

## Temporal Coherence, Natural Image Sequences, and the Visual Cortex

Jarmo Hurri, Aapo Hyvärinen

We show that two important properties of the primary visual cortex emerge when the principle of temporal coherence is applied to natural image sequences. The properties are simple-cell-like receptive fields and complex-cell-like pooling of simple cell outputs, which emerge when we apply two different approaches to temporal coherence. In the first approach we extract receptive fields whose outputs are as temporally coherent as possible. This approach yields simple-cell-like receptive fields (oriented, localized, multiscale). Thus, temporal coherence is an alternative to sparse coding in modeling the emergence of simple cell receptive fields. The second approach is based on a two-layer statistical generative model of natural image sequences. In addition to modeling the temporal coherence of individual simple cells, this model includes inter-cell temporal dependencies. Estimation of this model from natural data yields both simple-cell-like receptive fields, and complex-cell-like pooling of simple cell outputs. In this completely unsupervised learning, both layers of the generative model are estimated simultaneously from scratch. This is a significant improvement on earlier statistical models of early vision, where only one layer has been learned, and others have been fixed a priori.

\*\*\*\*\*

## Nonparametric Representation of Policies and Value Functions: A Trajectory-Based Approach

Christopher Atkeson, Jun Morimoto

A longstanding goal of reinforcement learning is to develop non-parametric representations of policies and value functions that support rapid learning without suffering from interference or the curse of dimensionality. We have developed a trajectory-based approach, in which policies and value functions are represented nonparametrically along trajectories. These trajectories, policies, and value functions are updated as the value function becomes more accurate or as a model of the task is updated. We have applied this approach to periodic tasks such as hopping and walking, which required handling discount factors and discontinuities in the task dynamics, and using function approximation to represent value functions at discontinuities. We also describe extensions of the approach to make the policies more robust to modeling error and sensor noise.

\*\*\*\*\*

## A Differential Semantics for Jointree Algorithms

James Park, Adnan Darwiche

A new approach to inference in belief networks has been recently proposed, which is based on an algebraic representation of belief networks using multi{linear functions. According to this approach, the key computational question is that of representing multi{linear functions compactly, since inference reduces to a simple process of evaluating and differentiating such functions. We show here that mainstream inference algorithms based on jointrees are a special case of this approach in a very precise sense. We use this result to prove new properties of jointree algorithms, and then discuss some of its practical and theoretical implications.

\*\*\*\*\*

## "Name That Song!" A Probabilistic Approach to Querying on Music and Text

Brochu Eric, Nando de Freitas

We present a novel, flexible statistical approach for modelling music and text jointly. The approach is based on multi-modal mixture models and maximum a posteriori estimation using EM. The learned models can be used to browse databases with documents containing music and text, to search for music using queries consisting of music and text (lyrics and other contextual information), to annotate text documents with music, and to automatically recommend or identify similar songs.

\*\*\*\*\*

## Automatic Derivation of Statistical Algorithms: The EM Family and Beyond

Bernd Fischer, Johann Schumann, Wray Buntine, Alexander Gray

Machine learning has reached a point where many probabilistic methods can be understood as variations, extensions and combinations of a much smaller set of ab

stract themes, e.g., as different instances of the EM algorithm. This enables the systematic derivation of algorithms customized for different models. Here, we describe the AUTO BAYES system which takes a high-level statistical model specification, uses powerful symbolic techniques based on schema-based program synthesis and computer algebra to derive an efficient specialized algorithm for learning that model, and generates executable code implementing that algorithm. This capability is far beyond that of code collections such as Matlab toolboxes or even tools for model-independent optimization such as BUGS for Gibbs sampling: complex new algorithms can be generated without new programming, algorithms can be highly specialized and tightly crafted for the exact structure of the model and data, and efficient and commented code can be generated for different languages or systems. We present automatically-derived algorithms ranging from closed-form solutions of Bayesian textbook problems to recently-proposed EM algorithms for clustering, regression, and a multinomial form of PCA.

\*\*\*\*\*

#### Going Metric: Denoising Pairwise Data

Volker Roth, Julian Laub, Klaus-Robert Müller, Joachim Buhmann

Pairwise data in empirical sciences typically violate metricity, either due to noise or due to fallible estimates, and therefore are hard to analyze by conventional machine learning technology. In this paper we therefore study ways to work around this problem. First, we present an alternative embedding to multi-dimensional scaling (MDS) that allows us to apply a variety of classical machine learning and signal processing algorithms. The class of pairwise grouping algorithms which share the shift-invariance property is statistically invariant under this embedding procedure, leading to identical assignments of objects to clusters. Based on this new vectorial representation, denoising methods are applied in a second step. Both steps provide a theoretically well controlled setup to translate from pairwise data to the respective denoised metric representation. We demonstrate the practical usefulness of our theoretical reasoning by discovering structure in protein sequence data bases, visibly improving performance upon existing automatic methods.

\*\*\*\*\*

#### String Kernels, Fisher Kernels and Finite State Automata

Craig Saunders, Alexei Vinokourov, John Shawe-taylor

In this paper we show how the generation of documents can be thought of as a k-stage Markov process, which leads to a Fisher kernel from which the n-gram and string kernels can be re-constructed. The Fisher kernel view gives a more flexible insight into the string kernel and suggests how it can be parametrised in a way that reflects the statistics of the training corpus. Furthermore, the probabilistic modelling approach suggests extending the Markov process to consider sub-sequences of varying length, rather than the standard fixed-length approach used in the string kernel. We give a procedure for determining which sub-sequences are informative features and hence generate a Finite State Machine model, which can again be used to obtain a Fisher kernel. By adjusting the parametrisation we can also influence the weighting received by the features. In this way we are able to obtain a logarithmic weighting in a Fisher kernel. Finally, experiments are reported comparing the different kernels using the standard Bag of Words kernel as a baseline.

\*\*\*\*\*

#### Extracting Relevant Structures with Side Information

Gal Chechik, Naftali Tishby

The problem of extracting the relevant aspects of data, in face of multiple conflicting structures, is inherent to modeling of complex data. Extracting structure in one random variable that is relevant for another variable has been principally addressed recently via the information bottleneck method [15]. However, such auxiliary variables often contain more information than is actually required

due to structures that are irrelevant for the task. In many other cases it is in fact easier to specify what is irrelevant than what is, for the task at hand. Identifying the relevant structures, however, can thus be considerably improved by also minimizing the information about another, irrelevant, variable. In this paper we give a general formulation of this problem and derive its formal, as well as algorithmic, solution. Its operation is demonstrated in a synthetic example and in two real world problems in the context of text categorization and face images. While the original information bottleneck problem is related to rate distortion theory, with the distortion measure replaced by the relevant information, extracting relevant features while removing irrelevant ones is related to rate distortion with side information.

\*\*\*\*\*

#### Classifying Patterns of Visual Motion - a Neuromorphic Approach

Jakob Heinzle, Alan A. Stocker

We report a system that classifies and can learn to classify patterns of visual motion on-line. The complete system is described by the dynamics of its physical network architectures. The combination of the following properties makes the system novel: Firstly, the front-end of the system consists of an aVLSI optical flow chip that collectively computes 2-D global visual motion in real-time [1]. Secondly, the complexity of the classification task is significantly reduced by mapping the continuous motion trajectories to sequences of 'motion events'. And thirdly, all the network structures are simple and with the exception of the optical flow chip based on a Winner-Take-All (WTA) architecture. We demonstrate the application of the proposed generic system for a contactless man-machine interface that allows to write letters by visual motion. Regarding the low complexity of the system, its robustness and the already existing front-end, a complete aVLSI system-on-chip implementation is realistic, allowing various applications in mobile electronic devices.

\*\*\*\*\*

#### How the Poverty of the Stimulus Solves the Poverty of the Stimulus

Willem Zuidema

Language acquisition is a special kind of learning problem because the outcome of learning of one generation is the input for the next. That makes it possible for languages to adapt to the particularities of the learner. In this paper, I show that this type of language change has important consequences for models of the evolution and acquisition of syntax.

\*\*\*\*\*

#### Fast Transformation-Invariant Factor Analysis

Anitha Kannan, Nebojsa Jojic, Brendan Frey

Dimensionality reduction techniques such as principal component analysis and factor analysis are used to discover a linear mapping between high dimensional data samples and points in a lower dimensional subspace. In [6], Jojic and Frey introduced mixture of transformation-invariant component analyzers (MTCA) that can account for global transformations such as translations and rotations, perform clustering and learn local appearance deformations by dimensionality reduction. However, due to enormous computational requirements of the EM algorithm for learning the model,  $O(d)$  is the dimensionality of a data sample, MTCA was not practical for most applications. In this paper, we demonstrate how fast Fourier transforms can reduce the computation to the order of  $O(1)$ . With this speedup, we show the effectiveness of MTCA derived in various applications - tracking, video textures, clustering video sequences, object recognition, and object detection in images.

\*\*\*\*\*

#### Spikernels: Embedding Spiking Neurons in Inner-Product Spaces

Lavi Shpigelman, Yoram Singer, Rony Paz, Eilon Vaadia

Inner-product operators, often referred to as kernels in statistical learning, define a mapping from some input space into a feature space. The focus of this paper is the construction of biologically-motivated kernels for cortical activities. The kernels we derive, termed Spikernels, map spike count sequences into an abstract vector space in which we can perform various prediction tasks. We dis

cuss in detail the derivation of Spikernels and describe an efficient algorithm for computing their value on any two sequences of neural population spike counts. We demonstrate the merits of our modeling approach using the Spikernel and various standard kernels for the task of predicting hand movement velocities from cortical recordings. In all of our experiments all the kernels we tested out perform the standard scalar product used in regression with the Spikernel consistently achieving the best performance.

\*\*\*\*\*

#### Graph-Driven Feature Extraction From Microarray Data Using Diffusion Kernels and Kernel CCA

Jean-philippe Vert, Minoru Kanehisa

We present an algorithm to extract features from high-dimensional gene expression profiles, based on the knowledge of a graph which links together genes known to participate to successive reactions in metabolic pathways. Motivated by the intuition that biologically relevant features are likely to exhibit smoothness with respect to the graph topology, the algorithm involves encoding the graph and the set of expression profiles into kernel functions, and performing a generalized form of canonical correlation analysis in the corresponding reproducible kernel Hilbert spaces. Function prediction experiments for the genes of the yeast *S. Cerevisiae* validate this approach by showing a consistent increase in performance when a state-of-the-art classifier uses the vector of features instead of the original expression profile to predict the functional class of a gene.

\*\*\*\*\*

#### Spike Timing-Dependent Plasticity in the Address Domain

R. Vogelstein, Francesco Tenore, Ralf Philipp, Miriam Adlerstein, David Goldberg, Gert Cauwenberghs

Address-event representation (AER), originally proposed as a means to communicate sparse neural events between neuromorphic chips, has proven efficient in implementing large-scale networks with arbitrary, configurable synaptic connectivity. In this work, we further extend the functionality of AER to implement arbitrary, configurable synaptic plasticity in the address domain. As proof of concept, we implement a biologically inspired form of spike timing-dependent plasticity (STDP) based on relative timing of events in an AER framework. Experimental results from an analog VLSI integrate-and-fire network demonstrate address domain learning in a task that requires neurons to group correlated inputs.

\*\*\*\*\*

#### Efficient Learning Equilibrium

Ronen Brafman, Moshe Tennenholtz

We introduce efficient learning equilibrium (ELE), a normative approach to learning in non cooperative settings. In ELE, the learning algorithms themselves are required to be in equilibrium. In addition, the learning algorithms arrive at a desired value after polynomial time, and deviations from a prescribed ELE become irrational after polynomial time. We prove the existence of an ELE in the perfect monitoring setting, where the desired value is the expected payoff in a Nash equilibrium. We also show that an ELE does not always exist in the imperfect monitoring case. Yet, it exists in the special case of common-interest games. Finally, we extend our results to general stochastic games.

\*\*\*\*\*

#### Using Tarjan's Red Rule for Fast Dependency Tree Construction

Dan Pelleg, Andrew Moore

We focus on the problem of efficient learning of dependency trees. It is well-known that given the pairwise mutual information coefficients, a minimum-weight spanning tree algorithm solves this problem exactly and in polynomial time. However, for large data-sets it is the construction of the correlation matrix that dominates the running time. We have developed a new spanning-tree algorithm which is capable of exploiting partial knowledge about edge weights. The partial knowledge we maintain is a probabilistic confidence interval on the coefficients, which we derive by examining just a small sample of the data. The algorithm is able to flag the need to shrink an interval, which translates to inspection of more

data for the particular attribute pair. Experimental results show running time that is near-constant in the number of records, without significant loss in accuracy of the generated trees. Interestingly, our spanning-tree algorithm is based solely on Tarjan's red-edge rule, which is generally considered a guaranteed recipe for bad performance.

\*\*\*\*\*

#### Approximate Inference and Protein-Folding

Chen Yanover, Yair Weiss

Side-chain prediction is an important subtask in the protein-folding problem. We show that finding a minimal energy side-chain configuration is equivalent to performing inference in an undirected graphical model. The graphical model is relatively sparse yet has many cycles. We used this equivalence to assess the performance of approximate inference algorithms in a real-world setting. Specifically we compared belief propagation (BP), generalized BP (GBP) and naive mean field (MF). In cases where exact inference was possible, max-product BP always found the global minimum of the energy (except in few cases where it failed to converge), while other approximation algorithms of similar complexity did not. In the full protein data set, max-product BP always found a lower energy configuration than the other algorithms, including a widely used protein-folding software (SCWRL).

\*\*\*\*\*

#### Learning a Forward Model of a Reflex

Bernd Porr, Florentin Wörgötter

We develop a systems theoretical treatment of a behavioural system that interacts with its environment in a closed loop situation such that its motor actions influence its sensor inputs. The simplest form of a feedback is a reflex. Reflexes occur always "too late"; i.e., only after a (unpleasant, painful, dangerous) reflex-eliciting sensor event has occurred. This defines an objective problem which can be solved if another sensor input exists which can predict the primary reflex and can generate an earlier reaction. In contrast to previous approaches, our linear learning algorithm allows for an analytical proof that this system learns to apply feed-forward control with the result that slow feedback loops are replaced by their equivalent feed-forward controller creating a forward model. In other words, learning turns the reactive system into a pro-active system. By means of a robot implementation we demonstrate the applicability of the theoretical results which can be used in a variety of different areas in physics and engineering.

\*\*\*\*\*

#### Information Regularization with Partially Labeled Data

Martin Szummer, Tommi Jaakkola

Classification with partially labeled data requires using a large number of unlabeled examples (or an estimated marginal  $P(x)$ ), to further constrain the conditional  $P(y|x)$  beyond a few available labeled examples. We formulate a regularization approach to linking the marginal and the conditional in a general way. The regularization penalty measures the information that is implied about the labels over covering regions. No parametric assumptions are required and the approach remains tractable even for continuous marginal densities  $P(x)$ . We develop algorithms for solving the regularization problem for finite covers, establish a limiting differential equation, and exemplify the behavior of the new regularization approach in simple cases.

\*\*\*\*\*

#### Intrinsic Dimension Estimation Using Packing Numbers

Balázs Kégl

We propose a new algorithm to estimate the intrinsic dimension of data sets. The method is based on geometric properties of the data and requires neither parametric assumptions on the data generating model nor input parameters to set. The method is compared to a similar, widely-used algorithm from the same family of geometric techniques. Experiments show that our method is more robust in terms of the data generating distribution and more reliable in the presence of noise

.  
\*\*\*\*\*

#### Mismatch String Kernels for SVM Protein Classification

Eleazar Eskin, Jason Weston, William Noble, Christina Leslie

We introduce a class of string kernels, called mismatch kernels, for use with support vector machines (SVMs) in a discriminative approach to the protein classification problem. These kernels measure sequence similarity based on shared occurrences of  $\ell$ -length subsequences, counted with up to  $m$  mismatches, and do not rely on any generative model for the positive training sequences. We compute the kernels efficiently using a mismatch tree data structure and report experiments on a benchmark SCOP dataset, where we show that the mismatch kernel used with an SVM classifier performs as well as the Fisher kernel, the most successful method for remote homology detection, while achieving considerable computational savings

.  
\*\*\*\*\*

#### Adaptive Caching by Refetching

Robert B. Gramacy, Manfred K. K. Warmuth, Scott Brandt, Ismail Ari

We are constructing caching policies that have 13-20% lower miss rates than the best of twelve baseline policies over a large variety of request streams. This represents an improvement of 49-63% over Least Recently Used, the most commonly implemented policy. We achieve this not by designing a specific new policy but by using on-line Machine Learning algorithms to dynamically shift between the standard policies based on their observed miss rates. A thorough experimental evaluation of our techniques is given, as well as a discussion of what makes caching an interesting on-line learning problem.

\*\*\*\*\*

#### Timing and Partial Observability in the Dopamine System

Nathaniel Daw, Aaron C. Courville, David Touretzky

According to a series of influential models, dopamine (DA) neurons signal reward prediction error using a temporal-difference (TD) algorithm. We address a problem not convincingly solved in these accounts: how to maintain a representation of cues that predict delayed consequences. Our new model uses a TD rule grounded in partially observable semi-Markov processes, a formalism that captures two largely neglected features of DA experiments: hidden state and temporal variability. Previous models predicted rewards using a tapped delay line representation of sensory inputs; we replace this with a more active process of inference about the underlying state of the world. The DA system can then learn to map these inferred states to reward predictions using TD. The new model can explain previously vexing data on the responses of DA neurons in the face of temporal variability. By combining statistical model-based learning with a physiologically grounded TD theory, it also brings into contact with physiology some insights about behavior that had previously been confined to more abstract psychological models.

\*\*\*\*\*

#### Multiple Cause Vector Quantization

David Ross, Richard Zemel

We propose a model that can learn parts-based representations of high-dimensional data. Our key assumption is that the dimensions of the data can be separated into several disjoint subsets, or factors, which take on values independently of each other. We assume each factor has a small number of discrete states, and model it using a vector quantizer. The selected states of each factor represent the multiple causes of the input. Given a set of training examples, our model learns the association of data dimensions with factors, as well as the states of each VQ. Inference and learning are carried out efficiently via variational algorithms. We present applications of this model to problems in image decomposition, collaborative filtering, and text classification.

\*\*\*\*\*

#### Unsupervised Color Constancy

Kinh Tieu, Erik Miller

In [1] we introduced a linear statistical model of joint color changes in images due to variation in lighting and certain non-geometric camera parameters. We

did this by measuring the mappings of colors in one image of a scene to colors in another image of the same scene under different lighting conditions. Here we increase the flexibility of this color flow model by allowing flow coefficients to vary according to a low order polynomial over the image. This allows us to better fit smoothly varying lighting conditions as well as curved surfaces without ending our model with too much capacity. We show results on image matching and shadow removal and detection.

\*\*\*\*\*

#### Value-Directed Compression of POMDPs

Pascal Poupart, Craig Boutilier

We examine the problem of generating state-space compressions of POMDPs in a way that minimally impacts decision quality. We analyze the impact of compressions on decision quality, observing that compressions that allow accurate policy evaluation (prediction of expected future reward) will not affect decision quality. We derive a set of sufficient conditions that ensure accurate prediction in this respect, illustrate interesting mathematical properties these confer on lossless linear compressions, and use these to derive an iterative procedure for finding good linear lossy compressions. We also elaborate on how structured representations of a POMDP can be used to find such compressions.

\*\*\*\*\*

#### Constraint Classification for Multiclass Classification and Ranking

Sariel Har-Peled, Dan Roth, Dav Zimak

The constraint classification framework captures many flavors of multiclass classification including winner-take-all multiclass classification, multilabel classification and ranking. We present a meta-algorithm for learning in this framework that learns via a single linear classifier in high dimension. We discuss distribution independent as well as margin-based generalization bounds and present empirical and theoretical evidence showing that constraint classification benefits over existing methods of multiclass classification.

\*\*\*\*\*

#### Neural Decoding of Cursor Motion Using a Kalman Filter

W Wu, M. Black, Y. Gao, M. Serruya, A. Shaikhouni, J. Donoghue, Elie Bienenstock

The direct neural control of external devices such as computer displays or prosthetic limbs requires the accurate decoding of neural activity representing continuous movement. We develop a real-time control system using the spiking activity of approximately 40 neurons recorded with an electrode array implanted in the arm area of primary motor cortex. In contrast to previous work, we develop a control-theoretic approach that explicitly models the motion of the hand and the probabilistic relationship between this motion and the mean firing rates of the cells in 70 bins. We focus on a realistic cursor control task in which the subject must move a cursor to "hit" randomly placed targets on a computer monitor. Encoding and decoding of the neural data is achieved with a Kalman filter which has a number of advantages over previous linear filtering techniques. In particular, the Kalman filter reconstructions of hand trajectories in off-line experiments are more accurate than previously reported results and the model provides insights into the nature of the neural coding of movement.

\*\*\*\*\*

#### On the Dirichlet Prior and Bayesian Regularization

Harald Steck, Tommi Jaakkola

A common objective in learning a model from data is to recover its network structure, while the model parameters are of minor interest. For example, we may wish to recover regulatory networks from high-throughput data sources. In this paper we examine how Bayesian regularization using a product of independent Dirichlet priors over the model parameters affects the learned model structure in a domain with discrete variables. We show that a small scale parameter - often interpreted as "equivalent sample size" or "prior strength" - leads to a strong regularization of the model structure (sparse graph) given a sufficiently large dataset. In particular, the empty graph is obtained in the limit of a vanishing scale parameter. This is diametrically opposite to what one may expect

ect in this limit, namely the complete graph from an (unregularized) maximum likelihood estimate. Since the prior affects the parameters as expected, the scale parameter balances a trade-off between regularizing the parameters vs. the structure of the model. We demonstrate the benefits of optimizing this trade-off in the sense of predictive accuracy.

\*\*\*\*\*

#### Scaling of Probability-Based Optimization Algorithms

J. Shapiro

Population-based Incremental Learning is shown require very sensitive scaling of its learning rate. The learning rate must scale with the system size in a problem-dependent way. This is shown in two problems: the needle-in-a haystack, in which the learning rate must vanish exponentially in the system size, and in a smooth function in which the learning rate must vanish like the square root of the system size. Two methods are proposed for removing this sensitivity. A learning dynamics which obeys detailed balance is shown to give consistent performance over the entire range of learning rates. An analog of mutation is shown to require a learning rate which scales as the inverse system size, but is problem independent.

\*\*\*\*\*

#### Forward-Decoding Kernel-Based Phone Recognition

Shantanu Chakrabartty, Gert Cauwenberghs

Forward decoding kernel machines (FDKM) combine large-margin classifiers with hidden Markov models (HMM) for maximum a posteriori (MAP) adaptive sequence estimation. State transitions in the sequence are conditioned on observed data using a kernel-based probability model trained with a recursive scheme that deals effectively with noisy and partially labeled data. Training over very large data sets is accomplished using a sparse probabilistic support vector machine (SVM) model based on quadratic entropy, and an on-line stochastic steepest descent algorithm. For speaker-independent continuous phone recognition, FDKM trained over 177,080 samples of the TIMIT database achieves 80.6% recognition accuracy over the full test set, without use of a prior phonetic language model.

\*\*\*\*\*

#### Optoelectronic Implementation of a FitzHugh-Nagumo Neural Model

Alexandre Romariz, Kelvin Wagner

An optoelectronic implementation of a spiking neuron model based on the FitzHugh-Nagumo equations is presented. A tunable semiconductor laser source and a spectral filter provide a nonlinear mapping from driver voltage to detected signal. Linear electronic feedback completes the implementation, which allows either electronic or optical input signals. Experimental results for a single system and numeric results of model interaction confirm that important features of spiking neural models can be implemented through this approach.

\*\*\*\*\*

#### Margin-Based Algorithms for Information Filtering

Nicolò Cesa-bianchi, Alex Conconi, Claudio Gentile

In this work, we study an information filtering model where the relevance labels associated to a sequence of feature vectors are realizations of an unknown probabilistic linear function. Building on the analysis of a restricted version of our model, we derive a general filtering rule based on the margin of a ridge regression estimator. While our rule may observe the label of a vector only by classifying the vector as relevant, experiments on a real-world document filtering problem show that the performance of our rule is close to that of the on-line classifier which is allowed to observe all labels. These empirical results are complemented by a theoretical analysis where we consider a randomized variant of our rule and prove that its expected number of mistakes is never much larger than that of the optimal filtering rule which knows the hidden linear model.

\*\*\*\*\*

#### Half-Lives of EigenFlows for Spectral Clustering

Chakra Chennubhotla, Allan Jepson

Using a Markov chain perspective of spectral clustering we present an algorithm



to automatically find the number of stable clusters in a dataset. The Markov chain's behaviour is characterized by the spectral properties of the matrix of transition probabilities, from which we derive eigenvectors along with their half-lives. An eigenvector describes the flow of probability mass due to the Markov chain, and it is characterized by its eigenvalue, or equivalently, by the half-life of its decay as the Markov chain is iterated. A ideal stable cluster is one with zero eigenvector and infinite half-life. The key insight in this paper is that bottlenecks between weakly coupled clusters can be identified by computing the sensitivity of the eigenvector's half-life to variations in the edge weights. We propose a novel EIGENCUTS algorithm to perform clustering that removes these identified bottlenecks in an iterative fashion.

\*\*\*\*\*

The RA Scanner: Prediction of Rheumatoid Joint Inflammation Based on Laser Imaging

Anton Schwaighofer, Volker Tresp, Peter Mayer, Alexander Scheel, Gerhard Müller  
We describe the RA scanner, a novel system for the examination of patients suffering from rheumatoid arthritis. The RA scanner is based on a novel laser-based imaging technique which is sensitive to the optical characteristics of finger joint tissue. Based on the laser images, finger joints are classified according to whether the inflammatory status has improved or worsened. To perform the classification task, various linear and kernel-based systems were implemented and their performances were compared. Special emphasis was put on measures to reliably perform parameter tuning and evaluation, since only a very small data set was available. Based on the results presented in this paper, it was concluded that the RA scanner permits a reliable classification of pathological finger joints, thus paving the way for a further development from prototype to product stage.

\*\*\*\*\*

Optimality of Reinforcement Learning Algorithms with Linear Function Approximation

Ralf Schoknecht

There are several reinforcement learning algorithms that yield approximate solutions for the problem of policy evaluation when the value function is represented with a linear function approximator. In this paper we show that each of the solutions is optimal with respect to a specific objective function. Moreover, we characterise the different solutions as images of the optimal exact value function under different projection operations. The results presented here will be useful for comparing the algorithms in terms of the error they achieve relative to the error of the optimal approximate solution.

\*\*\*\*\*

Evidence Optimization Techniques for Estimating Stimulus-Response Functions

Maneesh Sahani, Jennifer Linden

An essential step in understanding the function of sensory nervous systems is to characterize as accurately as possible the stimulus-response function (SRF) of the neurons that relay and process sensory information. One increasingly common experimental approach is to present a rapidly varying complex stimulus to the animal while recording the responses of one or more neurons, and then to directly estimate a functional transformation of the input that accounts for the neuronal firing. The estimation techniques usually employed, such as Wiener filtering or other correlation-based estimation of the Wiener or Volterra kernels, are equivalent to maximum likelihood estimation in a Gaussian-output-noise regression model. We explore the use of Bayesian evidence-optimization techniques to condition these estimates. We show that by learning hyperparameters that control the smoothness and sparsity of the transfer function it is possible to improve dramatically the quality of SRF estimates, as measured by their success in predicting responses to novel input.

\*\*\*\*\*

Binary Coding in Auditory Cortex

Michael Deweese, Anthony Zador

Cortical neurons have been reported to use both rate and temporal codes. Here we describe a novel mode in which each neuron generates exactly 0 or 1 action po

tentials, but not more, in response to a stimulus. We used cell-attached recording, which ensured single-unit isolation, to record responses in rat auditory cortex to brief tone pips. Surprisingly, the majority of neurons exhibited binary behavior with few multi-spike responses; several dramatic examples consisted of exactly one spike on 100% of trials, with no trial-to-trial variability in spike count. Many neurons were tuned to stimulus frequency. Since individual trials yielded at most one spike for most neurons, the information about stimulus frequency was encoded in the population, and would not have been accessible to later stages of processing that only had access to the activity of a single unit. These binary units allow a more efficient population code than is possible with conventional rate coding units, and are consistent with a model of cortical processing in which synchronous packets of spikes propagate stably from one neuronal population to the next.

\*\*\*\*\*

#### Learning Attractor Landscapes for Learning Motor Primitives

Auke Ijspeert, Jun Nakanishi, Stefan Schaal

Many control problems take place in continuous state-action spaces, e.g., as in manipulator robotics, where the control objective is often defined as finding a desired trajectory that reaches a particular goal state. While reinforcement learning offers a theoretical framework to learn such control policies from scratch, its applicability to higher dimensional continuous state-action spaces remains rather limited to date. Instead of learning from scratch, in this paper we suggest to learn a desired complex control policy by transforming an existing simple canonical control policy. For this purpose, we represent canonical policies in terms of differential equations with well-defined attractor properties. By nonlinearly transforming the canonical attractor dynamics using techniques from nonparametric regression, almost arbitrary new nonlinear policies can be generated without losing the stability properties of the canonical system. We demonstrate our techniques in the context of learning a set of movement skills for a humanoid robot from demonstrations of a human teacher. Policies are acquired rapidly, and, due to the properties of well formulated differential equations, can be reused and modified on-line under dynamic changes of the environment. The linear parameterization of nonparametric regression moreover lends itself to recognize and classify previously learned movement skills. Evaluations in simulations and on an actual 30 degree-of-freedom humanoid robot exemplify the feasibility and robustness of our approach.

\*\*\*\*\*

#### Combining Dimensions and Features in Similarity-Based Representations

Daniel Navarro, Michael Lee

This paper develops a new representational model of similarity data that combines continuous dimensions with discrete features. An algorithm capable of learning these representations is described, and a Bayesian model selection approach for choosing the appropriate number of dimensions and features is developed. The approach is demonstrated on a classic data set that considers the similarities between the numbers 0 through 9.

\*\*\*\*\*

#### Bayesian Monte Carlo

Zoubin Ghahramani, Carl Rasmussen

We investigate Bayesian alternatives to classical Monte Carlo methods for evaluating integrals. Bayesian Monte Carlo (BMC) allows the incorporation of prior knowledge, such as smoothness of the integrand, into the estimation. In a simple problem we show that this outperforms any classical importance sampling method. We also attempt more challenging multidimensional integrals involved in computing marginal likelihoods of statistical models (a.k.a. partition functions and model evidences). We find that Bayesian Monte Carlo outperformed Annealed Importance Sampling, although for very high dimensional problems or problems with massive multimodality BMC may be less adequate. One advantage of the Bayesian approach to Monte Carlo is that samples can be drawn from any distribution. This allows for the possibility of active design of sample points so as to maximise information gain.

\*\*\*\*\*

## A Model for Learning Variance Components of Natural Images

Yan Karklin, Michael Lewicki

We present a hierarchical Bayesian model for learning efficient codes of higher-order structure in natural images. The model, a non-linear generalization of independent component analysis, replaces the standard assumption of independence for the joint distribution of coefficients with a distribution that is adapted to the variance structure of the coefficients of an efficient image basis. This offers a novel description of higher-order image structure and provides a way to learn coarse-coded, sparse-distributed representations of abstract image properties such as object location, scale, and texture.

\*\*\*\*\*

## Effective Dimension and Generalization of Kernel Learning

Tong Zhang

We investigate the generalization performance of some learning problems in Hilbert function Spaces. We introduce a concept of scale-sensitive effective data dimension, and show that it characterizes the convergence rate of the underlying learning problem. Using this concept, we can naturally extend results for parametric estimation problems in finite dimensional spaces to non-parametric kernel learning methods. We derive upper bounds on the generalization performance and show that the resulting convergent rates are optimal under various circumstances.

\*\*\*\*\*

## Nash Propagation for Loopy Graphical Games

Luis E. Ortiz, Michael Kearns

We introduce NashProp, an iterative and local message-passing algorithm for computing Nash equilibria in multi-player games represented by arbitrary undirected graphs. We provide a formal analysis and experimental evidence demonstrating that NashProp performs well on large graphical games with many loops, often converging in just a dozen iterations on graphs with hundreds of nodes. NashProp generalizes the tree algorithm of (Kearns et al. 2001), and can be viewed as similar in spirit to belief propagation in probabilistic inference, and thus complements the recent work of (Vickrey and Koller 2002), who explored a junction tree approach. Thus, as for probabilistic inference, we have at least two promising general-purpose approaches to equilibria computation in graphs.

\*\*\*\*\*

## Neuromorphic Bistable VLSI Synapses with Spike-Timing-Dependent Plasticity

Giacomo Indiveri

We present analog neuromorphic circuits for implementing bistable synapses with spike-timing-dependent plasticity (STDP) properties. In these types of synapses, the short-term dynamics of the synaptic efficacies are governed by the relative timing of the pre- and post-synaptic spikes, while on long time scales the efficacies tend asymptotically to either a potentiated state or to a depressed one. We fabricated a prototype VLSI chip containing a network of integrate and fire neurons interconnected via bistable STDP synapses. Test results from this chip demonstrate the synapse's STDP learning properties, and its long-term bistable characteristics.

\*\*\*\*\*

## Source Separation with a Sensor Array using Graphical Models and Subband Filtering

Hagai Attias

Source separation is an important problem at the intersection of several fields, including machine learning, signal processing, and speech technology. Here we describe new separation algorithms which are based on probabilistic graphical models with latent variables. In contrast with existing methods, these algorithms exploit detailed models to describe source properties. They also use subband filtering ideas to model the reverberant environment, and employ an explicit model for background and sensor noise. We leverage variational techniques to keep the computational complexity per EM iteration linear in the number of frames.

\*\*\*\*\*

## Dopamine Induced Bistability Enhances Signal Processing in Spiny Neurons

Aaron Gruber, Sara Solla, James Houk

Single unit activity in the striatum of awake monkeys shows a marked dependence on the expected reward that a behavior will elicit. We present a computational model of spiny neurons, the principal neurons of the striatum, to assess the hypothesis that direct neuromodulatory effects of dopamine through the activation of D1 receptors mediate the reward dependency of spiny neuron activity. Dopamine release results in the amplification of key ion currents, leading to the emergence of bistability, which not only modulates the peak firing rate but also introduces a temporal and state dependence of the model's response, thus improving the detectability of temporally correlated inputs.

\*\*\*\*\*

## An Information Theoretic Approach to the Functional Classification of Neurons

Elad Schneidman, William Bialek, Michael Li

A population of neurons typically exhibits a broad diversity of responses to sensory inputs. The intuitive notion of functional classification is that cells can be clustered so that most of the diversity is captured by the identity of the clusters rather than by individuals within clusters. We show how this intuition can be made precise using information theory, without any need to introduce a metric on the space of stimuli or responses. Applied to the retinal ganglion cells of the salamander, this approach recovers classical results, but also provides clear evidence for subclasses beyond those identified previously. Further, we find that each of the ganglion cells is functionally unique, and that even within the same subclass only a few spikes are needed to reliably distinguish between cells.

\*\*\*\*\*

## An Asynchronous Hidden Markov Model for Audio-Visual Speech Recognition

Samy Bengio

This paper presents a novel Hidden Markov Model architecture to model the joint probability of pairs of asynchronous sequences describing the same event. It is based on two other Markovian models, namely Asynchronous Input / Output Hidden Markov Models and Pair Hidden Markov Models. An EM algorithm to train the model is presented, as well as a Viterbi decoder that can be used to obtain the optimal state sequence as well as the alignment between the two sequences. The model has been tested on an audio-visual speech recognition task using the M2VTS database and yielded robust performances under various noise conditions.

\*\*\*\*\*

## Shape Recipes: Scene Representations that Refer to the Image

William Freeman, Antonio Torralba

The goal of low-level vision is to estimate an underlying scene, given an observed image. Real-world scenes (eg, albedos or shapes) can be very complex, conventionally requiring high dimensional representations which are hard to estimate and store. We propose a low-dimensional representation, called a scene recipe, that relies on the image itself to describe the complex scene configurations. Shape recipes are an example: these are the regression coefficients that predict the bandpassed shape from image data. We describe the benefits of this representation, and show two uses illustrating their properties: (1) we improve stereo shape estimates by learning shape recipes at low resolution and applying them at full resolution; (2) Shape recipes implicitly contain information about lighting and materials and we use them for material segmentation.

\*\*\*\*\*

## Real-Time Particle Filters

Cody Kwok, Dieter Fox, Marina Meila

Particle filters estimate the state of dynamical systems from sensor information. In many real time applications of particle filters, however, sensor information arrives at a significantly higher rate than the update rate of the filter. The prevalent approach to dealing with such situations is to update the particle filter as often as possible and to discard sensor information that cannot be processed

d in time. In this paper we present real-time particle filters, which make use of all sensor information even when the filter update rate is below the update rate of the sensors. This is achieved by representing posteriors as mixtures of sample sets, where each mixture component integrates one observation arriving during a filter update. The weights of the mixture components are set so as to minimize the approximation error introduced by the mixture representation. Thereby, our approach focuses computational resources (samples) on valuable sensor information. Experiments using data collected with a mobile robot show that our approach yields strong improvements over other approaches.

\*\*\*\*\*

#### Critical Lines in Symmetry of Mixture Models and its Application to Component Splitting

Kenji Fukumizu, Shotaro Akaho, Shun-ichi Amari

We show the existence of critical points as lines for the likelihood function of mixture-type models. They are given by embedding of a critical point for models with less components. A sufficient condition that the critical line gives local maxima or saddle points is also derived. Based on this fact, a component-split method is proposed for a mixture of Gaussian components, and its effectiveness is verified through experiments.

\*\*\*\*\*

#### Manifold Parzen Windows

Pascal Vincent, Yoshua Bengio

The similarity between objects is a fundamental element of many learning algorithms. Most non-parametric methods take this similarity to be fixed, but much recent work has shown the advantages of learning it, in particular to exploit the local invariances in the data or to capture the possibly non-linear manifold on which most of the data lies. We propose a new non-parametric kernel density estimation method which captures the local structure of an underlying manifold through the leading eigenvectors of regularized local covariance matrices. Experiments in density estimation show significant improvements with respect to Parzen density estimators. The density estimators can also be used within Bayes classifiers, yielding classification rates similar to SVMs and much superior to the Parzen classifier.

\*\*\*\*\*

#### Parametric Mixture Models for Multi-Labeled Text

Naonori Ueda, Kazumi Saito

We propose probabilistic generative models, called parametric mixture models (PMMs), for multiclass, multi-labeled text categorization problem. Conventionally, the binary classification approach has been employed, in which whether or not text belongs to a category is judged by the binary classifier for every category. In contrast, our approach can simultaneously detect multiple categories of text using PMMs. We derive efficient learning and prediction algorithms for PMMs. We also empirically show that our method could significantly outperform the conventional binary methods when applied to multi-labeled text categorization using real World Wide Web pages.

\*\*\*\*\*

#### Transductive and Inductive Methods for Approximate Gaussian Process Regression

Anton Schwaighofer, Volker Tresp

Gaussian process regression allows a simple analytical treatment of exact Bayesian inference and has been found to provide good performance, yet scales badly with the number of training data. In this paper we compare several approaches towards scaling Gaussian processes regression to large data sets: the subset of representers method, the reduced rank approximation, online Gaussian processes, and the Bayesian committee machine. Furthermore we provide theoretical insight into some of our experimental results. We found that subset of representers methods can give good and particularly fast predictions for data sets with high and medium noise levels. On complex low noise data sets, the Bayesian committee machine achieves significantly better accuracy, yet at a higher computational cost.

\*\*\*\*\*

#### Adapting Codes and Embeddings for Polychotomies

Gunnar Rätsch, Sebastian Mika, Alex Smola

In this paper we consider formulations of multi-class problems based on a generalized notion of a margin and using output coding. This includes, but is not restricted to, standard multi-class SVM formulations. Differently from many previous approaches we learn the code as well as the embedding function. We illustrate how this can lead to a formulation that allows for solving a wider range of problems with for instance many classes or even "missing classes". To keep our optimization problems tractable we propose an algorithm capable of solving them using two-class classifiers, similar in spirit to Boosting.

\*\*\*\*\*

Topographic Map Formation by Silicon Growth Cones

Brian Taba, Kwabena A. Boahen

We describe a self-configuring neuromorphic chip that uses a model of activity-dependent axon remodeling to automatically wire topographic maps based solely on input correlations. Axons are guided by growth cones, which are modeled in analog VLSI for the first time. Growth cones migrate up neurotrophin gradients, which are represented by charge diffusing in transistor channels. Virtual axons move by rerouting address-events. We refined an initially gross topographic projection by simulating retinal wave input.

\*\*\*\*\*

Analysis of Information in Speech Based on MANOVA

Sachin Kajarekar, Hynek Hermansky

We propose analysis of information in speech using three sources - language (phone), speaker and channel. Information in speech is measured as mutual information between the source and the set of features extracted from speech signal. We assume that distribution of features can be modeled using Gaussian distribution. The mutual information is computed using the results of analysis of variability in speech. We observe similarity in the results of phone variability and phone information, and show that the results of the proposed analysis have more meaningful interpretations than the analysis of variability.

\*\*\*\*\*

Discriminative Binaural Sound Localization

Ehud Ben-reuven, Yoram Singer

Time difference of arrival (TDOA) is commonly used to estimate the azimuth of a source in a microphone array. The most common methods to estimate TDOA are based on finding extrema in generalized cross-correlation waveforms. In this paper we apply microphone array techniques to a manikin head. By considering the entire cross-correlation waveform we achieve azimuth prediction accuracy that exceeds extrema locating methods. We do so by quantizing the azimuthal angle and treating the prediction problem as a multiclass categorization task. We demonstrate the merits of our approach by evaluating the various approaches on Sony's AIBO robot.

\*\*\*\*\*

Fractional Belief Propagation

Wim Wiegerinck, Tom Heskes

We consider loopy belief propagation for approximate inference in probabilistic graphical models. A limitation of the standard algorithm is that clique marginals are computed as if there were no loops in the graph. To overcome this limitation, we introduce fractional belief propagation. Fractional belief propagation is formulated in terms of a family of approximate free energies, which includes the Bethe free energy and the naive mean-field free energy as special cases. Using the linear response correction of the clique marginals, the scale parameters can be tuned. Simulation results illustrate the potential merits of the approach.

\*\*\*\*\*

Stability-Based Model Selection

Tilman Lange, Mikio Braun, Volker Roth, Joachim Buhmann

Model selection is linked to model assessment, which is the problem of comparing different models, or model parameters, for a specific learning task. For supervi-

sed learning, the standard practical technique is cross-validation, which is not applicable for semi-supervised and unsupervised settings. In this paper, a new model assessment scheme is introduced which is based on a notion of stability. The stability measure yields an upper bound to cross-validation in the supervised case, but extends to semi-supervised and unsupervised problems. In the experimental part, the performance of the stability measure is studied for model order selection in comparison to standard techniques in this area.

\*\*\*\*\*

#### Hidden Markov Model of Cortical Synaptic Plasticity: Derivation of the Learning Rule

Michael Eisele, Kenneth Miller

Cortical synaptic plasticity depends on the relative timing of pre- and postsynaptic spikes and also on the temporal pattern of presynaptic spikes and of postsynaptic spikes. We study the hypothesis that cortical synaptic plasticity does not associate individual spikes, but rather whole firing episodes, and depends only on when these episodes start and how long they last, but as little as possible on the timing of individual spikes. Here we present the mathematical background for such a study. Standard methods from hidden Markov models are used to determine what "firing episodes" are. Estimating the probability of being in such an episode requires not only the knowledge of past spikes, but also of future spikes. We show how to construct a causal learning rule, which depends only on past spikes, but associates pre- and postsynaptic firing episodes as if it also knew future spikes. We also show that this learning rule agrees with some features of synaptic plasticity in superficial layers of rat visual cortex (Froemke and Dan, Nature 416:433, 2002).

\*\*\*\*\*

#### Automatic Alignment of Local Representations

Yee Teh, Sam Roweis

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.

\*\*\*\*\*

#### Field-Programmable Learning Arrays

Seth Bridges, Miguel Figueroa, Chris Diorio, David Hsu

This paper introduces the Field-Programmable Learning Array, a new paradigm for rapid prototyping of learning primitives and machine-learning algorithms in silicon. The FPLA is a mixed-signal counterpart to the all-digital Field-Programmable Gate Array in that it enables rapid prototyping of algorithms in hardware. Unlike the FPGA, the FPLA is targeted directly for machine learning by providing local, parallel, on-line analog learning using floating-gate MOS synapse transistors. We present a prototype FPLA chip comprising an array of reconfigurable computational blocks and local interconnect. We demonstrate the viability of this architecture by mapping several learning circuits onto the prototype chip.

\*\*\*\*\*

#### Boosting Density Estimation

Saharon Rosset, Eran Segal

Several authors have suggested viewing boosting as a gradient descent search for a good fit in function space. We apply gradient-based boosting methodology to the unsupervised learning problem of density estimation. We show convergence properties of the algorithm and prove that a strength of weak learnability property applies to this problem as well. We illustrate the potential of this approach through experiments with boosting Bayesian networks to learn density models.

\*\*\*\*\*

#### Support Vector Machines for Multiple-Instance Learning

Stuart Andrews, Ioannis Tsochantaridis, Thomas Hofmann

This paper presents two new formulations of multiple-instance learning as a maximum margin problem. The proposed extensions of the Support Vector Machine (SVM) learning approach lead to mixed integer quadratic programs that can be solved heuristically. Our generalization of SVMs makes a state-of-

-the-art classification technique, including non-linear classification via kernels, available to an area that up to now has been largely dominated by special purpose methods. We present experimental results on a pharmaceutical data set and on applications in automated image indexing and document categorization.

\*\*\*\*\*

#### Bias-Optimal Incremental Problem Solving

Jürgen Schmidhuber

Given is a problem sequence and a probability distribution (the bias) on programs computing solution candidates. We present an optimally fast way of incrementally solving each task in the sequence. Bias shifts are computed by program prefixes that modify the distribution on their suffixes by reusing successful code for previous tasks (stored in non-modifiable memory). No tested program gets more runtime than its probability times the total search time. In illustrative experiments, ours becomes the first general system to learn a universal solver for arbitrary disk Towers of Hanoi tasks (minimal solution size). It demonstrates the advantages of incremental learning by profiting from previously solved, simpler tasks involving samples of a simple context free language.

\*\*\*\*\*

#### Rate Distortion Function in the Spin Glass State: A Toy Model

Tatsuto Murayama, Masato Okada

We applied statistical mechanics to an inverse problem of linear mapping to investigate the physics of optimal lossy compressions. We used the replica symmetry breaking technique with a toy model to demonstrate Shannon's result. The rate distortion function, which is widely known as the theoretical limit of the compression with a fidelity criterion, is derived. Numerical study shows that sparse constructions of the model provide suboptimal compressions.

\*\*\*\*\*

#### Adaptive Nonlinear System Identification with Echo State Networks

Herbert Jaeger

Echo state networks (ESN) are a novel approach to recurrent neural network training. An ESN consists of a large, fixed, recurrent "reservoir" network, from which the desired output is obtained by training suitable output connection weights. Determination of optimal output weights becomes a linear, uniquely solvable task of MSE minimization. This article reviews the basic ideas and describes an online adaptation scheme based on the RLS algorithm known from adaptive linear systems. As an example, a 10-th order NARMA system is adaptively identified. The known benefits of the RLS algorithms carryover from linear systems to nonlinear ones; specifically, the convergence rate and misadjustment can be determined at design time.

\*\*\*\*\*

#### Real Time Voice Processing with Audiovisual Feedback: Toward Autonomous Agents with Perfect Pitch

Lawrence Saul, Daniel Lee, Charles Isbell, Yann Cun

We have implemented a real time front end for detecting voiced speech and estimating its fundamental frequency. The front end performs the signal processing for voice-driven agents that attend to the pitch contours of human speech and provide continuous audiovisual feedback. The algorithm we use for pitch tracking has several distinguishing features: it makes no use of FFTs or autocorrelation at the pitch period; it updates the pitch incrementally on a sample-by-sample basis; it avoids peak picking and does not require interpolation in time or frequency to obtain high resolution estimates; and it works reliably over a four octave range, in real time, without the need for postprocessing to produce smooth contours. The algorithm is based on two simple ideas in neural computation: the introduction of a purposeful nonlinearity, and the error signal of a least squares fit. The pitch tracker is used in two real time multimedia applications: a voice-to-MIDI player that synthesizes electronic music from vocalized melodies, and an audiovisual Karaoke machine with multimodal feedback. Both applications run on a laptop and display the user's pitch scrolling across the screen as he or she



he sings into the computer.

\*\*\*\*\*

## An Impossibility Theorem for Clustering

Jon Kleinberg

Although the study of clustering is centered around an intuitively compelling goal, it has been very difficult to develop a unified framework for reasoning about it at a technical level, and profoundly diverse approaches to clustering abound in the research community. Here we suggest a formal perspective on the difficulty in finding such a unification, in the form of an impossibility theorem: for a set of three simple properties, we show that there is no clustering function satisfying all three. Relaxations of these properties expose some of the interesting (and unavoidable) trade-offs at work in well-studied clustering techniques such as single-linkage, sum-of-pairs, k-means, and k-median.

\*\*\*\*\*

## Visual Development Aids the Acquisition of Motion Velocity Sensitivities

Robert Jacobs, Melissa Dominguez

We consider the hypothesis that systems learning aspects of visual perception may benefit from the use of suitably designed developmental progressions during training. Four models were trained to estimate motion velocities in sequences of visual images. Three of the models were "developmental models" in the sense that the nature of their input changed during the course of training. They received a relatively impoverished visual input early in training, and the quality of this input improved as training progressed. One model used a coarse-to-multiscale developmental progression (i.e. it received coarse-scale motion features early in training and finer-scale features were added to its input as training progressed), another model used a fine-to-multiscale progression, and the third model used a random progression. The final model was non-developmental in the sense that the nature of its input remained the same throughout the training period. The simulation results show that the coarse-to-multiscale model performed best. Hypotheses are offered to account for this model's superior performance. We conclude that suitably designed developmental sequences can be useful to systems learning to estimate motion velocities. The idea that visual development can aid visual learning is a viable hypothesis in need of further study.

\*\*\*\*\*

## On the Complexity of Learning the Kernel Matrix

Olivier Bousquet, Daniel Herrmann

We investigate data based procedures for selecting the kernel when learning with Support Vector Machines. We provide generalization error bounds by estimating the Rademacher complexities of the corresponding function classes. In particular we obtain a complexity bound for function classes induced by kernels with given eigenvectors, i.e., we allow to vary the spectrum and keep the eigenvectors fixed. This bound is only a logarithmic factor bigger than the complexity of the function class induced by a single kernel. However, optimizing the margin over such classes leads to overfitting. We thus propose a suitable way of constraining the class. We use an efficient algorithm to solve the resulting optimization problem, present preliminary experimental results, and compare them to an alignment-based approach.

\*\*\*\*\*

## Clustering with the Fisher Score

Koji Tsuda, Motoaki Kawanabe, Klaus-Robert Müller

Recently the Fisher score (or the Fisher kernel) is increasingly used as a feature extractor for classification problems. The Fisher score is a vector of parameter derivatives of loglikelihood of a probabilistic model. This paper gives a theoretical analysis about how class information is preserved in the space of the Fisher score, which turns out that the Fisher score consists of a few important dimensions with class information and many nuisance dimensions. When we perform clustering with the Fisher score, K-Means type methods are obviously inappropriate because they make use of all dimensions. So we will develop a novel but simple clustering algorithm specialized for the Fisher score, which can exploit im

- portant dimensions. This algorithm is successfully tested in experiments with artificial data and real data (amino acid sequences).

\*\*\*\*\*

#### Artefactual Structure from Least-Squares Multidimensional Scaling

Nicholas Hughes, David Lowe

We consider the problem of illusory or artefactual structure from the visualization of high-dimensional structureless data. In particular we examine the role of the distance metric in the use of topographic mappings based on the statistical field of multidimensional scaling. We show that the use of a squared Euclidean metric (i.e. the STRESS measure) gives rise to an annular structure when the input data is drawn from a high-dimensional isotropic distribution, and we provide a theoretical justification for this observation.

\*\*\*\*\*

#### Adaptation and Unsupervised Learning

Peter Dayan, Maneesh Sahani, Gregoire Deback

Adaptation is a ubiquitous neural and psychological phenomenon, with a wealth of instantiations and implications. Although a basic form of plasticity, it has, bar some notable exceptions, attracted computational theory of only one main variety. In this paper, we study adaptation from the perspective of factor analysis, a paradigmatic technique of unsupervised learning. We use factor analysis to re-interpret a standard view of adaptation, and apply our new model to some recent data on adaptation in the domain of face discrimination.

\*\*\*\*\*

#### Learning in Zero-Sum Team Markov Games Using Factored Value Functions

Michail G. Lagoudakis, Ronald Parr

We present a new method for learning good strategies in zero-sum Markov games in which each side is composed of multiple agents collaborating against an opposing team of agents. Our method requires full observability and communication during learning, but the learned policies can be executed in a distributed manner. The value function is represented as a factored linear architecture and its structure determines the necessary computational resources and communication bandwidth. This approach permits a tradeoff between simple representations with little or no communication between agents and complex, computationally intensive representations with extensive coordination between agents. Thus, we provide a principled means of using approximation to combat the exponential blowup in the joint action space of the participants. The approach is demonstrated with an example that shows the efficiency gains over naive enumeration.

\*\*\*\*\*

#### Convergence Properties of Some Spike-Triggered Analysis Techniques

Liam Paninski

We analyze the convergence properties of three spike-triggered data analysis techniques. All of our results are obtained in the setting of a (possibly multidimensional) linear-nonlinear (LN) cascade model for stimulus-driven neural activity. We start by giving exact rate of convergence results for the common spike-triggered average (STA) technique. Next, we analyze a spike-triggered covariance method, variants of which have been recently exploited successfully by Bialek, Simoncelli, and colleagues. These first two methods suffer from extraneous conditions on their convergence; therefore, we introduce an estimator for the LN model parameters which is designed to be consistent under general conditions. We provide an algorithm for the computation of this estimator and derive its rate of convergence. We close with a brief discussion of the efficiency of these estimators and an application to data recorded from the primary motor cortex of awake, behaving primates.

\*\*\*\*\*

#### Learning to Take Concurrent Actions

Khashayar Rohanimanesh, Sridhar Mahadevan

We investigate a general semi-Markov Decision Process (SMDP) framework for modeling concurrent decision making, where agents learn optimal plans over concurrent temporally extended actions. We introduce three types of parallel termination schemes { all, any and continue } and theoretically and experimentally compare them

em.

\*\*\*\*\*

#### Handling Missing Data with Variational Bayesian Learning of ICA

Kwokleung Chan, Te-Won Lee, Terrence J. Sejnowski

Missing data is common in real-world datasets and is a problem for many estimation techniques. We have developed a variational Bayesian method to perform Independent Component Analysis (ICA) on high-dimensional data containing missing entries. Missing data are handled naturally in the Bayesian framework by integrating the generative density model. Modeling the distributions of the independent sources with mixture of Gaussians allows sources to be estimated with different kurtosis and skewness. The variational Bayesian method automatically determines the dimensionality of the data and yields an accurate density model for the observed data without overfitting problems. This allows direct probability estimation of missing values in the high dimensional space and avoids dimension reduction preprocessing which is not feasible with missing data.

\*\*\*\*\*

#### Ranking with Large Margin Principle: Two Approaches

Amnon Shashua, Anat Levin

We discuss the problem of ranking  $k$  instances with the use of a "large margin" principle. We introduce two main approaches: the first is the "fixed margin" policy in which the margin of the closest neighboring classes is being maximized - which turns out to be a direct generalization of SVM to ranking learning. The second approach allows for  $k - 1$  different margins where the sum of margins is maximized. This approach is shown to reduce to LI-SVM when the number of classes  $k = 2$ . Both approaches are optimal in size of  $2l$  where  $l$  is the total number of training examples. Experiments performed on visual classification and "collaborative filtering" show that both approaches outperform existing ordinal regression algorithms applied for ranking and multi-class SVM applied to general multi-class classification.

\*\*\*\*\*

#### A Bilinear Model for Sparse Coding

David Grimes, Rajesh P. N. Rao

Recent algorithms for sparse coding and independent component analysis (ICA) have demonstrated how localized features can be learned from natural images. However, these approaches do not take image transformations into account. As a result, they produce image codes that are redundant because the same feature is learned at multiple locations. We describe an algorithm for sparse coding based on a bilinear generative model of images. By explicitly modeling the interaction between image features and their transformations, the bilinear approach helps reduce redundancy in the image code and provides a basis for transformation-invariant vision. We present results demonstrating bilinear sparse coding of natural images. We also explore an extension of the model that can capture spatial relationships between the independent features of an object, thereby providing a new framework for parts-based object recognition.

\*\*\*\*\*

#### Data-Dependent Bounds for Bayesian Mixture Methods

Ron Meir, Tong Zhang

We consider Bayesian mixture approaches, where a predictor is constructed by forming a weighted average of hypotheses from some space of functions. While such procedures are known to lead to optimal predictors in several cases, where sufficiently accurate prior information is available, it has not been clear how they perform when some of the prior assumptions are violated. In this paper we establish data-dependent bounds for such procedures, extending previous randomized approaches such as the Gibbs algorithm to a fully Bayesian setting. The guarantees established in this work enable the utilization of Bayesian mixture approaches in agnostic settings, where the usual assumptions of the Bayesian paradigm fail to hold. Moreover, the bounds derived can be directly applied to non-Bayesian mixture approaches such as Bagging and Boosting.

\*\*\*\*\*

#### Information Diffusion Kernels

Guy Lebanon, John Lafferty

A new family of kernels for statistical learning is introduced that exploits the geometric structure of statistical models. Based on the heat equation on the Riemannian manifold defined by the Fisher information metric, information diffusion kernels generalize the Gaussian kernel of Euclidean space, and provide a natural way of combining generative statistical modeling with non-parametric discriminative learning. As a special case, the kernels give a new approach to applying kernel-based learning algorithms to discrete data. Bounds on covering numbers for the new kernels are proved using spectral theory in differential geometry, and experimental results are presented for text classification.

\*\*\*\*\*

Boosted Dyadic Kernel Discriminants

Baback Moghaddam, Gregory Shakhnarovich

We introduce a novel learning algorithm for binary classification with hyperplane discriminants based on pairs of training points from opposite classes (dyadic hypercuts). This algorithm is further extended to nonlinear discriminants using kernel functions satisfying Mercer's conditions. An ensemble of simple dyadic hypercuts is learned incrementally by means of a confidence-rated version of AdaBoost, which provides a sound strategy for searching through the finite set of hypercut hypotheses. In experiments with real-world datasets from the UCI repository, the generalization performance of the hypercut classifiers was found to be comparable to that of SVMs and k-NN classifiers. Furthermore, the computational cost of classification (at run time) was found to be similar to, or better than, that of SVM. Similarly to SVMs, boosted dyadic kernel discriminants tend to maximize the margin (via AdaBoost). In contrast to SVMs, however, we offer an on-line and incremental learning machine for building kernel discriminants whose complexity (number of kernel evaluations) can be directly controlled (traded off for accuracy).

\*\*\*\*\*

Derivative Observations in Gaussian Process Models of Dynamic Systems

E. Solak, R. Murray-smith, W. Leithead, D. Leith, Carl Rasmussen

Gaussian processes provide an approach to nonparametric modelling which allows a straightforward combination of function and derivative observations in an empirical model. This is of particular importance in identification of nonlinear dynamic systems from experimental data. 1) It allows us to combine derivative information, and associated uncertainty with normal function observations into the learning and inference process. This derivative information can be in the form of priors specified by an expert or identified from perturbation data close to equilibrium. 2) It allows a seamless fusion of multiple local linear models in a consistent manner, inferring consistent models and ensuring that integrability constraints are met. 3) It improves dramatically the computational efficiency of Gaussian process models for dynamic system identification, by summarising large quantities of near-equilibrium data by a handful of linearisations, reducing the training set size - traditionally a problem for Gaussian process models.

\*\*\*\*\*

Global Versus Local Methods in Nonlinear Dimensionality Reduction

Vin Silva, Joshua Tenenbaum

Recently proposed algorithms for nonlinear dimensionality reduction fall broadly into two categories which have different advantages and disadvantages: global (Isomap [1]), and local (Locally Linear Embedding [2], Laplacian Eigenmaps [3]). We present two variants of Isomap which combine the advantages of the global approach with what have previously been exclusive advantages of local methods: computational sparsity and the ability to invert conformal maps.

\*\*\*\*\*

Learning Graphical Models with Mercer Kernels

Francis Bach, Michael Jordan

We present a class of algorithms for learning the structure of graphical models from data. The algorithms are based on a measure known as the kernel generalized variance (KGV), which essentially allows us to treat all variables on an equal footing as Gaussians in a feature space obtained from Mercer kernels. Thus we ar

able to learn hybrid graphs involving discrete and continuous variables of arbitrary type. We explore the computational properties of our approach, showing how to use the kernel trick to compute the relevant statistics in linear time. We illustrate our framework with experiments involving discrete and continuous data.

\*\*\*\*\*

#### Stochastic Neighbor Embedding

Geoffrey E. Hinton, Sam Roweis

We describe a probabilistic approach to the task of placing objects, described by high-dimensional vectors or by pairwise dissimilarities, in a low-dimensional space in a way that preserves neighbor identities. A Gaussian is centered on each object in the high-dimensional space and the densities under this Gaussian (or the given dissimilarities) are used to define a probability distribution over all the potential neighbors of the object. The aim of the embedding is to approximate this distribution as well as possible when the same operation is performed on the low-dimensional "images" of the objects. A natural cost function is a sum of Kullback-Leibler divergences, one per object, which leads to a simple gradient for adjusting the positions of the low-dimensional images. Unlike other dimensionality reduction methods, this probabilistic framework makes it easy to represent each object by a mixture of widely separated low-dimensional images. This allows ambiguous objects, like the document count vector for the word "bank", to have versions close to the images of both "river" and "finance" without forcing the images of outdoor concepts to be located close to those of corporate concepts.

\*\*\*\*\*

#### Learning in Spiking Neural Assemblies

David Barber

We consider a statistical framework for learning in a class of networks of spiking neurons. Our aim is to show how optimal local learning rules can be readily derived once the neural dynamics and desired functionality of the neural assembly have been specified, in contrast to other models which assume (sub-optimal) learning rules. Within this framework we derive local rules for learning temporal sequences in a model of spiking neurons and demonstrate its superior performance to correlation (Hebbian) based approaches. We further show how to include mechanisms such as synaptic depression and outline how the framework is readily extensible to learning in networks of highly complex spiking neurons. A stochastic quantal vesicle release mechanism is considered and implications on the complexity of learning discussed.

\*\*\*\*\*

#### A Model for Real-Time Computation in Generic Neural Microcircuits

Wolfgang Maass, Thomas Natschläger, Henry Markram

Henry Markram Brain Mind Institute

\*\*\*\*\*

#### Categorization Under Complexity: A Unified MDL Account of Human Learning of Regular and Irregular Categories

David Fass, Jacob Feldman

We present an account of human concept learning—that is, learning of categories from examples—based on the principle of minimum description length (MDL). In support of this theory, we tested a wide range of two-dimensional concept types, including both regular (simple) and highly irregular (complex) structures, and found the MDL theory to give a good account of subjects' performance. This suggests that the intrinsic complexity of a concept (that is, its description length) systematically influences its learnability.

\*\*\*\*\*

#### Learning with Multiple Labels

Rong Jin, Zoubin Ghahramani

In this paper, we study a special kind of learning problem in which each training instance is given a set of (or distribution over) candidate class labels and only one of the candidate labels is the correct one. Such a problem can occur, e.g., in an information retrieval setting w

here a set of words is associated with an image, or if classes labels are organized hierarchically. We propose a novel discriminative approach for handling the ambiguity of class labels in the training examples. The experiments with the proposed approach over five different UCI datasets show that our approach is able to find the correct label among the set of candidate labels and actually achieve performance close to the case when each training instance is given a single correct label. In contrast, naive methods degrade rapidly as more ambiguity is introduced into the labels.

\*\*\*\*\*

#### Branching Law for Axons

Dmitri Chklovskii, Armen Stepanyants

What determines the caliber of axonal branches? We pursue the hypothesis that the axonal caliber has evolved to minimize signal propagation delays, while keeping arbor volume to a minimum. We show that for a general cost function the optimal diameters of mother ( $d_0$ ) and daughter ( $d_1, d_2$ ) branches at a bifurcation obey a branching law:  $d_i \propto v_i^{1/(1+V)}$  (where  $v_i$  is the conduction speed of branch  $i$ ). This law is based on the fact that the conduction speed scales with the axon diameter to the power  $V$  ( $V = 1$  for myelinated axons and  $V = 0.5$  for unmyelinated axons). We test the branching law on the available experimental data and find a reasonable agreement.

\*\*\*\*\*

#### Dyadic Classification Trees via Structural Risk Minimization

Clayton Scott, Robert Nowak

Classification trees are one of the most popular types of classifiers, with ease of implementation and interpretation being among their attractive features. Despite the widespread use of classification trees, theoretical analysis of their performance is scarce. In this paper, we show that a new family of classification trees, called dyadic classification trees (DCTs), are near optimal (in a minimax sense) for a very broad range of classification problems. This demonstrates that other schemes (e.g., neural networks, support vector machines) cannot perform significantly better than DCTs in many cases. We also show that this near optimal performance is attained with linear (in the number of training data) complexity growing and pruning algorithms. Moreover, the performance of DCTs on benchmark datasets compares favorably to that of standard CART, which is generally more computationally intensive and which does not possess similar near optimality properties. Our analysis stems from theoretical results on structural risk minimization, on which the pruning rule for DCTs is based.

\*\*\*\*\*

#### PAC-Bayes & Margins

John Langford, John Shawe-Taylor

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.

\*\*\*\*\*

#### Dynamical Causal Learning

David Danks, Thomas Griffiths, Joshua Tenenbaum

theories of human causal

\*\*\*\*\*

#### Maximally Informative Dimensions: Analyzing Neural Responses to Natural Signals

Tatyana Sharpee, Nicole Rust, William Bialek

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.

\*\*\*\*\*

#### Fast Exact Inference with a Factored Model for Natural Language Parsing

Dan Klein, Christopher D. Manning

We present a novel generative model for natural language tree structures in which

h semantic (lexical dependency) and syntactic (PCFG) structures are scored with separate models. This factorization provides conceptual simplicity, straightforward opportunities for separately improving the component models, and a level of performance comparable to similar, non-factored models. Most importantly, unlike other modern parsing models, the factored model admits an extremely effective A\* parsing algorithm, which enables efficient, exact inference.

\*\*\*\*\*

#### Morton-Style Factorial Coding of Color in Primary Visual Cortex

Javier Movellan, Thomas Wachtler, Thomas D. Albright, Terrence Sejnowski

We introduce the notion of Morton-style factorial coding and illustrate how it may help understand information integration and perceptual coding in the brain.

We show that by focusing on average responses one may miss the existence of factorial coding mechanisms that become only apparent when analyzing spike count histograms. We show evidence suggesting that the classical/non-classical receptive field organization in the cortex effectively enforces the development of Morton-style factorial codes. This may provide some cues to help understand perceptual coding in the brain and to develop new unsupervised learning algorithms. While methods like ICA (Bell & Sejnowski, 1997) develop independent codes, in Morton-style coding the goal is to make two or more external aspects of the world become independent when conditioning on internal representations.

\*\*\*\*\*

#### One-Class LP Classifiers for Dissimilarity Representations

Elzbieta Pekalska, David M.J. Tax, Robert Duin

Problems in which abnormal or novel situations should be detected can be approached by describing the domain of the class of typical examples. These applications come from the areas of machine diagnostics, fault detection, illness identification or, in principle, refer to any problem where little knowledge is available outside the typical class. In this paper we explain why proximities are natural representations for domain descriptors and we propose a simple one-class classifier for dissimilarity representations. By the use of linear programming an efficient one-class description can be found, based on a small number of prototype objects. This classifier can be made (1) more robust by transforming the dissimilarities and (2) cheaper to compute by using a reduced representation set. Finally, a comparison to a comparable one-class classifier by Campbell and Bennett is given.

\*\*\*\*\*

#### Regularized Greedy Importance Sampling

Finnegan Southey, Dale Schuurmans, Ali Ghodsi

Greedy importance sampling is an unbiased estimation technique that reduces the variance of standard importance sampling by explicitly searching for modes in the estimation objective. Previous work has demonstrated the feasibility of implementing this method and proved that the technique is unbiased in both discrete and continuous domains. In this paper we present a reformulation of greedy importance sampling that eliminates the free parameters from the original estimator, and introduces a new regularization strategy that further reduces variance without compromising unbiasedness. The resulting estimator is shown to be effective for difficult estimation problems arising in Markov random field inference.

In particular, improvements are achieved over standard MCMC estimators when the distribution has multiple peaked modes.

\*\*\*\*\*

#### Modeling Midazolam's Effect on the Hippocampus and Recognition Memory

Kenneth Malmberg, René Zeelenberg, Richard Shiffrin

'~lidazolam

\*\*\*\*\*

#### Fast Kernels for String and Tree Matching

Alex Smola, S.v.n. Vishwanathan

In this paper we present a new algorithm suitable for matching discrete objects such as strings and trees in linear time, thus obviating dynamic programming with quadratic time complexity. Furthermore, prediction cost in many cases can be reduced to linear cost in the length of the sequence to be classified

ied, regardless of the number of support vectors. This improvement on the currently available algorithms makes string kernels a viable alternative for the practitioner.

\*\*\*\*\*

Expected and Unexpected Uncertainty: ACh and NE in the Neocortex

Peter Dayan, Angela J. Yu

Inference and adaptation in noisy and changing, rich sensory environments are rife with a variety of specific sorts of variability. Experimental and theoretical studies suggest that these different forms of variability play different behavioral, neural and computational roles, and may be reported by different (notably neuromodulatory) systems. Here, we refine our previous theory of acetylcholine's role in cortical inference in the (oxymoronic) terms of expected uncertainty, and advocate a theory for norepinephrine in terms of unexpected uncertainty. We suggest that norepinephrine reports the radical divergence of bottom-up inputs from prevailing top-down interpretations, to influence inference and plasticity. We illustrate this proposal using an adaptive factor analysis model.

\*\*\*\*\*

Feature Selection and Classification on Matrix Data: From Large Margins to Small Covering Numbers

Sepp Hochreiter, Klaus Obermayer

We investigate the problem of learning a classification task for datasets which are described by matrices. Rows and columns of these matrices correspond to objects, where row and column objects may belong to different sets, and the entries in the matrix express the relationships between them. We interpret the matrix elements as being produced by an unknown kernel which operates on object pairs and we show that - under mild assumptions - these kernels correspond to dot products in some (unknown) feature space. Minimizing a bound for the generalization error of a linear classifier which has been obtained using covering numbers we derive an objective function for model selection according to the principle of structural risk minimization. The new objective function has the advantage that it allows the analysis of matrices which are not positive definite, and not even symmetric or square. We then consider the case that row objects are interpreted as features. We suggest an additional constraint, which imposes sparseness on the row objects and show, that the method can then be used for feature selection. Finally, we apply this method to data obtained from DNA microarrays, where "column" objects correspond to samples, "row" objects correspond to genes and matrix elements correspond to expression levels. Benchmarks are conducted using standard one-gene classification and support vector machines and K-nearest neighbors after standard feature selection. Our new method extracts a sparse set of genes and provides superior classification results.

\*\*\*\*\*

Learning Semantic Similarity

Jaz Kandola, Nello Cristianini, John Shawe-taylor

The standard representation of text documents as bags of words suffers from well known limitations, mostly due to its inability to exploit semantic similarity between terms. Attempts to incorporate some notion of term similarity include latent semantic indexing [8], the use of semantic networks [9], and probabilistic methods [5]. In this paper we propose two methods for inferring semantic similarity from a corpus. The first one defines word-similarity based on document-similarity and viceversa, giving rise to a system of equations whose equilibrium point we use to obtain a semantic similarity measure. The second method models semantic relations by means of a diffusion process on a graph defined by lexicon and co-occurrence information. Both approaches produce valid kernel functions parametrised by a real number. The paper shows how the alignment measure can be used to successfully perform model selection over this parameter. Combined with the use of support vector machines we obtain positive results.

\*\*\*\*\*

How Linear are Auditory Cortical Responses?

Maneesh Sahani, Jennifer Linden



By comparison to some other sensory cortices, the functional properties of cells in the primary auditory cortex are not yet well understood. Recent attempts to obtain a generalized description of auditory cortical responses have often relied upon characterization of the spectrotemporal receptive field (STRF), which amounts to a model of the stimulus-response function (SRF) that is linear in the spectrogram of the stimulus. How well can such a model account for neural responses at the very first stages of auditory cortical processing? To answer this question, we develop a novel methodology for evaluating the fraction of stimulus-related response power in a population that can be captured by a given type of SRF model. We use this technique to show that, in the thalamo-recipient layers of primary auditory cortex, STRF models account for no more than 40% of the stimulus-related power in neural responses.

\*\*\*\*\*

#### Incremental Gaussian Processes

Joaquin Quiñero Candela, Ole Winther

In this paper, we consider Tipping's relevance vector machine (RVM) [1] and formalize an incremental training strategy as a variant of the expectation-maximization (EM) algorithm that we call Subspace EM (SSEM). Working with a subset of active basis functions, the sparsity of the RVM solution will ensure that the number of basis functions and thereby the computational complexity is kept low. We also introduce a mean field approach to the intractable classification model that is expected to give a very good approximation to exact Bayesian inference and contains the Laplace approximation as a special case. We test the algorithms on two large data sets with  $O(10^3)$  (cid:0)  $10^4$  examples. The results indicate that Bayesian learning of large data sets, e.g. the MNIST database is realistic.

\*\*\*\*\*

#### Dynamic Bayesian Networks with Deterministic Latent Tables

David Barber

The application of latent/hidden variable Dynamic Bayesian Networks is constrained by the complexity of marginalising over latent variables. For this reason either small latent dimensions or Gaussian latent conditional tables linearly dependent on past states are typically considered in order that inference is tractable. We suggest an alternative approach in which the latent variables are modelled using deterministic conditional probability tables. This specialisation has the advantage of tractable inference even for highly complex non-linear/non-Gaussian visible conditional probability tables. This approach enables the consideration of highly complex latent dynamics whilst retaining the benefits of a tractable probabilistic model.

\*\*\*\*\*

#### Generalized<sup>2</sup> Linear<sup>2</sup> Models

Geoffrey J. Gordon

We introduce the Generalized<sup>2</sup> Linear<sup>2</sup> Model, a statistical estimator which combines features of nonlinear regression and factor analysis.

A (GL)<sup>2</sup>M approximately decomposes a rectangular matrix  $X$  into a simpler representation  $j(g(A)h(B))$ . Here  $A$  and  $B$  are low-rank matrices, while  $j$ ,  $g$ , and  $h$  are link functions. (GL)<sup>2</sup>Ms include many useful models as special cases, including principal components analysis, exponential-family peA, the infomax formulation of independent components analysis, linear regression, and generalized linear models. They also include new and interesting special cases, one of which we describe below. We also present an iterative procedure which optimizes the parameters of a (GL)<sup>2</sup>M. This procedure reduces to well-known algorithms for some of the special cases listed above; for other special cases, it is new.

\*\*\*\*\*

#### Spectro-Temporal Receptive Fields of Subthreshold Responses in Auditory Cortex

Christian K. Machens, Michael Wehr, Anthony Zador

How do cortical neurons represent the acoustic environment? This question is often addressed by probing with simple stimuli such as clicks or tone pips. Such stimuli have the advantage of yielding easily interpreted answers, but have the disadvantage that they may fail to uncover complex or higher-order neuronal responses.

onse properties. Here we adopt an alternative approach, probing neuronal responses with complex acoustic stimuli, including animal vocalizations and music. We have used in vivo whole cell methods in the rat auditory cortex to record subthreshold membrane potential fluctuations elicited by these stimuli. Whole cell recording reveals the total synaptic input to a neuron from all the other neurons in the circuit, instead of just its output—a sparse binary spike train—as in conventional single unit physiological recordings. Whole cell recording thus provides a much richer source of information about the neuron’s response. Many neurons responded robustly and reliably to the complex stimuli in our ensemble. Here we analyze the linear component—the spectro-temporal receptive field (STRF)—of the transformation from the sound (as represented by its time-varying spectrogram) to the neuron’s membrane potential. We find that the STRF has a rich dynamical structure, including excitatory regions positioned in general accord with the prediction of the simple tuning curve. We also find that in many cases, much of the neuron’s response, although deterministically related to the stimulus, cannot be predicted by the linear component, indicating the presence of as-yet-uncharacterized nonlinear response properties.

\*\*\*\*\*

#### Bayesian Models of Inductive Generalization

Neville Sanjana, Joshua Tenenbaum

We argue that human inductive generalization is best explained in a Bayesian framework, rather than by traditional models based on similarity computations. We go beyond previous work on Bayesian concept learning by introducing an unsupervised method for constructing flexible hypothesis spaces, and we propose a version of the Bayesian Occam’s razor that trades off priors and likelihoods to prevent under- or over-generalization in these flexible spaces. We analyze two published data sets on inductive reasoning as well as the results of a new behavioral study that we have carried out.

\*\*\*\*\*

#### Dynamical Constraints on Computing with Spike Timing in the Cortex

Arunava Banerjee, Alexandre Pouget

If the cortex uses spike timing to compute, the timing of the spikes must be robust to perturbations. Based on a recent framework that provides a simple criterion to determine whether a spike sequence produced by a generic network is sensitive to initial conditions, and numerical simulations of a variety of network architectures, we argue within the limits set by our model of the neuron, that it is unlikely that precise sequences of spike timings are used for computation under conditions typically found in the cortex.

\*\*\*\*\*

#### Bayesian Estimation of Time-Frequency Coefficients for Audio Signal Enhancement

Patrick Wolfe, Simon Godsill

The Bayesian paradigm provides a natural and effective means of exploiting prior knowledge concerning the time-frequency structure of sound signals such as speech and music—something which has often been overlooked in traditional audio signal processing approaches. Here, after constructing a Bayesian model and prior distributions capable of taking into account the time-frequency characteristics of typical audio waveforms, we apply Markov chain Monte Carlo methods in order to sample from the resultant posterior distribution of interest. We present speech enhancement results which compare favourably in objective terms with standard time-varying filtering techniques (and in several cases yield superior performance, both objectively and subjectively); moreover, in contrast to such methods, our results are obtained without an assumption of prior knowledge of the noise power.

\*\*\*\*\*

#### Real-Time Monitoring of Complex Industrial Processes with Particle Filters

Rubén Morales-Menéndez, Nando de Freitas, David Poole

This paper discusses the application of particle filtering algorithms to fault diagnosis in complex industrial processes. We consider two ubiquitous processes: an industrial dryer and a level tank. For these applications, we compared three

particle filtering variants: standard particle filtering, Rao-Blackwellised particle filtering and a version of Rao-Blackwellised particle filtering that does one-step look-ahead to select good sampling regions. We show that the overhead of the extra processing per particle of the more sophisticated methods is more than compensated by the decrease in error and variance.

\*\*\*\*\*

#### Retinal Processing Emulation in a Programmable 2-Layer Analog Array Processor CMOS Chip

R. Carmona, F. Jiménez-garrido, R. Dominguez-castro, S. Espejo, A. Rodríguez-vázquez

A bio-inspired model for an analog programmable array processor (APAP), based on studies on the vertebrate retina, has permitted the realization of complex programmable spatio-temporal dynamics in VLSI. This model mimics the way in which images are processed in the visual pathway, rendering a feasible alternative for the implementation of early vision applications in standard technologies. A prototype chip has been designed and fabricated in a 0.5( $\mu$ m) standard CMOS process. Computing power per area and power consumption is amongst the highest reported for a single chip. Design challenges, trade-offs and some experimental results are presented in this paper.

\*\*\*\*\*

#### Bayesian Image Super-Resolution

Michael Tipping, Christopher Bishop

The extraction of a single high-quality image from a set of low( $\times$ 3) resolution images is an important problem which arises in fields such as remote sensing, surveillance, medical imaging and the extraction of still images from video. Typical approaches are based on the use of cross-correlation to register the images followed by the inversion of the transformation from the unknown high resolution image to the observed low resolution images, using regularization to resolve the ill-posed nature of the inversion process. In this paper we develop a Bayesian treatment of the super-resolution problem in which the likelihood function for the image registration parameters is based on a marginalization over the unknown high-resolution image. This approach allows us to estimate the unknown point spread function, and is rendered tractable through the introduction of a Gaussian process prior over images. Results indicate a significant improvement over techniques based on MAP (maximum a-posteriori) point optimization of the high resolution image and associated registration parameters.

\*\*\*\*\*

#### Charting a Manifold

Matthew Brand

We construct a nonlinear mapping from a high-dimensional sample space to a low-dimensional vector space, effectively recovering a Cartesian coordinate system for the manifold from which the data is sampled. The mapping preserves local geometric relations in the manifold and is pseudo-invertible. We show how to estimate the intrinsic dimensionality of the manifold from samples, decompose the sample data into locally linear low-dimensional patches, merge these patches into a single low-dimensional coordinate system, and compute forward and reverse mappings between the sample and coordinate spaces. The objective functions are convex and their solutions are given in closed form.

\*\*\*\*\*

#### A Minimal Intervention Principle for Coordinated Movement

Emanuel Todorov, Michael Jordan

Behavioral goals are achieved reliably and repeatedly with movements rarely reproducible in their detail. Here we offer an explanation: we show that not only are variability and goal achievement compatible, but indeed that allowing variability in redundant dimensions is the optimal control strategy in the face of uncertainty. The optimal feedback control laws for typical motor tasks obey a "minimal intervention" principle: deviations from the average trajectory are only corrected when they interfere with the task goals. The resulting behavior exhibits ta

sk-constrained variability, as well as synergetic coupling among actuators—which is another unexplained empirical phenomenon.

\*\*\*\*\*

#### A Probabilistic Approach to Single Channel Blind Signal Separation

Gil-jin Jang, Te-Won Lee

We present a new technique for achieving source separation when given only a single channel recording. The main idea is based on exploiting the inherent time structure of sound sources by learning a priori sets of basis filters in time domain that encode the sources in a statistically efficient manner. We derive a learning algorithm using a maximum likelihood approach given the observed single channel data and sets of basis filters. For each time point we infer the source signals and their contribution factors. This inference is possible due to the prior knowledge of the basis filters and the associated coefficient densities. A flexible model for density estimation allows accurate modeling of the observation and our experimental results exhibit a high level of separation performance for mixtures of two music signals as well as the separation of two voice signals.

\*\*\*\*\*

#### A Digital Antennal Lobe for Pattern Equalization: Analysis and Design

Alex Holub, Gilles Laurent, Pietro Perona

Re-mapping patterns in order to equalize their distribution may greatly simplify both the structure and the training of classifiers. Here, the properties of one such map obtained by running a few steps of discrete-time dynamical system are explored. The system is called 'Digital Antennal Lobe' (DAL) because it is inspired by recent studies of the antennal lobe, a structure in the olfactory system of the grasshopper. The pattern-spreading properties of the DAL as well as its average behavior as a function of its (few) design parameters are analyzed by extending previous results of Van Vreeswijk and Sompolinsky. Furthermore, a technique for adapting the parameters of the initial design in order to obtain opportune noise-rejection behavior is suggested. Our results are demonstrated with a number of simulations.

\*\*\*\*\*

#### Knowledge-Based Support Vector Machine Classifiers

Glenn Fung, Olvi Mangasarian, Jude Shavlik

Prior knowledge in the form of multiple polyhedral sets, each belonging to one of two categories, is introduced into a reformulation of a linear support vector machine classifier. The resulting formulation leads to a linear program that can be solved efficiently. Real world examples, from DNA sequencing and breast cancer prognosis, demonstrate the effectiveness of the proposed method. Numerical results show improvement in test set accuracy after the incorporation of prior knowledge into ordinary, data-based linear support vector machine classifiers. One experiment also shows that a linear classifier, based solely on prior knowledge, far outperforms the direct application of prior knowledge rules to classified data. Keywords: use and refinement of prior knowledge, support vector machines, linear programming

\*\*\*\*\*

#### Conditional Models on the Ranking Poset

Guy Lebanon, John Lafferty

A distance-based conditional model on the ranking poset is presented for use in classification and ranking. The model is an extension of the Mallows model, and generalizes the classifier combination methods used by several ensemble learning algorithms, including error correcting output codes, discrete AdaBoost, logistic regression and cranking. The algebraic structure of the ranking poset leads to a simple Bayesian interpretation of the conditional model and its special cases. In addition to a unifying view, the framework suggests a probabilistic interpretation for error correcting output codes and an extension beyond the binary coding scheme.

\*\*\*\*\*

#### A Formulation for Minimax Probability Machine Regression

Thomas Strohmann, Gregory Grudic

We formulate the regression problem as one of maximizing the minimum probability, symbolized by (cid:10), that future predicted outputs of the regression model will be within some (cid:6)" bound of the true regression function. Our formulation is unique in that we obtain a direct estimate of this lower probability bound (cid:10). The proposed framework, minimax probability machine regression (MPMR), is based on the recently described minimax probability machine classification algorithm [Lanckriet et al.] and uses Mercer Kernels to obtain nonlinear regression models. MPMR is tested on both toy and real world data, verifying the accuracy of the (cid:10) bound, and the efficacy of the regression models.

\*\*\*\*\*

Multiclass Learning by Probabilistic Embeddings

Ofer Dekel, Yoram Singer

We describe a new algorithmic framework for learning multiclass categorization problems. In this framework a multiclass predictor is composed of a pair of embeddings that map both instances and labels into a common space. In this space each instance is assigned the label it is nearest to. We outline and analyze an algorithm, termed Bunching, for learning the pair of embeddings from labeled data.

A key construction in the analysis of the algorithm is the notion of probabilistic output codes, a generalization of error correcting output codes (ECOC). Furthermore, the method of multiclass categorization using ECOC is shown to be an instance of Bunching. We demonstrate the advantage of Bunching over ECOC by comparing their performance on numerous categorization problems.

\*\*\*\*\*

Improving a Page Classifier with Anchor Extraction and Link Analysis

William W. Cohen

Most text categorization systems use simple models of documents and document collections. In this paper we describe a technique that improves a simple web page classifier's performance on pages from a new, unseen web site, by exploiting link structure within a site as well as page structure within hub pages. On real-world test cases, this technique significantly and substantially improves the accuracy of a bag-of-words classifier, reducing error rate by about half, on average.

The system uses a variant of co-training to exploit unlabeled data from a new site. Pages are labeled using the base classifier; the results are used by a restricted wrapper-learner to propose potential "main-category anchor wrappers"; and finally, these wrappers are used as features by a third learner to find a categorization of the site that implies a simple hub structure, but which also largely agrees with the original bag-of-words classifier.

\*\*\*\*\*

VIBES: A Variational Inference Engine for Bayesian Networks

Christopher Bishop, David Spiegelhalter, John Winn

In recent years variational methods have become a popular tool for approximate inference and learning in a wide variety of probabilistic models. For each new application, however, it is currently necessary (cid:12)rst to derive the variational update equations, and then to implement them in application-specific (cid:12)c code. Each of these steps is both time consuming and error prone. In this paper we describe a general purpose inference engine called VIBES ('Variational Inference for Bayesian Networks') which allows a wide variety of probabilistic models to be implemented and solved variationally without recourse to coding. New models are specified (cid:12)ed either through a simple script or via a graphical interface analogous to a drawing package. VIBES then automatically generates and solves the variational equations. We illustrate the power and (cid:13)exibility of VIBES using examples from Bayesian mixture modelling.

\*\*\*\*\*

A Convergent Form of Approximate Policy Iteration

Theodore Perkins, Doina Precup

We study a new, model-free form of approximate policy iteration which uses Sarsa updates with linear state-action value function approximation for policy evaluation, and a "policy improvement operator" to generate a new policy based on the learned state-action values. We prove that if the policy improvement operator pr

duces  $\epsilon$ -soft policies and is Lipschitz continuous in the action values, with a constant that is not too large, then the approximate policy iteration algorithm converges to a unique solution from any initial policy. To our knowledge, this is the first convergence result for any form of approximate policy iteration under similar computational-resource assumptions.

\*\*\*\*\*

#### Mean Field Approach to a Probabilistic Model in Information Retrieval

Bin Wu, K. Wong, David Bodoff

We study an explicit parametric model of documents, queries, and relevance assessment for Information Retrieval (IR). Mean-field methods are applied to analyze the model and derive efficient practical algorithms to estimate the parameters in the problem. The hyperparameters are estimated by a fast approximate leave-one-out cross-validation procedure based on the cavity method. The algorithm is further evaluated on several benchmark databases by comparing with standard algorithms in IR.

\*\*\*\*\*

#### Exponential Family PCA for Belief Compression in POMDPs

Nicholas Roy, Geoffrey J. Gordon

Geoffrey Gordon

\*\*\*\*\*

#### Dynamic Structure Super-Resolution

Amos J. Storkey

The problem of super-resolution involves generating feasible higher resolution images, which are pleasing to the eye and realistic, from a given low resolution image. This might be attempted by using simple (cid:12)lters for smoothing out the high resolution blocks or through applications where substantial prior information is used to imply the textures and shapes which will occur in the images.

In this paper we describe an approach which lies between the two extremes. It is a generic unsupervised method which is usable in all domains, but goes beyond simple smoothing methods in what it achieves. We use a dynamic tree-like architecture to model the high resolution data. Approximate conditioning on the low resolution image is achieved through a mean (cid:12)eld approach.

\*\*\*\*\*

#### Independent Components Analysis through Product Density Estimation

Trevor Hastie, Rob Tibshirani

We present a simple direct approach for solving the ICA problem, using density estimation and maximum likelihood. Given a candidate orthogonal frame, we model each of the coordinates using a semi-parametric density estimate based on cubic splines. Since our estimates have two continuous derivatives, we can easily run a second order search for the frame parameters. Our method performs very favorably when compared to state-of-the-art techniques.

\*\*\*\*\*

#### Discriminative Densities from Maximum Contrast Estimation

Peter Meinicke, Thorsten Twellmann, Helge Ritter

We propose a framework for classifier design based on discriminative densities for representation of the differences of the class-conditional distributions in a way that is optimal for classification. The densities are selected from a parametrized set by constrained maximization of some objective function which measures the average (bounded) difference, i.e. the contrast between discriminative densities. We show that maximization of the contrast is equivalent to minimization of an approximation of the Bayes risk. Therefore using suitable classes of probability density functions, the resulting maximum contrast classifiers (MCCs) can approximate the Bayes rule for the general multiclass case. In particular for a certain parametrization of the density functions we obtain MCCs which have the same functional form as the well-known Support Vector Machines (SVMs). We show that MCC-training in general requires some nonlinear optimization but under certain conditions the problem is concave and can be tackled by a single linear program. We indicate the close relation between SVM- and MCC-training and in particular we show that Linear Programming Machines can be viewed as an approximate

e realization of MCCs. In the experiments on benchmark data sets, the MCC shows a competitive classification performance.

\*\*\*\*\*

#### Interpreting Neural Response Variability as Monte Carlo Sampling of the Posterior

Patrik Hoyer, Aapo Hyvärinen

The responses of cortical sensory neurons are notoriously variable, with the number of spikes evoked by identical stimuli varying significantly from trial to trial. This variability is most often interpreted as 'noise', purely detrimental to the sensory system. In this paper, we propose an alternative view in which the variability is related to the uncertainty, about world parameters, which is inherent in the sensory stimulus. Specifically, the responses of a population of neurons are interpreted as stochastic samples from the posterior distribution in a latent variable model. In addition to giving theoretical arguments supporting such a representational scheme, we provide simulations suggesting how some aspects of response variability might be understood in this framework.

\*\*\*\*\*

#### Combining Features for BCI

Guido Dornhege, Benjamin Blankertz, Gabriel Curio, Klaus-Robert Müller

Recently, interest is growing to develop an effective communication interface connecting the human brain to a computer, the 'Brain-Computer Interface' (BCI). One motivation of BCI research is to provide a new communication channel substituting normal motor output in patients with severe neuromuscular disabilities. In the last decade, various neurophysiological cortical processes, such as slow potential shifts, movement related potentials (MRPs) or event-related desynchronization (ERD) of spontaneous EEG rhythms, were shown to be suitable for BCI, and, consequently, different independent approaches of extracting BCI-relevant EEG-features for single-trial analysis are under investigation. Here, we present a method and systematically compare several concepts for combining such EEG-features to improve the single-trial classification. Feature combinations are evaluated on movement imagination experiments with 3 subjects where EEG-features are based on either MRPs or ERD, or both. Those combination methods that incorporate the assumption that the single EEG-features are physiologically mutually independent outperform the plain method of 'adding' evidence where the single-feature vectors are simply concatenated. These results strengthen the hypothesis that MRP and ERD reflect at least partially independent aspects of cortical processes and open a new perspective to boost BCI effectiveness.

\*\*\*\*\*

#### A Probabilistic Model for Learning Concatenative Morphology

Matthew Snover, Michael Brent

This paper describes a system for the unsupervised learning of morphological suffixes and stems from word lists. The system is composed of a generative probability model and hill-climbing and directed search algorithms. By extracting and examining morphologically rich subsets of an input lexicon, the directed search identifies highly productive paradigms. The hill-climbing algorithm then further maximizes the probability of the hypothesis. Quantitative results are shown by measuring the accuracy of the morphological relations identified. Experiments in English and Polish, as well as comparisons with another recent unsupervised morphology learning algorithm demonstrate the effectiveness of this technique.

\*\*\*\*\*

#### The Decision List Machine

Marina Sokolova, Mario Marchand, Nathalie Japkowicz, John Shawe-taylor

We introduce a new learning algorithm for decision lists to allow features that are constructed from the data and to allow a trade-off between accuracy and complexity. We bound its generalization error in terms of the number of errors and the size of the classifier it finds on the training data. We also compare its performance on some natural data sets with the set covering machine and the support vector machine.

\*\*\*\*\*

#### How to Combine Color and Shape Information for 3D Object Recognition: Kernels do

the Trick

B. Caputo, Gy. Dorkó

This paper presents a kernel method that allows to combine color and shape information for appearance-based object recognition. It doesn't require to define a new common representation, but use the power of kernels to combine different representations together in an effective manner. These results are achieved using results of statistical mechanics of spin glasses combined with Markov random fields via kernel functions. Experiments show an increase in recognition rate up to 5.92% with respect to conventional strategies.

\*\*\*\*\*

Reconstructing Stimulus-Driven Neural Networks from Spike Times

Duane Nykamp

We present a method to distinguish direct connections between two neurons from common input originating from other, unmeasured neurons. The distinction is computed from the spike times of the two neurons in response to a white noise stimulus. Although the method is based on a highly idealized linear-nonlinear approximation of neural response, we demonstrate via simulation that the approach can work with a more realistic, integrate-and-fire neuron model. We propose that the approach exemplified by this analysis may yield viable tools for reconstructing stimulus-driven neural networks from data gathered in neurophysiology experiments.

\*\*\*\*\*

Selectivity and Metaplasticity in a Unified Calcium-Dependent Model

Luk Chong Yeung, Brian Blais, Leon Cooper, Harel Shouval

A unified, biophysically motivated Calcium-Dependent Learning model has been shown to account for various rate-based and spike time-dependent paradigms for inducing synaptic plasticity. Here, we investigate the properties of this model for a multi-synapse neuron that receives inputs with different spike-train statistics. In addition, we present a physiological form of metaplasticity, an activity-driven regulation mechanism, that is essential for the robustness of the model. A neuron thus implemented develops stable and selective receptive fields, given various input statistics

\*\*\*\*\*

Learning to Detect Natural Image Boundaries Using Brightness and Texture

David Martin, Charles Fowlkes, Jitendra Malik

The goal of this work is to accurately detect and localize boundaries in natural scenes using local image measurements. We formulate features that respond to characteristic changes in brightness and texture associated with natural boundaries. In order to combine the information from these features in an optimal way, a classifier is trained using human labeled images as ground truth. We present precision-recall curves showing that the resulting detector outperforms existing approaches.

\*\*\*\*\*

Concurrent Object Recognition and Segmentation by Graph Partitioning

Stella X. Yu, Ralph Gross, Jianbo Shi

Segmentation and recognition have long been treated as two separate processes. We propose a mechanism based on spectral graph partitioning that readily combine the two processes into one. A part-based recognition system detects object patches, supplies their partial segmentations as well as knowledge about the spatial configurations of the object. The goal of patch grouping is to find a set of patches that conform best to the object configuration, while the goal of pixel grouping is to find a set of pixels that have the best low-level feature similarity. Through pixel-patch interactions and between-patch competition encoded in the solution space, these two processes are realized in one joint optimization problem. The globally optimal partition is obtained by solving a constrained eigenvalue problem. We demonstrate that the resulting object segmentation eliminates false positives for the part detection, while overcoming occlusion and weak contours for the low-level edge detection.

\*\*\*\*\*



## Approximate Linear Programming for Average-Cost Dynamic Programming

Benjamin Roy, Daniela Farias

This paper extends our earlier analysis on approximate linear programming as an approach to approximating the cost-to-go function in a discounted-cost dynamic program [6]. In this paper, we consider the average-cost criterion and a version of approximate linear programming that generates approximations to the optimal average cost and differential cost function. We demonstrate that a naive version of approximate linear programming prioritizes approximation of the optimal average cost and that this may not be well-aligned with the objective of deriving a policy with low average cost. For that, the algorithm should aim at producing a good approximation of the differential cost function. We propose a two-phase variant of approximate linear programming that allows for external control of the relative accuracy of the approximation of the differential cost function over different portions of the state space via state-relevance weights. Performance bounds suggest that the new algorithm is compatible with the objective of optimizing performance and provide guidance on appropriate choices for state-relevance weights.

\*\*\*\*\*

## Learning Sparse Topographic Representations with Products of Student-t Distributions

Max Welling, Simon Osindero, Geoffrey E. Hinton

We propose a model for natural images in which the probability of an image is proportional to the product of the probabilities of some filter outputs. We encourage the system to find sparse features by using a Student-t distribution to model each filter output. If the t-distribution is used to model the combined outputs of sets of neurally adjacent filters, the system learns a topographic map in which the orientation, spatial frequency and location of the filters change smoothly across the map. Even though maximum likelihood learning is intractable in our model, the product form allows a relatively efficient learning procedure that works well even for highly overcomplete sets of filters. Once the model has been learned it can be used as a prior to derive the "iterated Wiener filter" for the purpose of denoising images.

\*\*\*\*\*

## Margin Analysis of the LVQ Algorithm

Koby Crammer, Ran Gilad-bachrach, Amir Navot, Naftali Tishby

Prototypes based algorithms are commonly used to reduce the computational complexity of Nearest-Neighbour (NN) classifiers. In this paper we discuss theoretical and algorithmical aspects of such algorithms. On the theory side, we present margin based generalization bounds that suggest that these kinds of classifiers can be more accurate than the 1-NN rule. Furthermore, we derived a training algorithm that selects a good set of prototypes using large margin principles. We also show that the 20 years old Learning Vector Quantization (LVQ) algorithm emerges naturally from our framework.

\*\*\*\*\*

## Feature Selection by Maximum Marginal Diversity

Nuno Vasconcelos

We address the question of feature selection in the context of visual recognition. It is shown that, besides efficient from a computational standpoint, the infomax principle is nearly optimal in the minimum Bayes error sense. The concept of marginal diversity is introduced, leading to a generic principle for feature selection (the principle of maximum marginal diversity) of extreme computational simplicity. The relationships between infomax and the maximization of marginal diversity are identified, uncovering the existence of a family of classification procedures for which near optimal (in the Bayes error sense) feature selection does not require combinatorial search. Examination of this family in light of recent studies on the statistics of natural images suggests that visual recognition problems are a subset of it.

\*\*\*\*\*

## A Neural Edge-Detection Model for Enhanced Auditory Sensitivity in Modulated Noise

Alon Fishbach, Bradford May

Psychophysical data suggest that temporal modulations of stimulus amplitude envelopes play a prominent role in the perceptual segregation of concurrent sounds. In particular, the detection of an unmodulated signal can be significantly improved by adding amplitude modulation to the spectral envelope of a competing masking noise. This perceptual phenomenon is known as "Comodulation Masking Release" (CMR). Despite the obvious influence of temporal structure on the perception of complex auditory scenes, the physiological mechanisms that contribute to CMR and auditory streaming are not well known. A recent physiological study by Nelken and colleagues has demonstrated an enhanced cortical representation of auditory signals in modulated noise. Our study evaluates these CMR-like response patterns from the perspective of a hypothetical auditory edge-detection neuron. It is shown that this simple neural model for the detection of amplitude transients can reproduce not only the physiological data of Nelken et al., but also, in light of previous results, a variety of physiological and psychoacoustical phenomena that are related to the perceptual segregation of concurrent sounds.

\*\*\*\*\*

Concentration Inequalities for the Missing Mass and for Histogram Rule Error

Luis E. Ortiz, David McAllester

This paper gives distribution-free concentration inequalities for the missing mass and the error rate of histogram rules. Negative association methods can be used to reduce these concentration problems to concentration questions about independent sums. Although the sums are independent, they are highly heterogeneous. Such highly heterogeneous independent sums cannot be analyzed using standard concentration inequalities such as Hoeffding's inequality, the Azuma-Hoeffding bound, Bernstein's inequality, Bennett's inequality, or McDiarmid's theorem.

\*\*\*\*\*

Learning Sparse Multiscale Image Representations

Phil Sallee, Bruno Olshausen

We describe a method for learning sparse multiscale image representations using a sparse prior distribution over the basis function coefficients. The prior consists of a mixture of a Gaussian and a Dirac delta function, and thus encourages coefficients to have exact zero values. Coefficients for an image are computed by sampling from the resulting posterior distribution with a Gibbs sampler. The learned basis is similar to the Steerable Pyramid basis, and yields slightly higher SNR for the same number of active coefficients. Denoising using the learned image model is demonstrated for some standard test images, with results that compare favorably with other denoising methods.

\*\*\*\*\*

The Effect of Singularities in a Learning Machine when the True Parameters Do Not Lie on such Singularities

Sumio Watanabe, Shun-ichi Amari

A lot of learning machines with hidden variables used in information science have singularities in their parameter spaces. At singularities, the Fisher information matrix becomes degenerate, resulting that the learning theory of regular statistical models does not hold. Recently, it was proven that, if the true parameter is contained in singularities, then the coefficient of the Bayes generalization error is equal to the pole of the zeta function of the Kullback information. In this paper, under the condition that the true parameter is almost but not contained in singularities, we show two results. (1) If the dimension of the parameter from inputs to hidden units is not larger than three, then there exists a region of true parameters where the generalization error is larger than those of regular models, however, if otherwise, then for any true parameter, the generalization error is smaller than those of regular models. (2) The symmetry of the generalization error and the training error does not hold in singular models in general.

\*\*\*\*\*

Learning About Multiple Objects in Images: Factorial Learning without Factorial Search

Christopher K. I. Williams, Michalis Titsias

We consider data which are images containing views of multiple objects. Our task is to learn about each of the objects present in the images. This task can be approached as a factorial learning problem, where each image must be explained by instantiating a model for each of the objects present with the correct instantiation parameters. A major problem with learning a factorial model is that as the number of objects increases, there is a combinatorial explosion of the number of configurations that need to be considered. We develop a method to extract object models sequentially from the data by making use of a robust statistical method, thus avoiding the combinatorial explosion, and present results showing successful extraction of objects from real images.

\*\*\*\*\*

Distance Metric Learning with Application to Clustering with Side-Information

Eric Xing, Michael Jordan, Stuart J. Russell, Andrew Ng

@cs.berkeley.edu

\*\*\*\*\*

Location Estimation with a Differential Update Network

Ali Rahimi, Trevor Darrell

Given a set of hidden variables with an a-priori Markov structure, we derive an online algorithm which approximately updates the posterior as pairwise measurements between the hidden variables become available. The update is performed using Assumed Density Filtering: to incorporate each pairwise measurement, we compute the optimal Markov structure which represents the true posterior and use it as a prior for incorporating the next measurement. We demonstrate the resulting algorithm by calculating globally consistent trajectories of a robot as it navigates along a 2D trajectory. To update a trajectory of length  $t$ , the update takes  $O(t)$ . When all conditional distributions are linear-Gaussian, the algorithm can be thought of as a Kalman Filter which simplifies the state covariance matrix after incorporating each measurement.

\*\*\*\*\*

Hyperkernels

Cheng Ong, Robert C. Williamson, Alex Smola

We consider the problem of choosing a kernel suitable for estimation using a Gaussian Process estimator or a Support Vector Machine. A novel solution is presented which involves defining a Reproducing Kernel Hilbert Space on the space of kernels itself. By utilizing an analog of the classical representer theorem, the problem of choosing a kernel from a parameterized family of kernels (e.g. of varying width) is reduced to a statistical estimation problem akin to the problem of minimizing a regularized risk functional. Various classical settings for model or kernel selection are special cases of our framework.

\*\*\*\*\*

Prediction and Semantic Association

Thomas Griffiths, Mark Steyvers

We explore the consequences of viewing semantic association as the result of attempting to predict the concepts likely to arise in a particular context. We argue that the success of existing accounts of semantic representation comes as a result of indirectly addressing this problem, and show that a closer correspondence to human data can be obtained by taking a probabilistic approach that explicitly models the generative structure of language.

\*\*\*\*\*

A Maximum Entropy Approach to Collaborative Filtering in Dynamic, Sparse, High-Dimensional Domains

Dmitry Pavlov, David Pennock

We develop a maximum entropy (maxent) approach to generating recommendations in the context of a user's current navigation stream, suitable for environments where data is sparse, high-dimensional, and dynamic—conditions typical of many recommendation applications. We address sparsity and dimensionality reduction by first clustering items based on user access patterns so as to attempt to minimize the a priori probability that recommendations will cross cluster boundaries and then recommending only within clusters. We address the inherent dynamic nature

re of the problem by explicitly modeling the data as a time series; we show how this representational expressivity fits naturally into a maxent framework. We conduct experiments on data from ResearchIndex, a popular online repository of over 470,000 computer science documents. We show that our maxent formulation outperforms several competing algorithms in offline tests simulating the recommendation of documents to ResearchIndex users.

\*\*\*\*\*

#### Annealing and the Rate Distortion Problem

Albert Parker, Tomáš Gedeon, Alexander Dimitrov

In this paper we introduce methodology to determine the bifurcation structure of optima for a class of similar cost functions from Rate Distortion Theory, Deterministic Annealing, Information Distortion and the Information Bottleneck Method. We also introduce a numerical algorithm which uses the explicit form of the bifurcating branches to find optima at a bifurcation point.

\*\*\*\*\*

#### Self Supervised Boosting

Max Welling, Richard Zemel, Geoffrey E. Hinton

Boosting algorithms and successful applications thereof abound for classification and regression learning problems, but not for unsupervised learning. We propose a sequential approach to adding features to a random field model by training them to improve classification performance between the data and an equal-sized sample of "negative examples" generated from the model's current estimate of the data density. Training in each boosting round proceeds in three stages: First we sample negative examples from the model's current Boltzmann distribution. Next, a feature is trained to improve classification performance between data and negative examples. Finally, a coefficient is learned which determines the importance of this feature relative to ones already in the pool. Negative examples only need to be generated once to learn each new feature. The validity of the approach is demonstrated on binary digits and continuous synthetic data.

\*\*\*\*\*

#### An Estimation-Theoretic Framework for the Presentation of Multiple Stimuli

Christian Furich

A framework is introduced for assessing the encoding accuracy and the discriminational ability of a population of neurons upon simultaneous presentation of multiple stimuli. Minimal square estimation errors are obtained from a Fisher information analysis in an abstract compound space comprising the features of all stimuli. Even for the simplest case of linear superposition of responses and Gaussian tuning, the symmetries in the compound space are very different from those in the case of a single stimulus. The analysis allows for a quantitative description of attentional effects and can be extended to include neural nonlinearities such as nonclassical receptive fields.

\*\*\*\*\*

#### Convergent Combinations of Reinforcement Learning with Linear Function Approximation

Ralf Schoknecht, Artur Merke

Convergence for iterative reinforcement learning algorithms like TD(0) depends on the sampling strategy for the transitions. However, in practical applications it is convenient to take transition data from arbitrary sources without losing convergence. In this paper we investigate the problem of repeated synchronous updates based on a fixed set of transitions. Our main theorem yields sufficient conditions of convergence for combinations of reinforcement learning algorithms and linear function approximation. This allows to analyse if a certain reinforcement learning algorithm and a certain function approximator are compatible. For the combination of the residual gradient algorithm with grid-based linear interpolation we show that there exists a universal constant learning rate such that the iteration converges independently of the concrete transition data.

\*\*\*\*\*

#### Informed Projections

David Cohn

Low rank approximation techniques are widespread in pattern recognition research – they include Latent Semantic Analysis (LSA), Probabilistic LSA, Principal Components Analysis (PCA), the Generative Aspect Model, and many forms of bibliometric analysis. All make use of a low-dimensional manifold onto which data are projected. Such techniques are generally “unsupervised,” which allows them to model data in the absence of labels or categories. With many practical problems, however, some prior knowledge is available in the form of context. In this paper, I describe a principled approach to incorporating such information, and demonstrate its application to PCA-based approximations of several data sets.

\*\*\*\*\*

Automatic Acquisition and Efficient Representation of Syntactic Structures

Zach Solan, Eytan Ruppin, David Horn, Shimon Edelman

The distributional principle according to which morphemes that occur in identical contexts belong, in some sense, to the same category [1] has been advanced as a means for extracting syntactic structures from corpus data. We extend this principle by applying it recursively, and by using mutual information for estimating category coherence. The resulting model learns, in an unsupervised fashion, highly structured, distributed representations of syntactic knowledge from corpora. It also exhibits promising behavior in tasks usually thought to require representations anchored in a grammar, such as systematicity.

\*\*\*\*\*

Application of Variational Bayesian Approach to Speech Recognition

Shinji Watanabe, Yasuhiro Minami, Atsushi Nakamura, Naonori Ueda

In this paper, we propose a Bayesian framework, which constructs shared-state triphone HMMs based on a variational Bayesian approach, and recognizes speech based on the Bayesian prediction classification; variational Bayesian estimation and clustering for speech recognition (VBEC). An appropriate model structure with high recognition performance can be found within a VBEC framework. Unlike conventional methods, including BIC or MDL criterion based on the maximum likelihood approach, the proposed model selection is valid in principle, even when there are insufficient amounts of data, because it does not use an asymptotic assumption. In isolated word recognition experiments, we show the advantage of VBEC over conventional methods, especially when dealing with small amounts of data.

\*\*\*\*\*

Stable Fixed Points of Loopy Belief Propagation Are Local Minima of the Bethe Free Energy

Tom Heskes

We extend recent work on the connection between loopy belief propagation and the Bethe free energy. Constrained minimization of the Bethe free energy can be turned into an unconstrained saddle-point problem. Both converging double-loop algorithms and standard loopy belief propagation can be interpreted as attempts to solve this saddle-point problem. Stability analysis then leads us to conclude that stable fixed points of loopy belief propagation must be (local) minima of the Bethe free energy. Perhaps surprisingly, the converse need not be the case: minima can be unstable fixed points. We illustrate this with an example and discuss implications.

\*\*\*\*\*

Identity Uncertainty and Citation Matching

Hanna Pasula, Bhaskara Marthi, Brian Milch, Stuart J. Russell, Ilya Shpitser

Identity uncertainty is a pervasive problem in real-world data analysis. It arises whenever objects are not labeled with unique identifiers or when those identifiers may not be perceived perfectly. In such cases, two observations may or may not correspond to the same object. In this paper, we consider the problem in the context of citation matching—the problem of deciding which citations correspond to the same publication. Our approach is based on the use of a relational probability model to define a generative model for the domain, including models of author and title corruption and a probabilistic citation grammar. Identity uncertainty is handled by extending standard models to incorporate probabilities over

the possible mappings between terms in the language and objects in the domain. Inference is based on Markov chain Monte Carlo, augmented with specific methods for generating efficient proposals when the domain contains many objects. Results on several citation data sets show that the method outperforms current algorithms for citation matching. The declarative, relational nature of the model also means that our algorithm can determine object characteristics such as author names by combining multiple citations of multiple papers.

\*\*\*\*\*

#### Fast Sparse Gaussian Process Methods: The Informative Vector Machine

Neil Lawrence, Matthias Seeger, Ralf Herbrich

We present a framework for sparse Gaussian process (GP) methods which uses forward selection with criteria based on information-theoretic principles, previously suggested for active learning. Our goal is not only to learn  $d$  sparse predictors (which can be evaluated in  $O(d)$  rather than  $O(n)$ ,  $d \ll n$ ,  $n$  the number of training points), but also to perform training under strong restrictions on time and memory requirements. The scaling of our method is at most  $O(n^2)$ , and in large real-world classification experiments we show that it can match prediction performance of the popular support vector machine (SVM), yet can be significantly faster in training. In contrast to the SVM, our approximation produces estimates of predictive probabilities ('error bars'), allows for Bayesian model selection and is less complex in implementation.

\*\*\*\*\*

#### Recovering Articulated Model Topology from Observed Rigid Motion

Leonid Taycher, John L. III, Trevor Darrell

Accurate representation of articulated motion is a challenging problem for machine perception. Several successful tracking algorithms have been developed that model human body as an articulated tree. We propose a learning-based method for creating such articulated models from observations of multiple rigid motions. This paper is concerned with recovering topology of the articulated model, when the rigid motion of constituent segments is known. Our approach is based on finding the Maximum Likelihood tree shaped factorization of the joint probability density function (PDF) of rigid segment motions. The topology of graphical model formed from this factorization corresponds to topology of the underlying articulated body. We demonstrate the performance of our algorithm on both synthetic and real motion capture data.

\*\*\*\*\*

#### Learning to Perceive Transparency from the Statistics of Natural Scenes

Anat Levin, Assaf Zomet, Yair Weiss

Certain simple images are known to trigger a percept of transparency: the input image  $I$  is perceived as the sum of two images  $I(x; y) = I_1(x; y) + I_2(x; y)$ . This percept is puzzling. First, why do we choose the "more complicated" description with two images rather than the "simpler" explanation  $I(x; y) = I_1(x; y) + 0$ ? Second, given the infinite number of ways to express  $I$  as a sum of two images, how do we compute the "best" decomposition? Here we suggest that transparency is the rational percept of a system that is adapted to the statistics of natural scenes. We present a probabilistic model of images based on the qualitative statistics of derivative filters and "corner detectors" in natural scenes and use this model to find the most probable decomposition of a novel image. The optimization is performed using loopy belief propagation. We show that our model computes perceptually "correct" decompositions on synthetic images and discuss its application to real images.

\*\*\*\*\*

#### Replay, Repair and Consolidation

Szabolcs Káli, Peter Dayan

A standard view of memory consolidation is that episodes are stored temporarily in the hippocampus, and are transferred to the neocortex through replay. Various recent experimental challenges to the idea of transfer, particularly for human memory, are forcing its re-evaluation. However, although there is independent neurophysiological evidence for replay, short of transfer, there are few theoretical ideas for what it might be doing. We suggest and demonstrate two important

computational roles associated with neocortical indices.

\*\*\*\*\*

#### Inferring a Semantic Representation of Text via Cross-Language Correlation Analysis

Alexei Vinokourov, Nello Cristianini, John Shawe-Taylor

The problem of learning a semantic representation of a text document from data is addressed, in the situation where a corpus of unlabeled paired documents is available, each pair being formed by a short English document and its French translation. This representation can then be used for any retrieval, categorization or clustering task, both in a standard and in a cross-lingual setting. By using kernel functions, in this case simple bag-of-words inner products, each part of the corpus is mapped to a high-dimensional space. The correlations between the two spaces are then learnt by using kernel Canonical Correlation Analysis. A set of directions is found in the first and in the second space that are maximally correlated. Since we assume the two representations are completely independent apart from the semantic content, any correlation between them should reflect some semantic similarity. Certain patterns of English words that relate to a specific meaning should correlate with certain patterns of French words corresponding to the same meaning, across the corpus. Using the semantic representation obtained in this way we first demonstrate that the correlations detected between the two versions of the corpus are significantly higher than random, and hence that a representation based on such features does capture statistical patterns that should reflect semantic information. Then we use such representation both in cross-language and in single-language retrieval tasks, observing performance that is consistently and significantly superior to LSI on the same data.

\*\*\*\*\*

#### Cluster Kernels for Semi-Supervised Learning

Olivier Chapelle, Jason Weston, Bernhard Schölkopf

We propose a framework to incorporate unlabeled data in kernel classifier, based on the idea that two points in the same cluster are more likely to have the same label. This is achieved by modifying the eigenspectrum of the kernel matrix. Experimental results assess the validity of this approach.

\*\*\*\*\*

#### A Hierarchical Bayesian Markovian Model for Motifs in Biopolymer Sequences

Eric Xing, Michael Jordan, Richard Karp, Stuart J. Russell

We propose a dynamic Bayesian model for motifs in biopolymer sequences which captures rich biological prior knowledge and positional dependencies in motif structure in a principled way. Our model posits that the position-specific multinomial parameters for monomer distribution are distributed as a latent Dirichlet-mixture random variable, and the position-specific Dirichlet component is determined by a hidden Markov process. Model parameters can be fit on training motifs using a variational EM algorithm within an empirical Bayesian framework. Variational inference is also used for detecting hidden motifs. Our model improves over previous models that ignore biological priors and positional dependence. It has much higher sensitivity to motifs during detection and a notable ability to distinguish genuine motifs from false recurring patterns.

\*\*\*\*\*

#### Speeding up the Parti-Game Algorithm

Maxim Likhachev, Sven Koenig

In this paper, we introduce an efficient replanning algorithm for nonde-terministic domains, namely what we believe to be the first incremental heuristic minimax search algorithm. We apply it to the dynamic discretization of continuous domains, resulting in an efficient implementation of the parti-game reinforcement-learning algorithm for control in high-dimensional domains.

\*\*\*\*\*

#### Kernel-Based Extraction of Slow Features: Complex Cells Learn Disparity and Translation Invariance from Natural Images

Alistair Bray, Dominique Martinez

In Slow Feature Analysis (SFA [1]), it has been demonstrated that high-order invariant properties can be extracted by projecting in(cid:173) puts into a

nonlinear space and computing the slowest changing features in this space; this has been proposed as a simple general model for learning nonlinear invariances in the visual system. However, this method is highly constrained by the curse of dimensionality which limits it to simple theoretical simulations. This paper demonstrates that by using a different but closely-related objective function for extracting slowly varying features ([2, 3]), and then exploiting the kernel trick, this curse can be avoided. Using this new method we show that both the complex cell properties of translation invariance and disparity coding can be learnt simultaneously from natural images when complex cells are driven by simple cells also learnt from the image.

\*\*\*\*\*

#### Monaural Speech Separation

Guoning Hu, Deliang Wang

that deals with

\*\*\*\*\*

#### Multiplicative Updates for Nonnegative Quadratic Programming in Support Vector Machines

Fei Sha, Lawrence Saul, Daniel Lee

We derive multiplicative updates for solving the nonnegative quadratic programming problem in support vector machines (SVMs). The updates have a simple closed form, and we prove that they converge monotonically to the solution of the maximum margin hyperplane. The updates optimize the traditionally proposed objective function for SVMs. They do not involve any heuristics such as choosing a learning rate or deciding which variables to update at each iteration. They can be used to adjust all the quadratic programming variables in parallel with a guarantee of improvement at each iteration. We analyze the asymptotic convergence of the updates and show that the coefficients of non-support vectors decay geometrically to zero at a rate that depends on their margins. In practice, the updates converge very rapidly to good classifiers.

\*\*\*\*\*

#### Kernel Design Using Boosting

Koby Crammer, Joseph Keshet, Yoram Singer

The focus of the paper is the problem of learning kernel operators from empirical data. We cast the kernel design problem as the construction of an accurate kernel from simple (and less accurate) base kernels. We use the boosting paradigm to perform the kernel construction process. To do so, we modify the booster so as to accommodate kernel operators. We also devise an efficient weak-learner for simple kernels that is based on generalized eigen vector decomposition. We demonstrate the effectiveness of our approach on synthetic data and on the USPS datasets. On the USPS dataset, the performance of the Perceptron algorithm with learned kernels is systematically better than a fixed RBF kernel.

\*\*\*\*\*

#### Rational Kernels

Corinna Cortes, Patrick Haffner, Mehryar Mohri

We introduce a general family of kernels based on weighted transducers or rational relations, rational kernels, that can be used for analysis of variable-length sequences or more generally weighted automata, in applications such as computational biology or speech recognition. We show that rational kernels can be computed efficiently using a general algorithm of composition of weighted transducers and a general single-source shortest-distance algorithm. We also describe several general families of positive definite symmetric rational kernels. These general kernels can be combined with Support Vector Machines to form efficient and powerful techniques for spoken-dialog classification: highly complex kernels become easy to design and implement and lead to substantial improvements in the classification accuracy. We also show that the string kernels considered in applications to computational biology are all specific instances of rational kernels.

\*\*\*\*\*

#### Linear Combinations of Optic Flow Vectors for Estimating Self-Motion - a Real-World



## rld Test of a Neural Model

Matthias Franz, Javaan Chahl

The tangential neurons in the monkey brain are sensitive to the typical optic flow patterns generated during self-motion. In this study, we examine whether a simplified linear model of these neurons can be used to estimate self-motion from the optic flow. We present a theory for the construction of an estimator consisting of a linear combination of optic flow vectors that incorporates prior knowledge both about the distance distribution of the environment, and about the noise and self-motion statistics of the sensor. The estimator is tested on a gantry car carrying an omnidirectional vision sensor. The experiments show that the proposed approach leads to accurate and robust estimates of rotation rates, whereas translation estimates turn out to be less reliable.

\*\*\*\*\*

## Kernel Dependency Estimation

Jason Weston, Olivier Chapelle, Vladimir Vapnik, André Elisseeff, Bernhard Schölkopf

We consider the learning problem of finding a dependency between a general class of objects and another, possibly different, general class of objects.

The objects can be for example: vectors, images, strings, trees or graphs.

Such a task is made possible by employing similarity measures in both input and output spaces using `ker(cid:173)nel` functions, thus embedding the objects into vector spaces. We experimentally validate our approach on several tasks: mapping strings to strings, pattern recognition, and reconstruction from partial images.

\*\*\*\*\*

## Circuit Model of Short-Term Synaptic Dynamics

Shih-Chii Liu, Malte Boegershausen, Pascal Suter

We describe a model of short-term synaptic depression that is derived from a silicon circuit implementation. The dynamics of this circuit model are similar to the dynamics of some present theoretical models of short-term depression except that the recovery dynamics of the variable describing the depression is nonlinear and it also depends on the presynaptic frequency. The equations describing the steady-state and transient responses of this synaptic model fit the experimental results obtained from a fabricated silicon network consisting of leaky integrate-and-fire neurons and different types of synapses. We also show experimental data demonstrating the possible computational roles of depression. One possible role of a depressing synapse is that the input can quickly bring the neuron up to threshold when the membrane potential is close to the resting potential.

\*\*\*\*\*

## Coulomb Classifiers: Generalizing Support Vector Machines via an Analogy to Electrostatic Systems

Sepp Hochreiter, Michael C. Mozer, Klaus Obermayer

We introduce a family of classifiers based on a physical analogy to an electrostatic system of charged conductors. The family, called Coulomb classifiers, includes the two best-known support-vector machines (SVMs), the  $\{SVM$  and the  $C\{SVM$ . In the electrostatics analogy, a training example corresponds to a charged conductor at a given location in space, the classification function corresponds to the electrostatic potential function, and the training objective function corresponds to the Coulomb energy. The electrostatic framework provides not only a novel interpretation of existing algorithms and their interrelationships, but it suggests a variety of new methods for SVMs including kernels that bridge the gap between polynomial and radial-basis functions, objective functions that do not require positive-definite kernels, regularization techniques that allow for the construction of an optimal classifier in Minkowski space. Based on the framework, we propose novel SVMs and perform simulation studies to show that they are comparable or superior to standard SVMs. The experiments include classification tasks on data which are represented in terms of their pairwise proximities, where a Coulomb Classifier outperformed standard SVMs.

\*\*\*\*\*

### Robust Novelty Detection with Single-Class MPM

Laurent Ghaoui, Michael Jordan, Gert Lanckriet

the "single-class minimax probabil(cid:173)

\*\*\*\*\*

### Binary Tuning is Optimal for Neural Rate Coding with High Temporal Resolution

Matthias Bethge, David Rotermund, Klaus Pawelzik

Here we derive optimal gain functions for minimum mean square re(cid:173) constr  
uction from neural rate responses subjected to Poisson noise. The shape of these  
functions strongly depends on the length  $T$  of the time window within which spik  
es are counted in order to estimate the under(cid:173) lying firing rate. A phas  
e transition towards pure binary encoding occurs if the maximum mean spike count  
becomes smaller than approximately three provided the minimum firing rate is ze  
ro. For a particular function class, we were able to prove the existence of a se  
cond-order phase tran(cid:173) sition analytically. The critical decoding time w  
indow length obtained from the analytical derivation is in precise agreement wit  
h the numerical results. We conclude that under most circumstances relevant to i  
nforma(cid:173) tion processing in the brain, rate coding can be better ascribed  
to a binary (low-entropy) code than to the other extreme of rich analog coding.

\*\*\*\*\*

### Feature Selection in Mixture-Based Clustering

Martin Law, Anil Jain, Mrio Figueiredo

There exist many approaches to clustering, but the important issue of feature se  
lection, i.e., selecting the data attributes that are relevant for clustering, i  
s rarely addressed. Feature selection for clustering is dif(cid:173) cult due to the abse  
nce of class labels. We propose two approaches to feature selection in the conte  
xt of Gaussian mixture-based clustering. In the f(cid:173) rst one, instead of making hard  
selections, we estimate feature saliencies. An expectation-maximization (EM) al  
gorithm is derived for this task. The second approach extends Koller and Sahami'  
s mutual-information- based feature relevance criterion to the unsupervised case  
. Feature selec- tion is then carried out by a backward search scheme. This sche  
me can be classi(cid:173) ed as a "wrapper", since it wraps mixture estimation in an oute  
r layer that performs feature selection. Experimental results on synthetic and r  
eal data show that both methods have promising performance.

\*\*\*\*\*

### Minimax Differential Dynamic Programming: An Application to Robust Biped Walking

Jun Morimoto, Christopher Atkeson

We developed a robust control policy design method in high-dimensional state spa  
ce by using differential dynamic programming with a minimax criterion. As an exa  
mple, we applied our method to a simulated f(cid:173) ve link biped robot. The results sho  
w lower joint torques from the optimal con- trol policy compared to a hand-tuned  
PD servo controller. Results also show that the simulated biped robot can succe  
ssfully walk with unknown disturbances that cause controllers generated by stand  
ard differential dy- namic programming and the hand-tuned PD servo to fail. Lear  
ning to compensate for modeling error and previously unknown disturbances in con  
junction with robust control design is also demonstrated.

\*\*\*\*\*

### Theory-Based Causal Inference

Joshua Tenenbaum, Thomas Griffiths

People routinely make sophisticated causal inferences unconsciously, ef- fortles  
sly, and from very little data - often from just one or a few ob- servations. We  
argue that these inferences can be explained as Bayesian computations over a hy  
pothesis space of causal graphical models, shaped by strong top-down prior knowl  
edge in the form of intuitive theories. We present two case studies of our appro  
ach, including quantitative mod- els of human causal judgments and brief compari  
sons with traditional bottom-up models of inference.

\*\*\*\*\*

### Learning to Classify Galaxy Shapes Using the EM Algorithm

Sergey Kirshner, Igor Cadez, Padhraic Smyth, Chandrika Kamath

We describe the application of probabilistic model-based learning to the problem  
of automatically identifying classes of galaxies, based on both morphological a

nd pixel intensity characteristics. The EM algorithm can be used to learn how to spatially orient a set of galaxies so that they are geometrically aligned. We augment this "ordering-model" with a mixture model on objects, and demonstrate how classes of galaxies can be learned in an unsupervised manner using a two-level EM algorithm. The resulting models provide highly accurate classification of galaxies in cross-validation experiments.

\*\*\*\*\*

#### FloatBoost Learning for Classification

Stan Li, Zhenqiu Zhang, Heung-yeung Shum, Hongjiang Zhang

AdaBoost [3] minimizes an upper error bound which is an exponential function of the margin on the training set [14]. However, the ultimate goal in applications of pattern classification is always minimum error rate. On the other hand, AdaBoost needs an effective procedure for learning weak classifiers, which by itself is difficult especially for high dimensional data. In this paper, we present a novel procedure, called FloatBoost, for learning a better boosted classifier. FloatBoost uses a backtrack mechanism after each iteration of AdaBoost to remove weak classifiers which cause higher error rates. The resulting float-boosted classifier consists of fewer weak classifiers yet achieves lower error rates than AdaBoost in both training and test. We also propose a statistical model for learning weak classifiers, based on a stagewise approximation of the posterior using an overcomplete set of scalar features. Experimental comparisons of FloatBoost and AdaBoost are provided through a difficult classification problem, face detection, where the goal is to learn from training examples a highly nonlinear classifier to differentiate between face and nonface patterns in a high dimensional space. The results clearly demonstrate the promises made by FloatBoost over AdaBoost.

\*\*\*\*\*

#### Adaptive Scaling for Feature Selection in SVMs

Yves Grandvalet, Stéphane Canu

This paper introduces an algorithm for the automatic relevance determination of input variables in kernelized Support Vector Machines. Relevance is measured by scale factors defining the input space metric, and feature selection is performed by assigning zero weights to irrelevant variables. The metric is automatically tuned by the minimization of the standard SVM empirical risk, where scale factors are added to the usual set of parameters defining the classifier. Feature selection is achieved by constraints encouraging the sparsity of scale factors. The resulting algorithm compares favorably to state-of-the-art feature selection procedures and demonstrates its effectiveness on a demanding facial expression recognition problem.

\*\*\*\*\*

#### Adaptive Classification by Variational Kalman Filtering

Peter Sykacek, Stephen J. Roberts

We propose in this paper a probabilistic approach for adaptive inference of generalized nonlinear classification that combines the computational advantage of a parametric solution with the flexibility of sequential sampling techniques. We regard the parameters of the classifier as latent states in a first order Markov process and propose an algorithm which can be regarded as variational generalization of standard Kalman filtering. The variational Kalman filter is based on two novel lower bounds that enable us to use a non-degenerate distribution over the adaptation rate. An extensive empirical evaluation demonstrates that the proposed method is capable of inferring competitive classifiers both in stationary and non-stationary environments. Although we focus on classification, the algorithm is easily extended to other generalized nonlinear models.

\*\*\*\*\*

#### A Statistical Mechanics Approach to Approximate Analytical Bootstrap Averages

Dörthe Malzahn, Manfred Oppen

We apply the replica method of Statistical Physics combined with a variational method to the approximate analytical computation of bootstrap averages for estimating the generalization error. We demonstrate our approach on regression with Gaussian processes and compare our results with averages obtained by Monte-Carlo sampling.

\*\*\*\*\*

## Developing Topography and Ocular Dominance Using Two aVLSI Vision Sensors and a Neurotrophic Model of Plasticity

Terry Elliott, Jörg Kramer

A neurotrophic model for the co-development of topography and ocular dominance columns in the primary visual cortex has recently been proposed. In the present work, we test this model by driving it with the output of a pair of neuronal vision sensors stimulated by disparate moving patterns. We show that the temporal correlations in the spike trains generated by the two sensors elicit the development of refined topography and ocular dominance columns, even in the presence of significant amounts of spontaneous activity and fixed-pattern noise in the sensors.

\*\*\*\*\*

## Gaussian Process Priors with Uncertain Inputs Application to Multiple-Step Ahead Time Series Forecasting

Agathe Girard, Carl Rasmussen, Joaquin Quiñero Candela, Roderick Murray-Smith

We consider the problem of multi-step ahead prediction in time series analysis using the non-parametric Gaussian process model.  $n$ -step ahead forecasting of a discrete-time non-linear dynamic system can be performed by doing repeated one-step ahead predictions. For a state-space at time model of the form is based on the point estimates of the previous outputs. In this paper, we show how, using an analytical Gaussian approximation, we can formally incorporate the uncertainty about intermediate regressor values, thus updating the uncertainty on the current prediction.

\*\*\*\*\*

## Maximum Likelihood and the Information Bottleneck

Noam Slonim, Yair Weiss

\*\*\*\*\*

## The Stability of Kernel Principal Components Analysis and its Relation to the Process Eigenspectrum

Christopher Williams, John Shawe-taylor

In this paper we analyze the relationships between the eigenvalues of the  $m \times m$  Gram matrix  $K$  for a kernel  $k(\cdot, \cdot)$  corresponding to a sample  $X_1, \dots, X_m$  drawn from a density  $p(x)$  and the eigenvalues of the corresponding continuous eigenproblem. We bound the differences between the two spectra and provide a performance bound on kernel PCA.

\*\*\*\*\*

## Adaptive Quantization and Density Estimation in Silicon

David Hsu, Seth Bridges, Miguel Figueroa, Chris Diorio

We present the bump mixture model, a statistical model for analog data where the probabilistic semantics, inference, and learning rules derive from low-level transistor behavior. The bump mixture model relies on translinear circuits to perform probabilistic inference, and floating-gate devices to perform adaptation. This system is low power, asynchronous, and fully parallel, and supports various on-chip learning algorithms. In addition, the mixture model can perform several tasks such as probability estimation, vector quantization, classification, and clustering. We tested a fabricated system on clustering, quantization, and classification of handwritten digits and show performance comparable to the EM algorithm on mixtures of Gaussians.

\*\*\*\*\*

## Improving Transfer Rates in Brain Computer Interfacing: A Case Study

Peter Meinicke, Matthias Kaper, Florian Hoppe, Manfred Heumann, Helge Ritter

In this paper we present results of a study on brain computer interfacing. We adopted an approach of Farwell & Donchin [4], which we tried to improve in several aspects. The main objective was to improve the transfer rates based on offline analysis of EEG-data but within a more realistic setup closer to an online realization than in the original studies. The objective was achieved along two different tracks: on the one hand we used state-of-the-art machine learning techniques for signal classification and on the other hand we augmented the data space b

y using more electrodes for the interface. For the classification task we utilized SVMs and, as motivated by recent findings on the learning of discriminative densities, we accumulated the values of the classification function in order to combine several classifications, which finally lead to significantly improved rates as compared with techniques applied in the original work. In combination with the data space augmentation, we achieved competitive transfer rates at an average of 50.5 bits/min and with a maximum of 84.7 bits/min.

\*\*\*\*\*

#### Discriminative Learning for Label Sequences via Boosting

Yasemin Altun, Thomas Hofmann, Mark Johnson

This paper investigates a boosting approach to discriminative learning of label sequences based on a sequence rank loss function. The proposed method combines many of the advantages of boosting schemes with the efficiency of dynamic programming methods and is attractive both, conceptually and computationally. In addition, we also discuss alternative approaches based on the Hamming loss for label sequences. The sequence boosting algorithm offers an interesting alternative to methods based on HMMs and the more recently proposed Conditional Random Fields. Applications areas for the presented technique range from natural language processing and information extraction to computational biology. We include experiments on named entity recognition and part-of-speech tagging which demonstrate the validity and competitiveness of our approach.

\*\*\*\*\*

#### Reinforcement Learning to Play an Optimal Nash Equilibrium in Team Markov Games

Xiaofeng Wang, Tuomas Sandholm

Multiagent learning is a key problem in AI. In the presence of multiple Nash equilibria, even agents with non-conflicting interests may not be able to learn an optimal coordination policy. The problem is exacerbated if the agents do not know the game and independently receive noisy payoffs. So, multiagent reinforcement learning involves two inter-related problems: identifying the game and learning to play. In this paper, we present optimal adaptive learning, the first algorithm that converges to an optimal Nash equilibrium with probability 1 in any team Markov game. We provide a convergence proof, and show that the algorithm's parameters are easy to set to meet the convergence conditions.

\*\*\*\*\*

#### Using Manifold Structure for Partially Labeled Classification

Mikhail Belkin, Partha Niyogi

We consider the general problem of utilizing both labeled and unlabeled data to improve classification accuracy. Under the assumption that the data lie on a submanifold in a high dimensional space, we develop an algorithmic framework to classify a partially labeled data set in a principled manner. The central idea of our approach is that classification functions are naturally defined only on the submanifold in question rather than the total ambient space. Using the Laplace Beltrami operator one produces a basis for a Hilbert space of square integrable functions on the submanifold. To recover such a basis, only unlabeled examples are required. Once a basis is obtained, training can be performed using the labeled data set. Our algorithm models the manifold using the adjacency graph for the data and approximates the Laplace Beltrami operator by the graph Laplacian. Practical applications to image and text classification are considered.

\*\*\*\*\*

#### Recovering Intrinsic Images from a Single Image

Marshall Tappen, William Freeman, Edward Adelson

We present an algorithm that uses multiple cues to recover shading and reflectance intrinsic images from a single image. Using both color information and a classifier trained to recognize gray-scale patterns, each image derivative is classified as being caused by shading or a change in the surface's reflectance. Generalized Belief Propagation is then used to propagate information from areas where the correct classification is clear to areas where it is ambiguous. We also show results on real images.

\*\*\*\*\*

# A Note on the Representational Incompatibility of Function Approximation and Factored Dynamics

Eric Allender, Sanjeev Arora, Michael Kearns, Cristopher Moore, Alexander Russell

We establish a new hardness result that shows that the difficulty of planning in factored Markov decision processes is representational rather than just computational. More precisely, we give a fixed family of factored MDPs with linear rewards whose optimal policies and value functions simply cannot be represented succinctly in any standard parametric form. Previous hardness results indicated that computing good policies from the MDP parameters was difficult, but left open the possibility of succinct function approximation for any fixed factored MDP. Our result applies even to policies which yield a polynomially poor approximation to the optimal value, and highlights interesting connections with the complexity class of Arthur-Merlin games.

\*\*\*\*\*

# A Prototype for Automatic Recognition of Spontaneous Facial Actions

M.S. Bartlett, G.C. Littlewort, T.J. Sejnowski, J.R. Movellan

We present ongoing work on a project for automatic recognition of spontaneous facial actions. Spontaneous facial expressions differ substantially from posed expressions, similar to how continuous, spontaneous speech differs from isolated words produced on command. Previous methods for automatic facial expression recognition assumed images were collected in controlled environments in which the subjects deliberately faced the camera. Since people often nod or turn their heads, automatic recognition of spontaneous facial behavior requires methods for handling out-of-image-plane head rotations. Here we explore an approach based on 3-D warping of images into canonical views. We evaluated the performance of the approach as a front-end for a spontaneous expression recognition system using support vector machines and hidden Markov models. This system employed general purpose learning mechanisms that can be applied to recognition of any facial movement. The system was tested for recognition of a set of facial actions defined by the Facial Action Coding System (FACS). We showed that 3D tracking and warping followed by machine learning techniques directly applied to the warped images, is a viable and promising technology for automatic facial expression recognition. One exciting aspect of the approach presented here is that information about movement dynamics emerged out of filters which were derived from the statistics of images.

\*\*\*\*\*

# Exact MAP Estimates by (Hyper)tree Agreement

Martin J. Wainwright, Tommi Jaakkola, Alan Willsky

We describe a method for computing provably exact maximum a posteriori (MAP) estimates for a subclass of problems on graphs with cycles. The basic idea is to represent the original problem on the graph with cycles as a convex combination of tree-structured problems. A convexity argument then guarantees that the optimal value of the original problem (i.e., the log probability of the MAP assignment) is upper bounded by the combined optimal values of the tree problems. We prove that this upper bound is met with equality if and only if the tree problems share an optimal configuration in common. An important implication is that any such shared configuration must also be the MAP configuration for the original problem. Next we develop a tree-reweighted max-product algorithm for attempting to find convex combinations of tree-structured problems that share a common optimum. We give necessary and sufficient conditions for a fixed point to yield the exact MAP estimate. An attractive feature of our analysis is that it generalizes naturally to convex combinations of hypertree-structured distributions.

\*\*\*\*\*

# Prediction of Protein Topologies Using Generalized IOHMMs and RNNs

Gianluca Pollastri, Pierre Baldi, Alessandro Vullo, Paolo Frasconi

We develop and test new machine learning methods for the prediction of topological representations of protein structures in the form of coarse- or (cid:12)ne-grained contact or distance maps that are translation and rotation invariant.

The methods are based on generalized input-output hidden Markov models (GIOHMMs) and generalized recursive neural networks (GRNNs). The methods are used to predict topology directly in the (cid:12)ne-grained case and, in the coarse-grained case, indirectly by (cid:12)rst learning how to score candidate graphs and then using the scoring function to search the space of possible con(cid:12)guratio ns. Computer simulations show that the pre dictors achieve state-of-the-art performance.

\*\*\*\*\*