Bayesian Networks

Andrew Ng, Michael Jordan

We present a class of approximate inference algorithms for graphical models of the QMR-DT type. We give convergence rates for these algorithms and for the Jaakkola and Jordan (1999) algorithm, and verify these theoretical predictions empirically. We also present empirical results on the difficult QMR-DT network problem, obtaining performance of the new algorithms roughly comparable to the Jaakkola and Jordan algorithm.

Perceptual Organization Based on Temporal Dynamics

Xiuwen Liu, DeLiang Wang

A figure-ground segregation network is proposed based on a novel boundary pair representation. Nodes in the network are boundary segments obtained through local grouping. Each node is excitatorily coupled with the neighboring nodes that belong to the same region, and inhibitorily coupled with the corresponding paired node. Gestalt grouping rules are incorporated by modulating connections. The status of a node

ode represents its probability being figural and is updated according to a differential equation. The system solves the figure-ground segregation problem through temporal evolution. Different perceptual phenomena, such as modal and amodal completion, virtual contours, grouping and shape composition are then explained through local diffusion. The system eliminates combinatorial optimization and accounts for many psychophysical results with a fixed set of parameters.

Learning Informative Statistics: A Nonparametric Approach

John W. Fisher III, Alexander Ihler, Paul Viola

We discuss an information theoretic approach for categorizing and modeling dynamic processes. The approach can learn a compact and informative statistic which summarizes past states to predict future observations. Furthermore, the uncertainty of the prediction is characterized nonparametrically by a joint density over the learned statistic and present observation. We discuss the application of the technique to both noise driven dynamical systems and random processes sampled from a density which is conditioned on the past. In the first case we show results in which both the dynamics of random walk and the statistics of the driving noise are captured. In the second case we present results in which a summarizing statistic is learned on noisy random telegraph waves with differing dependencies on past states. In both cases the algorithm yields a principled approach for discriminating processes with differing dynamics and/or dependencies. The method is grounded in ideas from information theory and nonparametric statistics.

Rules and Similarity in Concept Learning

Joshua Tenenbaum

This paper argues that two apparently distinct modes of generalizing concepts - abstracting rules and computing similarity to exemplars - should both be seen as special cases of a more general Bayesian learning framework. Bayes explains the specific workings of these two modes - which rules are abstracted, how similarity is measured - as well as why generalization should appear rule- or similarity-based in different situations. This analysis also suggests why the rules/similarity distinction, even if not computationally fundamental, may still be useful at the algorithmic level as part of a principled approximation to fully Bayesian learning.

Support Vector Method for Novelty Detection

Bernhard Schölkopf, Robert C. Williamson, Alex Smola, John Shawe-Taylor, John Platt

Suppose you are given some dataset drawn from an underlying probability distribution P and you want to estimate a "simple" subset S of input space such that the probability that a test point drawn from P lies outside of S equals some a priori specified $1/\ell$ between 0 and 1. We propose a method to approach this problem by trying to estimate a function f which is positive on S and negative on the complement. The functional form of f is given by a kernel expansion in terms of a potentially small subset of the training data; it is regularized by controlling the length of the weight vector in an associated feature space. We provide a theoretical analysis of the statistical performance of our algorithm. The algorithm is a natural extension of the support vector algorithm to the case of unlabelled data.

Generalized Model Selection for Unsupervised Learning in High Dimensions

Shivakumar Vaithyanathan, Byron Dom

We describe a Bayesian approach to model selection in unsupervised learning that determines both the feature set and the number of clusters. We then evaluate this scheme (based on marginal likelihood) and one based on cross-validated likelihood. For the Bayesian scheme we derive a closed-form solution of the marginal likelihood by assuming appropriate forms of

the likelihood function and prior. Extensive experiments compare these approaches and all results are verified by comparison against ground truth.

In these experiments the Bayesian scheme using our objective function gave better results than cross-validation.

An Improved Decomposition Algorithm for Regression Support Vector Machines

Pavel Laskov

A new decomposition algorithm for training regression Support Vector Machines (SVM) is presented. The algorithm builds on the basic principles of decomposition proposed by Osuna et. al., and addresses the issue of optimal working set selection. The new criteria for testing optimality of a working set are derived. Based on these criteria, the principle of "maximal inconsistency" is proposed to form (approximately) optimal working sets. Experimental results show superior performance of the new algorithm in comparison with traditional training of regression SVM without decomposition. Similar results have been previously reported on decomposition algorithms for pattern recognition SVM. The new algorithm is also applicable to advanced SVM formulations based on regression, such as density estimation and integral equation SVM.

An Analog VLSI Model of Periodicity Extraction

André van Schaik

that extracts

From Coexpression to Coregulation: An Approach to Inferring Transcriptional Regulation among Gene Classes from Large-Scale Expression Data

Eric Mjolsness, Tobias Mann, Rebecca Castaño, Barbara Wold

small-scale gene

Data Visualization and Feature Selection: New Algorithms for Nongaussian Data

Howard Yang, John Moody

Data visualization and feature selection methods are proposed based on the joint mutual information and ICA. The visualization methods can find many good 2-D projections for high dimensional data interpretation, which cannot be easily found by the other existing methods. The new variable selection method is found to be better in eliminating redundancy in the inputs than other methods based on simple mutual information. The efficacy of the methods is illustrated on a radar signal analysis problem to find 2-D viewing coordinates for data visualization and to select inputs for a neural network classifier. Keywords: feature selection, joint mutual information, ICA, visualization, classification.

An Information-Theoretic Framework for Understanding Saccadic Eye Movements

Tai Sing Lee, Stella Yu

In this paper, we propose that information maximization can provide a unified framework for understanding saccadic eye movements. In this framework, the mutual information among the correlative representations of the retinal image, the priors constructed from our long term visual experience, and a dynamic short-term internal representation constructed from recent saccades provides a map for guiding eye navigation. By directing the eyes to locations of maximum complexity in neuronal ensemble responses at each step, the automatic saccadic eye movement system greedily collects information about the external world, while modifying the neural representations in the process. This framework attempts to connect several psychological phenomena, such as pop-out and inhibition of return, to long term visual experience and short term working memory. It also provides an interesting perspective on contextual computation and formation of neural representation in the visual system.

Noisy Neural Networks and Generalizations

Hava Siegelmann, Alexander Roitershtein, Asa Ben-Hur

In this paper we define a probabilistic computational model which generalizes many noisy neural network models, including the recent work of Maass and Sontag [5]. We identify weak ergodicity as the mechanism responsible for restriction of the computational power of probabilistic models to definite languages, independent of the characteristics of the noise: whether it is discrete or analog, or if it depends on the input or not, and independent of whether the variables are discrete or continuous. We give examples of weakly ergodic models including noisy computational systems with noise depending on the current state and inputs, aggregate models, and computational systems which update in continuous time.

The Nonnegative Boltzmann Machine

Oliver Downs, David MacKay, Daniel Lee

The nonnegative Boltzmann machine (NNBM) is a recurrent neural network model that can describe multimodal nonnegative data. Application of maximum likelihood estimation to this model gives a learning rule that is analogous to the binary Boltzmann machine. We examine the utility of the mean field approximation for the NNBM, and describe how Monte Carlo sampling techniques can be used to learn its parameters. Reflective slice sampling is particularly well-suited for this distribution, and can efficiently be implemented to sample the distribution. We illustrate learning of the NNBM on a translationally invariant distribution, as well as on a generative model for images of human faces.

Boosting Algorithms as Gradient Descent

Llew Mason, Jonathan Baxter, Peter Bartlett, Marcus Frean

We provide an abstract characterization of boosting algorithms as gradient descent on cost-functionals in an inner-product function space. We prove convergence of these functional-gradient-descent algorithms under quite weak conditions. Following previous theoretical results bounding the generalization performance of convex combinations of classifiers in terms of general cost functions of the margin, we present a new algorithm (DOOM II) for performing a gradient descent optimization of such cost functions. Experiments on several data sets from the UC Irvine repository demonstrate that DOOM II generally outperforms AdaBoost, especially in high noise situations, and that the overfitting behaviour of AdaBoost is predicted by our cost functions.

Local Probability Propagation for Factor Analysis

Brendan J. Frey

Ever since Pearl's probability propagation algorithm in graphs with cycles was shown to produce excellent results for error-correcting decoding a few years ago, we have been curious about whether local probability propagation could be used successfully for machine learning. One of the simplest adaptive models is the factor analyzer, which is a two-layer network that models bottom layer sensory inputs as a linear combination of top layer factors plus independent Gaussian sensor noise. We show that local probability propagation in the factor analyzer network usually takes just a few iterations to perform accurate inference, even in networks with 320 sensors and 80 factors. We derive an expression for the algorithm's fixed point and show that this fixed point matches the exact solution in a variety of networks, even when the fixed point is unstable. We also show that this method can be used successfully to perform inference for approximate EM and we give results on an online face recognition task. A simple way to encode input patterns is to suppose that each input can be well approximated by a linear combination of component vectors, where the amplitudes of the vectors are modulated to match the input. For a given training set, the most appropriate set of component vectors will depend on how we expect the modulation levels to behave and how

we measure the distance between the input and its approximation. These effects can be captured by a generative probabilistic model that specifies a distribution $p(z)$ over modulation levels $z = (z_1, \dots, z_K)$ and a distribution $p(x|z)$ over sensors $x = (x_1, \dots, x_N)^T$ given the modulation levels. Principal component analysis, independent component analysis and factor analysis can be viewed as maximum likelihood learning in a model of this type, where we assume that over the training set, the appropriate modulation levels are independent and the overall distortion is given by the sum of the individual sensor distortions.

A MCMC Approach to Hierarchical Mixture Modelling

Christopher Williams

There are many hierarchical clustering algorithms available, but these lack a firm statistical basis. Here we set up a hierarchical probabilistic mixture model, where data is generated in a hierarchical tree-structured manner. Markov chain Monte Carlo (MCMC) methods are demonstrated which can be used to sample from the posterior distribution over trees containing variable numbers of hidden units.

The Infinite Gaussian Mixture Model

Carl Rasmussen

In a Bayesian mixture model it is not necessary a priori to limit the number of components to be finite. In this paper an infinite Gaussian mixture model is presented which neatly sidesteps the difficult problem of finding the "right" number of mixture components. Inference in the model is done using an efficient parameter-free Markov Chain that relies entirely on Gibbs sampling.

Reconstruction of Sequential Data with Probabilistic Models and Continuity Constraints

Miguel Carreira-Perpiñán

We consider the problem of reconstructing a temporal discrete sequence of multidimensional real vectors when part of the data is missing, under the assumption that the sequence was generated by a continuous process. A particular case of this problem is multivariate regression, which is very difficult when the underlying mapping is one-to-many. We propose an algorithm based on a joint probability model of the variables of interest, implemented using a nonlinear latent variable model. Each point in the sequence is potentially reconstructed as any of the modes of the conditional distribution of the missing variables given the present variables (computed using an exhaustive mode search in a Gaussian mixture). Mode selection is determined by a dynamic programming search that minimises a geometric measure of the reconstructed sequence, derived from continuity constraints. We illustrate the algorithm with a toy example and apply it to a real-world inverse problem, the acoustic-to-articulatory mapping. The results show that the algorithm outperforms conditional mean imputation and multilayer perceptrons.

Learning the Similarity of Documents: An Information-Geometric Approach to Document Retrieval and Categorization

Thomas Hofmann

The project pursued in this paper is to develop from first information-geometric principles a general method for learning the similarity between text documents. Each individual document is modeled as a memoryless information source. Based on a latent class decomposition of the term-document matrix, a low-dimensional (curved) multinomial subfamily is learned. From this model a canonical similarity function - known as the Fisher kernel - is derived. Our approach can be applied for unsupervised and supervised learning problems alike. This in particular covers interesting cases where both, labeled and unlabeled data are

available. Experiments in automated indexing and text categorization verify the advantages of the proposed method.

Bayesian Reconstruction of 3D Human Motion from Single-Camera Video

Nicholas Howe, Michael Leventon, William Freeman

The three-dimensional motion of humans is underdetermined when the observation is limited to a single camera, due to the inherent 3D ambiguity of 2D video. We present a system that reconstructs the 3D motion of human subjects from single-camera video, relying on prior knowledge about human motion, learned from training data, to resolve those ambiguities. After initialization in 2D, the tracking and 3D reconstruction is automatic; we show results for several video sequences. The results show the power of treating 3D body tracking as an inference problem.

Mixture Density Estimation

Jonathan Li, Andrew Barron

Gaussian mixtures (or so-called radial basis function networks) for density estimation provide a natural counterpart to sigmoidal neural networks for function fitting and approximation. In both cases, it is possible to give simple expressions for the iterative improvement of performance as components of the network are introduced one at a time. In particular, for mixture density estimation we show that a k -component mixture estimated by maximum likelihood (or by an iterative likelihood improvement that we introduce) achieves log-likelihood within order $1/k$ of the log-likelihood achievable by any convex combination. Consequences for approximation and estimation using Kullback-Leibler risk are also given. A Minimum Description Length principle selects the optimal number of components k that minimizes the risk bound.

Information Capacity and Robustness of Stochastic Neuron Models

Elad Schneidman, Idan Segev, Naftali Tishby

The reliability and accuracy of spike trains have been shown to depend on the nature of the stimulus that the neuron encodes. Adding ion channel stochasticity to neuronal models results in a macroscopic behavior that replicates the input-dependent reliability and precision of real neurons. We calculate the amount of information that an ion channel based stochastic Hodgkin-Huxley (HH) neuron model can encode about a wide set of stimuli. We show that both the information rate and the information per spike of the stochastic model are similar to the values reported experimentally. Moreover, the amount of information that the neuron encodes is correlated with the amplitude of fluctuations in the input, and less so with the average firing rate of the neuron. We also show that for the HH ion channel density, the information capacity is robust to changes in the density of ion channels in the membrane, whereas changing the ratio between the Na^+ and K^+ ion channels has a considerable effect on the information that the neuron can encode. Finally, we suggest that neurons may maximize their information capacity by appropriately balancing the density of the different ion channels that underlie neuronal excitability.

Bayesian Transduction

Thore Graepel, Ralf Herbrich, Klaus Obermayer

Transduction is an inference principle that takes a training sample and aims at estimating the values of a function at given points contained in the so-called working sample as opposed to the whole of input space for induction. Transduction provides a confidence measure on single predictions rather than classifiers - a feature particularly important for risk-sensitive applications. The possibly infinite number of functions is reduced to a finite number of equivalence classes on the working sample.

A rigorous Bayesian analysis reveals that for standard classification loss we cannot benefit from considering more than one test point at a time

e. The probability of the label of a given test point is determined as the posterior measure of the corresponding subset of hypothesis space. We consider the PAC setting of binary classification by linear discriminant functions (perceptrons) in kernel space such that the probability of labels is determined by the volume ratio in version space. We suggest to sample this region by an ergodic billiard. Experimental results on real world data indicate that Bayesian Transduction compares favorably to the well-known Support Vector Machine, in particular if the posterior probability of labellings is used as a confidence measure to exclude test points of low confidence.

Constructing Heterogeneous Committees Using Input Feature Grouping: Application to Economic Forecasting

Yuansong Liao, John Moody

The committee approach has been proposed for reducing model uncertainty and improving generalization performance. The advantage of committees depends on (1) the performance of individual members and (2) the correlational structure of errors between members. This paper presents an input grouping technique for designing a heterogeneous committee. With this technique, all input variables are first grouped based on their mutual information. Statistically similar variables are assigned to the same group. Each member's input set is then formed by input variables extracted from different groups. Our designed committees have less error correlation between its members, since each member observes different input variable combinations. The individual member's feature sets contain less redundant information, because highly correlated variables will not be combined together. The member feature sets contain almost complete information, since each set contains a feature from each information group. An empirical study for a noisy and nonstationary economic forecasting problem shows that committees constructed by our proposed technique outperform committees formed using several existing techniques.

An Analysis of Turbo Decoding with Gaussian Densities

Paat Rusmevichientong, Benjamin Van Roy

We provide an analysis of the turbo decoding algorithm (TDA) in a setting involving Gaussian densities. In this context, we are able to show that the algorithm converges and that - somewhat surprisingly - though the density generated by the TDA may differ significantly from the desired posterior density, the means of these two densities coincide.

Probabilistic Methods for Support Vector Machines

Peter Sollich

I describe a framework for interpreting Support Vector Machines (SVMs) as maximum a posteriori (MAP) solutions to inference problems with Gaussian Process priors. This can provide intuitive guidelines for choosing a 'good' SVM kernel. It can also assign (by evidence maximization) optimal values to parameters such as the noise level C which cannot be determined unambiguously from properties of the MAP solution alone (such as cross-validation error). I illustrate this using a simple approximate expression for the SVM evidence. Once C has been determined, error bars on SVM predictions can also be obtained.

Neural System Model of Human Sound Localization

Craig Jin, Simon Carlile

This paper examines the role of biological constraints in the human auditory localization process. A psychophysical and neural system modeling approach was undertaken in which performance comparisons between competing models and a human subject explore the relevant biologically plausible "realism constraints". The directional acoustical cues, upon which sound localization is based, were derived from the human subject's head-

related transfer functions (HRTFs). Sound stimuli were generated by convolving bandpass noise with the HRTFs and were presented to both the subject and the model. The input stimuli to the model were processed using the Auditory Image Model of cochlear processing. The cochlear data was then analyzed by a time-delay neural network which integrated temporal and spectral information to determine the spatial location of the sound source. The combined cochlear model and neural network provided a system model of the sound localization process. Human-like localization performance was qualitatively achieved for broadband and bandpass stimuli when the model architecture incorporated frequency division (or tonotopicity), and was trained using variable bandwidth and center-frequency sounds.

A Neuromorphic VLSI System for Modeling the Neural Control of Axial Locomotion
Girish Patel, Edgar Brown, Stephen DeWeerth

We have developed and tested an analog/digital VLSI system that models the coordination of biological segmental oscillators underlying axial locomotion in animals such as leeches and lampreys. In its current form the system consists of a chain of twelve pattern generating circuits that are capable of arbitrary contralateral inhibitory synaptic coupling. Each pattern generating circuit is implemented with two independent silicon Morris-Lecar neurons with a total of 32 programmable (floating-gate based) inhibitory synapses, and an asynchronous address-event interaction connection element that provides synaptic connectivity and implements axonal delay. We describe and analyze the data from a set of experiments exploring the system behavior in terms of synaptic coupling.

Graded Grammaticality in Prediction Fractal Machines
Shan Parfitt, Peter Tiffo, Georg Dorffner

We introduce a novel method of constructing language models, which avoids some of the problems associated with recurrent neural networks. The method of creating a Prediction Fractal Machine (PFM) [1] is briefly described and some experiments are presented which demonstrate the suitability of PFMs for language modeling. PFMs distinguish reliably between minimal pairs, and their behavior is consistent with the hypothesis [4] that wellformedness is 'graded' not absolute. A discussion of their potential to offer fresh insights into language acquisition and processing follows.

Learning Sparse Codes with a Mixture-of-Gaussians Prior
Bruno Olshausen, K. Millman

We describe a method for learning an overcomplete set of basis functions for the purpose of modeling sparse structure in images. The sparsity of the basis function coefficients is modeled with a mixture-of-Gaussians distribution. One Gaussian captures non-active coefficients with a small-variance distribution centered at zero, while one or more other Gaussians capture active coefficients with a large-variance distribution. We show that when the prior is in such a form, there exist efficient methods for learning the basis functions as well as the parameters of the prior. The performance of the algorithm is demonstrated on a number of test cases and also on natural images. The basis functions learned on natural images are similar to those obtained with other methods, but the sparse form of the coefficient distribution is much better described. Also, since the parameters of the prior are adapted to the data, no assumption about sparse structure in the images need be made a priori, rather it is learned from the data.

Semiparametric Approach to Multichannel Blind Deconvolution of Nonminimum Phase Systems

Liqing Zhang, Shun-ichi Amari, Andrzej Cichocki

In this paper we discuss the semi parametric statistical model for blind deconvolution. First we introduce a Lie Group to the manifold of non(17

3) causal FIR filters. Then blind deconvolution problem is formulated in the framework of a semiparametric model, and a family of estimating functions is derived for blind deconvolution. A natural gradient learning algorithm is developed for training noncausal filters. Stability of the natural gradient algorithm is also analyzed in this framework.

Application of Blind Separation of Sources to Optical Recording of Brain Activity

Holger Schoner, Martin Stetter, Ingo Schiefl, John Mayhew, Jennifer Lund, Niall McLoughlin, Klaus Obermayer

In the analysis of data recorded by optical imaging from intrinsic signals (measurement of changes of light reflectance from cortical tissue) the removal of noise and artifacts such as blood vessel patterns is a serious problem. Often bandpass filtering is used, but the underlying assumption that a spatial frequency exists, which separates the mapping component from other components (especially the global signal), is questionable. Here we propose alternative ways of processing optical imaging data, using blind source separation techniques based on the spatial decorrelation of the data. We first perform benchmarks on artificial data in order to select the way of processing, which is most robust with respect to sensor noise. We then apply it to recordings of optical imaging experiments from macaque primary visual cortex. We show that our BSS technique is able to extract ocular dominance and orientation preference maps from single condition stacks, for data, where standard post-processing procedures fail. Artifacts, especially blood vessel patterns, can often be completely removed from the maps. In summary, our method for blind source separation using extended spatial decorrelation is a superior technique for the analysis of optical recording data.

Spectral Cues in Human Sound Localization

Craig Jin, Anna Corderoy, Simon Carlile, André van Schaik

The differential contribution of the monaural and interaural spectral cues to human sound localization was examined using a combined psychophysical and analytical approach. The cues to a sound's location were correlated on an individual basis with the human localization responses to a variety of spectrally manipulated sounds. The spectral cues derive from the acoustical filtering of an individual's auditory periphery which is characterized by the measured head-related transfer functions (HRTFs). Auditory localization performance was determined in virtual auditory space (VAS). Psychoacoustical experiments were conducted in which the amplitude spectra of the sound stimulus was varied independently at each ear while preserving the normal timing cues, an impossibility in the free-field environment. Virtual auditory noise stimuli were generated over earphones for a specified target direction such that there was a "false" flat spectrum at the left eardrum. Using the subject's HRTFs, the sound spectrum at the right eardrum was then adjusted so that either the true right monaural spectral cue or the true interaural spectral cue was preserved. All subjects showed systematic mislocalizations in both the true right and true interaural spectral conditions which was absent in their control localization performance. The analysis of the different cues along with the subjects' localization responses suggests there are significant differences in the use of the monaural and interaural spectral cues and that the auditory system's reliance on the spectral cues varies with the sound condition.

Robust Full Bayesian Methods for Neural Networks

Christophe Andrieu, João de Freitas, Arnaud Doucet

In this paper, we propose a full Bayesian model for neural networks. This model treats the model dimension (number of neurons), model parameters, regularization parameters and noise parameters as random variables that need to be estimated. We then propose a reversible jump Markov chain Monte Carlo

Monte Carlo (MCMC) method to perform the necessary computations. We find that the results are not only better than the previously reported ones, but also appear to be robust with respect to the prior specification. Moreover, we present a geometric convergence theorem for the algorithm.

A Neurodynamical Approach to Visual Attention

Gustavo Deco, Josef Zihl

The psychophysical evidence for "selective attention" originates mainly from visual search experiments. In this work, we formulate a hierarchical system of interconnected modules consisting in populations of neurons for modeling the underlying mechanisms involved in selective visual attention. We demonstrate that our neural system for visual search works across the visual field in parallel but due to the different intrinsic dynamics can show the two experimentally observed modes of visual attention, namely: the serial and the parallel search mode. In other words, neither explicit model of a focus of attention nor saliency maps are used. The focus of attention appears as an emergent property of the dynamic behavior of the system. The neural population dynamics are handled in the framework of the mean-field approximation. Consequently, the whole process can be expressed as a system of coupled differential equations.

Dynamics of Supervised Learning with Restricted Training Sets and Noisy Teachers

Anthony Coolen, C. Mace

We generalize a recent formalism to describe the dynamics of supervised learning in layered neural networks, in the regime where data recycling is inevitable, to the case of noisy teachers. Our theory generates reliable predictions for the evolution in time of training- and generalization errors, and extends the class of mathematically solvable learning processes in large neural networks to those situations where overfitting can occur.

Audio Vision: Using Audio-Visual Synchrony to Locate Sounds

John Hershey, Javier Movellan

Psychophysical and physiological evidence shows that sound localization of acoustic signals is strongly influenced by their synchrony with visual signals. This effect, known as ventriloquism, is at work when sound coming from the side of a TV set feels as if it were coming from the mouth of the actors. The ventriloquism effect suggests that there is important information about sound location encoded in the synchrony between the audio and video signals. In spite of this evidence, audiovisual synchrony is rarely used as a source of information in computer vision tasks. In this paper we explore the use of audio visual synchrony to locate sound sources. We developed a system that searches for regions of the visual landscape that correlate highly with the acoustic signals and tags them as likely to contain an acoustic source. We discuss our experience implementing the system, present results on a speaker localization task and discuss potential applications of the approach.

Predictive Sequence Learning in Recurrent Neocortical Circuits

Rajesh Rao, Terrence J. Sejnowski

Neocortical circuits are dominated by massive excitatory feedback: more than eighty percent of the synapses made by excitatory cortical neurons are onto other excitatory cortical neurons. Why is there such massive recurrent excitation in the neocortex and what is its role in cortical computation? Recent neurophysiological experiments have shown that the plasticity of recurrent neocortical synapses is governed by a temporally asymmetric Hebbian learning rule. We describe how such a rule may allow the cortex to modify recurrent synapses for prediction of input sequences. The goal is to predict the next cortical input from the recent past based on previous experience of similar input sequences. We show that a temporal difference learning rule for prediction used in conjunction with dendritic back-propagating action potentials

ntials reproduces the temporally asymmetric Hebbian plasticity observed physiologically. Biophysical simulations demonstrate that a network of cortical neurons can learn to predict moving stimuli and develop direction selective responses as a consequence of learning. The space-time response properties of model neurons are shown to be similar to those of direction selective cells in alert monkey VI.

Differentiating Functions of the Jacobian with Respect to the Weights

Gary Flake, Barak Pearlmutter

For many problems, the correct behavior of a model depends not only on its input-output mapping but also on properties of its Jacobian matrix, the matrix of partial derivatives of the model's outputs with respect to its inputs.

We introduce the J-prop algorithm, an efficient general method for computing the exact partial derivatives of a variety of simple functions of the Jacobian of a model with respect to its free parameters. The algorithm applies to any parameterized feedforward model, including nonlinear regression, multilayer perceptrons, and radial basis function networks.

Effects of Spatial and Temporal Contiguity on the Acquisition of Spatial Information

Thea Ghiselli-Crippa, Paul Munro

Spatial information comes in two forms: direct spatial information (for example, retinal position) and indirect temporal contiguity information, since objects encountered sequentially are in general spatially close. The acquisition of spatial information by a neural network is investigated here. Given a spatial layout of several objects, networks are trained on a prediction task. Networks using temporal sequences with no direct spatial information are found to develop internal representations that show distances correlated with distances in the external layout. The influence of spatial information is analyzed by providing direct spatial information to the system during training that is either consistent with the layout or inconsistent with it. This approach allows examination of the relative contributions of spatial and temporal contiguity.

Agglomerative Information Bottleneck

Noam Slonim, Naftali Tishby

We introduce a novel distributional clustering algorithm that maximizes the mutual information per cluster between data and given categories. This algorithm can be considered as a bottom up hard version of the recently introduced "Information Bottleneck Method". The algorithm is compared with the top-down soft version of the information bottleneck method and a relationship between the hard and soft results is established. We demonstrate the algorithm on the 20 Newsgroups data set. For a subset of two newsgroups we achieve compression by 3 orders of magnitudes losing only 10% of the original mutual information.

Spike-based Learning Rules and Stabilization of Persistent Neural Activity

Xiaohui Xie, H. Sebastian Seung

We analyze the conditions under which synaptic learning rules based on action potential timing can be approximated by learning rules based on firing rates. In particular, we consider a form of plasticity in which synapses depress when a presynaptic spike is followed by a postsynaptic spike, and potentiate with the opposite temporal ordering. Such differential anti-Hebbian plasticity can be approximated under certain conditions by a learning rule that depends on the time derivative of the postsynaptic firing rate. Such a learning rule acts to stabilize persistent neural activity patterns in recurrent neural networks.

Nonlinear Discriminant Analysis Using Kernel Functions

Volker Roth, Volker Steinhage

Fishers linear discriminant analysis (LDA) is a classical multivariate technique both for dimension reduction and classification. The data vectors are transformed into a low dimensional subspace such that the class centroids are spread out as much as possible. In this subspace LDA works as a simple prototype classifier with linear decision boundaries.

However, in many applications the linear boundaries do not adequately separate the classes. We present a nonlinear generalization of discriminant analysis that uses the kernel trick of representing dot products by kernel functions. The presented algorithm allows a simple formulation of the EM-algorithm in terms of kernel functions which leads to a unique concept for unsupervised mixture analysis, supervised discriminant analysis and semi-supervised discriminant analysis with partially unlabelled observations in feature spaces.

On Input Selection with Reversible Jump Markov Chain Monte Carlo Sampling
Peter Sykacek

In this paper we will treat input selection for a radial basis function (RBF) like classifier within a Bayesian framework. We approximate the a-posteriori distribution over both model coefficients and input subsets by samples drawn with Gibbs updates and reversible jump moves. Using some public datasets, we compare the classification accuracy of the method with a conventional ARD scheme. These datasets are also used to infer the a-posteriori probabilities of different input subsets.

Uniqueness of the SVM Solution

Christopher J. C. Burges, David Crisp

We give necessary and sufficient conditions for uniqueness of the support vector solution for the problems of pattern recognition and regression estimation, for a general class of cost functions. We show that if the solution is not unique, all support vectors are necessarily at bound, and we give some simple examples of non-unique solutions. We note that uniqueness of the primal (dual) solution does not necessarily imply uniqueness of the dual (primal) solution. We show how to compute the threshold b when the solution is unique, but when all support vectors are at bound, in which case the usual method for determining b does not work.

The Parallel Problems Server: an Interactive Tool for Large Scale Machine Learning

Charles Isbell, Parry Husbands

Imagine that you wish to classify data consisting of tens of thousands of examples residing in a twenty thousand dimensional space. How can one apply standard machine learning algorithms? We describe the Parallel Problems Server (PPServer) and MATLAB*P. In tandem they allow users of networked computers to work transparently on large data sets from within Matlab.

This work is motivated by the desire to bring the many benefits of scientific computing algorithms and computational power to machine learning researchers. We demonstrate the usefulness of the system on a number of tasks. For example, we perform independent components analysis on very large text corpora consisting of tens of thousands of documents, making minimal changes to the original Bell and Sejnowski Matlab source (Bell and Sejnowski, 1995). Applying ML techniques to data previously beyond their reach leads to interesting analyses of both data and algorithms.

Search for Information Bearing Components in Speech

Howard Yang, Hynek Hermansky

In this paper, we use mutual information to characterize the distributions of phonetic and speaker/channel information in a time-frequency space. The mutual information (MI) between the phonetic label and one feature, and the joint mutual information (JMI) between the phonetic label and two or three features are estimated. The Miller's bias formula

as for entropy and mutual information estimates are extended to include higher order terms. The MI and the JMI for speaker/channel recognition are also estimated. The results are complementary to those for phonetic classification. Our results show how the phonetic information is locally spread and how the speaker/channel information is globally spread in time and frequency.

An Oscillatory Correlation Framework for Computational Auditory Scene Analysis
Guy Brown, DeLiang Wang

A neural model is described which uses oscillatory correlation to segregate speech from interfering sound sources. The core of the model is a two-layer neural oscillator network. A sound stream is represented by a synchronized population of oscillators, and different streams are represented by desynchronized oscillator populations. The model has been evaluated using a corpus of speech mixed with interfering sounds, and produces an improvement in signal-to-noise ratio for every mixture.

Manifold Stochastic Dynamics for Bayesian Learning
Mark Zlochin, Yoram Baram

We propose a new Markov Chain Monte Carlo algorithm which is a generalization of the stochastic dynamics method. The algorithm performs exploration of the state space using its intrinsic geometric structure, facilitating efficient sampling of complex distributions. Applied to Bayesian learning in neural networks, our algorithm was found to perform at least as well as the best state-of-the-art method while consuming considerably less time.

Bifurcation Analysis of a Silicon Neuron
Girish Patel, Gennady Cymbalyuk, Ronald Calabrese, Stephen DeWeerth

We have developed a VLSI silicon neuron and a corresponding mathematical model that is a two state-variable system. We describe the circuit implementation and compare the behaviors observed in the silicon neuron and the mathematical model. We also perform bifurcation analysis of the mathematical model by varying the externally applied current and show that the behaviors exhibited by the silicon neuron under corresponding conditions are in good agreement to those predicted by the bifurcation analysis.

Speech Modelling Using Subspace and EM Techniques

Gavin Smith, João de Freitas, Tony Robinson, Mahesan Niranjan

The speech waveform can be modelled as a piecewise-stationary linear stochastic state space system, and its parameters can be estimated using an expectation-maximisation (EM) algorithm. One problem is the initialization of the EM algorithm. Standard initialisation schemes can lead to poor formant trajectories. But these trajectories however are important for vowel intelligibility. The aim of this paper is to investigate the suitability of subspace identification methods to initialise EM. The paper compares the subspace state space system identification (4SID) method with the EM algorithm.

The 4SID and EM methods are similar in that they both estimate a state sequence (but using Kalman filters and Kalman smoothers respectively), and then estimate parameters (but using least-squares and maximum likelihood respectively). The similarity of 4SID and EM motivates the use of 4SID to initialise EM. Also, 4SID is non-iterative and requires no initialisation, whereas EM is iterative and requires initialisation. However 4SID is sub-optimal compared to EM in a probabilistic sense. During experiments on real speech, 4SID methods compare favourably with conventional initialisation techniques.

They produce smoother formant trajectories, have greater frequency resolution, and produce higher likelihoods.

Optimal Kernel Shapes for Local Linear Regression

Dirk Ormoneit, Trevor Hastie

Local linear regression performs very well in many low-dimensional forecasting problems. In high-dimensional spaces, its performance typically decays due to the well-known "curse-of-dimensionality". A possible way to approach this problem is by varying the "shape" of the weighting kernel. In this work we suggest a new, data-driven method to estimating the optimal kernel shape. Experiments using an artificially generated data set and data from the UC Irvine repository show the benefits of kernel shaping.

Robust Neural Network Regression for Offline and Online Learning

Thomas Briegel, Volker Tresp

We replace the commonly used Gaussian noise model in nonlinear regression by a more flexible noise model based on the Student-t distribution. The degrees of freedom of the t-distribution can be chosen such that as special cases either the Gaussian distribution or the Cauchy distribution are realized. The latter is commonly used in robust regression.

Since the t-distribution can be interpreted as being an infinite mixture of Gaussians, parameters and hyperparameters such as the degrees of freedom of the t-distribution can be learned from the data based on an EM-learning algorithm. We show that modeling using the t-distribution leads to improved predictors on real world data sets. In particular, if outliers are present, the t-distribution is superior to the Gaussian noise model.

In effect, by adapting the degrees of freedom, the system can "learn" to distinguish between outliers and non-outliers. Especially for online learning tasks, one is interested in avoiding inappropriate weight changes due to measurement outliers to maintain stable online learning capability. We show experimentally that using the t-distribution as a noise model leads to stable online learning algorithms and outperforms state-of-the-art online learning methods like the extended Kalman filter algorithm.

Neural Network Based Model Predictive Control

Stephen Piche, James Keeler, Greg Martin, Gene Boe, Doug Johnson, Mark Gerules
Model Predictive Control (MPC), a control algorithm which uses an optimizer to solve for the optimal control moves over a future time horizon based upon a model of the process, has become a standard control technique in the process industries over the past two decades. In most industrial applications, a linear dynamic model developed using empirical data is used even though the process itself is often nonlinear. Linear models have been used because of the difficulty in developing a generic nonlinear model from empirical data and the computational expense often involved in using nonlinear models. In this paper, we present a generic neural network based technique for developing nonlinear dynamic models from empirical data and show that these models can be efficiently used in a model predictive control framework. This nonlinear MPC based approach has been successfully implemented in a number of industrial applications in the refining, petrochemical, paper and food industries. Performance of the controller on a nonlinear industrial process, a polyethylene reactor, is presented.

Coastal Navigation with Mobile Robots

Nicholas Roy, Sebastian Thrun

The problem that we address in this paper is how a mobile robot can plan in order to arrive at its goal with minimum uncertainty. Traditional motion planning algorithms often assume that a mobile robot can track its position reliably, however, in real world situations, reliable localization may not always be feasible. Partially Observable Markov Decision Processes (POMDPs) provide one way to maximize the certainty of reaching the goal state, but at the cost of computational intractability for large state spaces.

The method we propose explicitly models the uncertainty of the robot's position as a state variable, and generates trajectories through the augmented pose-uncertainty space. By minimizing the positional uncertainty

at the goal, the robot reduces the likelihood it becomes lost. We demonstrate experimentally that coastal navigation reduces the uncertainty at the goal, especially with degraded localization.

Boosting with Multi-Way Branching in Decision Trees

Yishay Mansour, David McAllester

It is known that decision tree learning can be viewed as a form of boosting. However, existing boosting theorems for decision tree learning allow only binary-branching trees and the generalization to multi-branching trees is not immediate. Practical decision tree algorithms, such as CART and C4.5, implement a trade-off between the number of branches and the improvement in tree quality as measured by an index function. Here we give a boosting justification for a particular quantitative trade-off curve. Our main theorem states, in essence, that if we require an improvement proportional to the log of the number of branches then top-down greedy construction of decision trees remains an effective boosting algorithm.

Neural Representation of Multi-Dimensional Stimuli

Christian Euriich, Stefan Wilke, Helmut Schwegler

Spatial information comes in two forms: direct spatial information (for example, retinal position) and indirect temporal contiguity information, since objects encountered sequentially are in general spatially close. The acquisition of spatial information by a neural network is investigated here. Given a spatial layout of several objects, networks are trained on a prediction task. Networks using temporal sequences with no direct spatial information are found to develop internal representations that show distances correlated with distances in the external layout. The influence of spatial information is analyzed by providing direct spatial information to the system during training that is either consistent with the layout or inconsistent with it. This approach allows examination of the relative contributions of spatial and temporal contiguity.

Model Selection for Support Vector Machines

Olivier Chapelle, Vladimir Vapnik

New functionals for parameter (model) selection of Support Vector Machines are introduced based on the concepts of the span of support vectors and rescaling of the feature space. It is shown that using these functionals, one can both predict the best choice of parameters of the model and the relative quality of performance for any value of parameter.

Memory Capacity of Linear vs. Nonlinear Models of Dendritic Integration

Panayiota Poirazi, Bartlett Mel

Previous biophysical modeling work showed that nonlinear interactions among nearby synapses located on active dendritic trees can provide a large boost in the memory capacity of a cell (Mel, 1992a, 1992b). The aim of our present work is to quantify this boost by estimating the capacity of (1) a neuron model with passive dendritic integration where inputs are combined linearly across the entire cell followed by a single global threshold, and (2) an active dendrite model in which a threshold is applied separately to the output of each branch, and the branch subtotals are combined linearly early. We focus here on the limiting case of binary-valued synaptic weights, and derive expressions which measure model capacity by estimating the number of distinct input-output functions available to both neuron types. We show that (1) the application of a fixed nonlinearity to each dendritic compartment substantially increases the model's flexibility, (2) for a neuron of realistic size, the capacity of the nonlinear cell can exceed that of the same-sized linear cell by more than an order of magnitude, and (3) the largest capacity boost occurs for cells with a relatively large number of dendritic subunits of relatively small size. We validated the analysis

by empirically measuring memory capacity with randomized two-class classification problems, where a stochastic delta rule was used to train both linear and nonlinear models. We found that large capacity boosts predicted for the nonlinear dendritic model were readily achieved in practice.

State Abstraction in MAXQ Hierarchical Reinforcement Learning

Thomas Dietterich

Many researchers have explored methods for hierarchical reinforcement learning (RL) with temporal abstractions, in which abstract actions are defined that can perform many primitive actions before terminating. However, little is known about learning with state abstractions, in which aspects of the state space are ignored. In previous work, we developed the MAXQ method for hierarchical RL. In this paper, we define five conditions under which state abstraction can be combined with the MAXQ value function decomposition. We prove that the MAXQ-Q learning algorithm converges under these conditions and show experimentally that state abstraction is important for the successful application of MAXQ-Q learning.

Modeling High-Dimensional Discrete Data with Multi-Layer Neural Networks

Yoshua Bengio, Samy Bengio

The curse of dimensionality is severe when modeling high-dimensional discrete data: the number of possible combinations of the variables explodes exponentially. In this paper we propose a new architecture for modeling high-dimensional data that requires resources (parameters and computations) that grow only at most as the square of the number of variables, using a multi-layer neural network to represent the joint distribution of the variables as the product of conditional distributions. The neural network can be interpreted as a graphical model without hidden variables, but in which the conditional distributions are tied through the hidden units. The connectivity of the neural network can be pruned by using dependency tests between the variables. Experiments on modeling the distribution of several discrete data sets show statistically significant improvements over other methods such as naive Bayes and comparable Bayesian networks, and show that significant improvements can be obtained by pruning the network.

A Generative Model for Attractor Dynamics

Richard Zemel, Michael C. Mozer

Attractor networks, which map an input space to a discrete output space, are useful for pattern completion. However, designing a net to have a given set of attractors is notoriously tricky; training procedures are CPU intensive and often produce spurious attractors and ill-conditioned attractor basins. These difficulties occur because each connection in the network participates in the encoding of multiple attractors. We describe an alternative formulation of attractor networks in which the encoding of knowledge is local, not distributed. Although localist attractor networks have similar dynamics to their distributed counterparts, they are much easier to work with and interpret. We propose a statistical formulation of localist attractor network dynamics, which yields a convergence proof and a mathematical interpretation of model parameters.

An Environment Model for Nonstationary Reinforcement Learning

Samuel Choi, Dit-Yan Yeung, Nevin Zhang

Reinforcement learning in nonstationary environments is generally regarded as an important and yet difficult problem. This paper partially addresses the problem by formalizing a subclass of nonstationary environments. The environment model, called hidden-mode Markov decision process (HM-MDP), assumes that environmental changes are always confined to a small number of hidden modes. A mode basically indexes a Markov decision process (MDP) and evolves with time according to a Markov chain. While HM-MDP is a special case of partially observable Markov decision processes (P

OMDP), modeling an HM-MDP environment via the more general POMDP model unnecessarily increases the problem complexity. A variant of the Baum-Welch algorithm is developed for model learning requiring less data and time.

Predictive Approaches for Choosing Hyperparameters in Gaussian Processes

S. Sundararajan, S. Keerthi

Gaussian Processes are powerful regression models specified by parameterized mean and covariance functions. Standard approaches to estimate these parameters (known by the name Hyperparameters) are Maximum Likelihood (ML) and Maximum A Posterior (MAP) approaches. In this paper, we propose and investigate predictive approaches, namely, maximization of Geisser's Surrogate Predictive Probability (GPP) and minimization of mean square error with respect to GPP (referred to as Geisser's Predictive mean square Error (GPE)) to estimate the hyperparameters. We also derive results for the standard Cross-Validation (CV) error and make a comparison. These approaches are tested on a number of problems and experimental results show that these approaches are strongly competitive to existing approaches.

Acquisition in Autoshaping

Sham Kakade, Peter Dayan

Quantitative data on the speed with which animals acquire behavioral responses during classical conditioning experiments should provide strong constraints on models of learning. However, most models have simply ignored these data; the few that have attempted to address them have failed by at least an order of magnitude. We discuss key data on the speed of acquisition, and show how to account for them using a statistically sound model of learning, in which differential reliabilities of stimuli play a crucial role.

Transductive Inference for Estimating Values of Functions

Olivier Chapelle, Vladimir Vapnik, Jason Weston

We introduce an algorithm for estimating the values of a function at a set of test points x_{e+1}, \dots, x_{l+m} given a set of training points $(x_i, y_i), \dots, (x_e, y_e)$ without estimating (as an intermediate step) the regression function. We demonstrate that this direct (transductive) way for estimating values of the regression (or classification in pattern recognition) can be more accurate than the traditional one based on two steps, first estimating the function and then calculating the values of this function at the points of interest.

Effective Learning Requires Neuronal Remodeling of Hebbian Synapses

Gal Chechik, Isaac Meilijson, Eytan Ruppin

This paper revisits the classical neuroscience paradigm of Hebbian learning. We find that a necessary requirement for effective associative memory learning is that the efficacies of the incoming synapses should be uncorrelated. This requirement is difficult to achieve in a robust manner by Hebbian synaptic learning, since it depends on network level information. Effective learning can yet be obtained by a neuronal process that maintains a zero sum of the incoming synaptic efficacies. This normalization drastically improves the memory capacity of associative networks, from an essentially bounded capacity to one that linearly scales with the network's size. It also enables the effective storage of patterns with heterogeneous coding levels in a single network. Such neuronal normalization can be successfully carried out by activity-dependent homeostasis of the neuron's synaptic efficacies, which was recently observed in cortical tissue. Thus, our findings strongly suggest that effective associative learning with Hebbian synapses alone is biologically implausible and that Hebbian synapses must be continuously remodeled by neuronally-driven regulatory processes in the brain.

The Relevance Vector Machine

Michael Tipping

The support vector machine (SVM) is a state-of-the-art technique for regression and classification, combining excellent generalisation properties with a sparse kernel representation. However, it does suffer from a number of disadvantages, notably the absence of probabilistic outputs, the requirement to estimate a trade-off parameter and the need to utilise 'Mercer' kernel functions. In this paper we introduce the Relevance Vector Machine (RVM), a Bayesian treatment of a generalised linear model of identical functional form to the SVM. The RVM suffers from none of the above disadvantages, and examples demonstrate that for comparable generalisation performance, the RVM requires dramatically fewer kernel functions.

Dual Estimation and the Unscented Transformation

Eric Wan, Rudolph van der Merwe, Alex Nelson

Dual estimation refers to the problem of simultaneously estimating the state of a dynamic system and the model which gives rise to the dynamics.

Algorithms include expectation-maximization (EM), dual Kalman filtering, and joint Kalman methods. These methods have recently been explored in the context of nonlinear modeling, where a neural network is used as the functional form of the unknown model. Typically, an extended Kalman filter (EKF) or smoother is used for the part of the algorithm that estimates the clean state given the current estimated model. An EKF may also be used to estimate the weights of the network. This paper points out the flaws in using the EKF, and proposes an improvement based on a new approach called the unscented transformation (UT) [3]. A substantial performance gain is achieved with the same order of computational complexity as that of the standard EKF. The approach is illustrated on several dual estimation methods.

Learning Statistically Neutral Tasks without Expert Guidance

Ton Weijters, Antal van den Bosch, Eric Postma

Eric Postma

Gaussian Fields for Approximate Inference in Layered Sigmoid Belief Networks

David Barber, Peter Sollich

Local "belief propagation" rules of the sort proposed by Pearl [15] are guaranteed to converge to the correct posterior probabilities in singly connected graphical models. Recently, a number of researchers have empirically demonstrated good performance of "loopy belief propagation" using these same rules on graphs with loops. Perhaps the most dramatic instance is the near Shannon-limit performance of "Turbo codes", whose decoding algorithm is equivalent to loopy belief propagation. Except for the case of graphs with a single loop, there has been little theoretical understanding of the performance of loopy propagation. Here we analyze belief propagation in networks with arbitrary topologies when the nodes in the graph describe jointly Gaussian random variables. We give an analytical formula relating the true posterior probabilities with those calculated using loopy propagation. We give sufficient conditions for convergence and show that when belief propagation converges it gives the correct posterior means for all graph topologies, not just networks with a single loop. The related "max-product" belief propagation algorithm finds the maximum posterior probability estimate for singly connected networks. We show that, even for non-Gaussian probability distributions, the convergence points of the max-product algorithm in loopy networks are maxima over a particular large local neighborhood of the posterior probability. These results help clarify the empirical performance results and motivate using the powerful belief propagation algorithm in a broader class of networks.

A Multi-class Linear Learning Algorithm Related to Winnow

Chris Mesterharm

In this paper, we present Committee, a new multi-class learning algorithm related to the Winnow family of algorithms. Committee is an algorithm for combining the predictions of a set of sub-experts in the on-line mistake-bounded model of learning. A sub-expert is a special type of attribute that predicts with a distribution over a finite number of classes. Committee learns a linear function of sub-experts and uses this function to make class predictions. We provide bounds for Committee that show it performs well when the target can be represented by a few relevant sub-experts. We also show how Committee can be used to solve more traditional problems composed of attributes. This leads to a natural extension that learns on multi-class problems that contain both traditional attributes and sub-experts.

Churn Reduction in the Wireless Industry

Michael C. Mozer, Richard Wolniewicz, David Grimes, Eric Johnson, Howard Kaushansky

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.

Recognizing Evoked Potentials in a Virtual Environment

Jessica Bayliss, Dana Ballard

Virtual reality (VR) provides immersive and controllable experimental environments. It expands the bounds of possible evoked potential (EP) experiments by providing complex, dynamic environments in order to study cognition without sacrificing environmental control. VR also serves as a safe dynamic testbed for brain-computer interface (BCI) research. However, there has been some concern about detecting EP signals in a complex VR environment. This paper shows that EPs exist at red, green, and yellow stop lights in a virtual driving environment. Experimental results show the existence of the P3 EP at "go" and "stop" lights and the contingent negative variation (CNV) EP at "slow down" lights. In order to test the feasibility of on-line recognition in VR, we looked at recognizing the P3 EP at red stop lights and the absence of this signal at yellow slow down lights. Recognition results show that the P3 may successfully be used to control the brakes of a VR car at stop lights.
