



CONTRIBUTED ARTICLE

Neural Network Exploration Using Optimal Experiment Design

DAVID A. COHN

Massachusetts Institute of Technology

(Received 21 November 1994; revised and accepted 14 November 1995)

Abstract—*I consider the question “How should one act when the only goal is to learn as much as possible?”. Building on the theoretical results of Fedorov (1972, *Theory of Optimal Experiments*, Academic Press) and MacKay (1992, *Neural Computation*, 4, 590–604), I apply techniques from optimal experiment design (OED) to guide the query/action selection of a neural network learner. I demonstrate that these techniques allow the learner to minimize its generalization error by exploring its domain efficiently and completely. I conclude that, while not a panacea, OED-based query/action selection has much to offer, especially in domains where its high computational costs can be tolerated.*

Copyright © 1996 Elsevier Science Ltd

Keywords—Active learning, Exploration, Optimal experiment design, Queries, Uncertainty.

1. INTRODUCTION

In many natural learning problems, the learner has the ability to act on its environment and gather data that will resolve its uncertainties. Most machine learning research, however, treats the learner as a passive recipient of data and ignores the role of this “active” component of learning. In this paper I employ techniques from the field of optimal experiment design (OED) to guide the actions of a learner, selecting actions/queries that are statistically expected to minimize its uncertainty and error.

Acknowledgements: Support is provided in part by a grant from the National Science Foundation under contract ASC-9217041. The author was also funded by ATR Human Information Processing Laboratories, Siemens Corporate Research and NSF grant CDA-9309300.

I am indebted to Michael I. Jordan and David J. C. MacKay for their help in making this research possible. Thanks are also due to the University of Toronto and the Xerion group for use of and assistance with the Xerion simulator, and to Cesare Alippi and two anonymous reviewers for useful comments on an earlier draft of this paper.

Requests for reprints should be sent to D. A. Cohn, Harlequin, Inc., One Cambridge Center, Cambridge, MA 02142, USA. ((617) 374-2442), E-mail: cohn@harlequin.com

¹ In some cases active selection of training data can sharply reduce worst case computational complexity from NP-complete to polynomial time (Baum & Lang, 1991), and in special cases to linear time.

1.1. Active Learning

Exploiting the active component of learning typically leads to improved generalization, usually at the cost of additional computation (see Figure 1) (Angluin, 1982; Cohn et al., 1990; Hwang et al., 1990).¹ There are two common situations where this tradeoff is desirable: in many situations the cost of taking an action outweighs the cost of the computation required to incorporate new information into the model. In these cases we wish to select queries judiciously, so that we can build a good model with the fewest data. This is the case if, for example, we are drilling oil wells or taking seismic measurements to locate buried waste. In other situations the data, although cheap, must be chosen carefully to ensure thorough exploration. Large amounts of data may be useless if they all come from an uninteresting part of the domain. This is the case with the learning control of a robot arm: exploring by generating random motor torques cannot be expected to give good coverage of the domain.

As computation becomes cheaper and faster, more problems fall within the realm where it is both desirable and practical to pursue active learning, expending more computation to ensure that one's exploration provides good data. The field of optimal experiment design, which is concerned with the statistics of gathering new data, provides a principled way to guide this exploration. This paper builds on the theoretical results of Fedorov (1972)

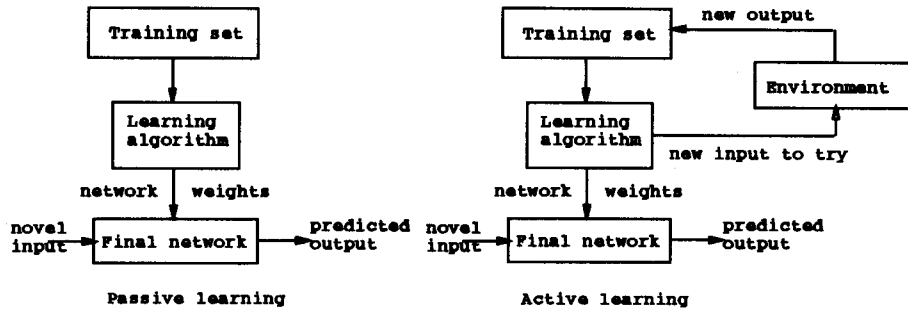


FIGURE 1. An active system will typically evaluate/train on its data iteratively, determining its next input based on the previous training examples. This iterative training may be computationally expensive, especially for learning systems like neural networks where good incremental algorithms are not available.

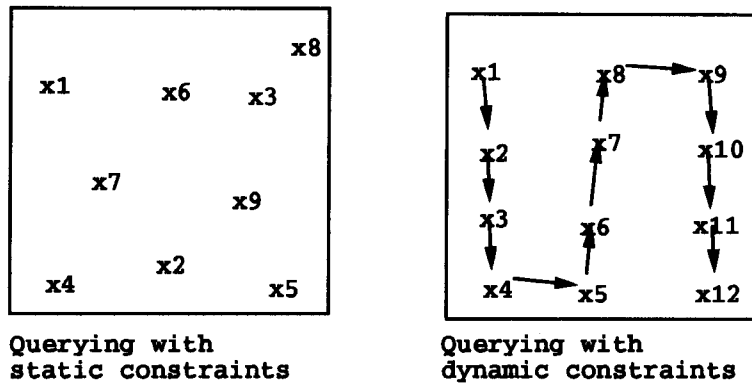


FIGURE 2. In problems with dynamic constraints, the set of candidate \tilde{x} can change after each query. The \tilde{x}_{t+1} accessible to the learner on the bottom depends on the choice made for \tilde{x}_t .

and MacKay (1992) to empirically demonstrate how OED may be applied to neural network learning, and to determine under what circumstances it is an effective approach.

The remainder of this section provides a formal problem definition, followed by a brief review of related work using optimal experiment design. Section 2 differentiates several classes of active learning problems for which OED is appropriate. Section 3 describes the theory behind optimal experiment design, and Section 4 demonstrates its application to several simple problems. Section 5 considers the computational costs of these experiments, and Section 6 concludes with a discussion of the results and implications for future work.

1.2. Problem Definition

I consider the problem of learning an input-output mapping $X \rightarrow Y$ from a set of m training examples $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$, where $x_i \in X$, $y_i \in Y$.

I denote the parameterized learner $f_w(\cdot)$, where $y = f_w(x)$ is the learner's output given input x and parameter vector w . The learner is trained by adjusting w to minimize the residual $S^2 = 1/(2m) \sum_{i=1}^m (f_w(x_i) - y_i)^T (f_w(x_i) - y_i)$ on the training set. Let \hat{w} be the parameter vector that minimizes S^2 .

Then $\hat{y} = f_{\hat{w}}(x)$ is the learner's "best guess" of the mapping $X \rightarrow Y$: given x , \hat{y} is an estimate of the corresponding y . For the purposes of this paper, $f_w(\cdot)$ will be taken to be a feedforward neural network with weight vector w .

At each time step, the learner is allowed to select a new training input \tilde{x} from a set of candidate inputs \tilde{X} . The selection of \tilde{x} may be viewed as a "query" (as to an oracle), as an "experiment", or simply as an "action". Having selected \tilde{x} , the learner observes the corresponding \tilde{y} , and incorporates the resulting new example (\tilde{x}, \tilde{y}) into the training set. The learner is updated, selects another new \tilde{x} and the process is repeated.

The goal is to choose examples that minimize the expectation of the learner's mean squared error over an input distribution $P(x)$. In contrast to some other learning paradigms (Valiant, 1984; Blumer et al., 1989), we will assume that $P(x)$ is known or that we have the means to cheaply draw unlabeled examples according to $P(x)$ so that we can approximate it.² In

² This assumption is reasonable in many situations; if we are attempting to learn to control a robot arm, for example, it is appropriate to assume that we know over what range we wish to control it. Similarly, if we are learning to classify speech, we can cheaply acquire unlabeled speech samples.

the context where $P(x)$ is known, it may be useful to think of it as a *weighting function* over X which we use when measuring performance.

Example 1: Mapping Buried Waste. Consider a mobile sensor array traversing a landscape to map out subsurface electromagnetic anomalies. Its location at time t serves as input x_t , and the instrument reading at that location is output y_t . At the next time step, it can choose its new input \tilde{x} from any location contiguous to its present position.

Example 2: Robot Arm Dynamics. Consider learning the dynamics of a robot arm. The input is the state-action triplet $x_t = \{\Theta_t, \dot{\Theta}_t, \tau_t\}$, where Θ_t and $\dot{\Theta}_t$ are the arm's joint angles and velocities, respectively, and τ_t is the torque applied at time t . The output $y_t = \{\Theta_{t+1}, \dot{\Theta}_{t+1}\}$ is the resulting state. Note that here, although we may specify an arbitrary torque τ_t , the rest of the input, $\{\Theta_t, \dot{\Theta}_t\}$ is determined by y_{t-1} .

While the above problem definition has wide-ranging application, it is by no means all-encompassing. For some learning problems, we are not interested in the entire mapping $X \rightarrow Y$, but in finding the x that maximizes y . In this case, we may rely on the broad literature of optimization and response surface techniques (Box & Draper, 1987). In other learning problems there may be additional constraints that must be considered, such as the need to avoid "failure" states. If the learner is required to perform as it learns (e.g., in a control task), we may also need to balance exploration and exploitation. In this paper, however, I assume that the cost of the query \tilde{x} is independent of \tilde{x} , and that the sole aim of active learning is to minimize E_{MSE} , the mean squared difference of the learner's output and the system's output, averaged over $P(x)$.

In addition to the "supervised" learning problems described in this paper there is the broad class of reinforcement learning problems (Barto et al., 1995). In reinforcement learning, the result of an action is not the target output value, but simply some reward or punishment which may or may not depend on the learner's entire previous sequence of actions. There has been some research on optimal exploration in reinforcement learning problems (e.g., Thrun (1992)), but the problem is still largely unsolved.

1.3. Related Work with Optimal Experiment Design

The literature on optimal experiment design is immense and dates back at least 50 years. I will just mention here a few closely related theoretical results and empirical studies; the interested reader should consult Atkinson and Donev (1992) for a survey of

results and applications using optimal experiment design.

A canonical description of the theory of OED is given in Fedorov (1972). MacKay (1992) showed that OED could be incorporated into a Bayesian framework for neural network data selection, and described several interesting optimization criteria. Sollich (1994) considers the theoretical generalization performance of linear networks given greedy vs globally optimal queries and varying assumptions on teacher distributions.

Empirically, optimal experiment design techniques have been successful when used for system identification tasks. In these cases a good parameterized model of the system is available, and learning involves finding the proper parameters. Ford et al. (1989) present a survey of such methods, with examples for a number on nonlinear parameterized problems. Armstrong (1989) used a form of OED to identify link masses and inertial moments of a robot arm, and found that automatically generated training trajectories provided a significant improvement over human-designed trajectories. Subrahmonia et al. (1992) successfully used experiment design to guide exploration of a sensor moving along the surface of an object parameterized as an unknown quadric.

Empirical work on using OED with neural networks is sparse. Plutowski and White (1993) successfully used it to filter an already-labeled data set for maximally informative points. Choueiki (1995) has successfully trained neural networks on quadratic surfaces with data drawn according to the D -optimality criterion, which is discussed in the Appendix.

2. LEARNING WITH STATIC AND DYNAMIC CONSTRAINTS

Different problems impose different constraints on the specification of \tilde{x} . These constraints may be classified as being either static or dynamic, and problems with dynamic constraints may be further divided according to whether or not the dynamics of the constraints are known *a priori*.

When a learner has *static* input constraints, its range of choices for \tilde{x} is fixed, regardless of previous actions. In some situations with static constraints, active learning may drive a learner's generalization error to zero as (2^{-m}) , in contrast to $O(1/m)$ for a learner using random examples³ (Baum & Haussler, 1989; Blumer et al., 1989; Cohn & Tesauro, 1992; Haussler, 1992), although the worst case bounds are no better than those for choosing at random (Eisenburg & Rivest, 1990).

Average case analysis indicates that on many

³ Consider cases where binary search is applicable.

domains with static constraints, the expected performance of active selection of training examples is significantly better than that of random sampling (Freund & Seung, 1993); these results have been supported by empirical studies (Cohn et al., 1990; Baum & Lang, 1991; Hwang et al., 1990). The above algorithms, however, are only suitable for classification problems, and in some cases require restricted architectures. The OED-based approach discussed in this paper is applicable to any network architecture whose output is differentiable with respect to its parameters, and may be used on both regression and classification problems.

When learning with *dynamic* input constraints, the set of available \tilde{x} may change on successive steps. Typically, these constraints represent some state of the system that is altered by the learner's actions. Training examples then describe a trajectory through state space (see Figure 2).

In problems where the dynamics are known, we may predict *a priori* what constraints we will face at time t , given an initial state and actions x_1, x_2, \dots, x_t . A more common, and more difficult problem is learning when the dynamics of the constraints are not known, and must be accommodated online. Learning the dynamics of a robot arm $\{\Theta_t, \dot{\Theta}_t, \tau_t\} \rightarrow \{\Theta_{t+1}, \dot{\Theta}_{t+1}\}$ is an example of this type of problem. At each time step t , the model input \tilde{x} is a state-action pair $\{\Theta_t, \dot{\Theta}_t, \tau_t\}$, where Θ_t and $\dot{\Theta}_t$ are constrained to be the learner's current state. Until the action is selected and taken, the learner does not know what its new state, and thus its new constraints will be (this is in fact exactly what it is attempting to learn).

In cases where the dynamics of the constraints are known *a priori*, we can plan a trajectory that will uniformly cover X in some prespecified number of steps. In general, though, we will have to resort to some online process to decide "what to try next". There are many successful heuristic exploration strategies (e.g., Thrun & Möller, 1992; Linden & Weber, 1993; Schmidhuber & Storck, 1993; Schaal & Atkeson, 1994), but for the most part, these methods are restricted to discrete state spaces. Continuous state and action spaces must be accommodated either through arbitrary discretization or through some form of on-line partitioning strategy, such as Moore's Parti-Game algorithm (Moore, 1994). The OED-based approach discussed in this paper is, by nature, applicable to domains with both continuous state and action spaces.

3. DATA SELECTION ACCORDING TO OED

One of the standard assumptions of OED is that the underlying system being learned has a structure that is compatible with the structure of the learner. In our case, this corresponds to assuming that there exist some set of parameters for our neural network that

will make the network's behavior closely model the system's behavior. Our stated goal is to find a set of parameters that minimize E_{MSE} .⁴ We will try to find these parameters by selecting training examples that optimally distinguish good parameter values from bad ones; specifically, we want examples such that any set of parameters which would result in a large E_{MSE} will have a large training error, and will be rejected.

We can measure the degree of constraint around a locally optimal solution as the amount by which the training error increases as network outputs are perturbed from that solution. In the terminology of Fedorov (1972), the inverse of this measure is the learner's "dispersion"—a large dispersion indicates that the estimate is relatively unconstrained. For problems where E_{MSE} has a statistical interpretation of a log likelihood (which is the most common rationale for using it as an error measure), the dispersion has the additional interpretation of the variance of the learner's estimate. This interpretation relies on a number of statistical assumptions, which are discussed in greater detail in Section 3.4.2.

The computations for deriving dispersions and variances are equivalent, so I will use the statistical terminology and notation in both cases, denoting as σ_y^2 the "variance" (dispersion) of the learner's output \hat{y} . Given the interpretation of variance, σ_y^2 has a direct quantitative link to E_{MSE} .

We define $P(x, y)$ to be the unknown joint distribution over x and y , in which case $P(x)$ is the known marginal distribution of x . The learner's output on input x , given training set \mathcal{D} is $f_{\tilde{w}}(x; \mathcal{D})$.⁵ We can then write the expected error of the learner as follows:

$$E_{\text{MSE}} = \int_x E_T[(f_{\tilde{w}}(x; \mathcal{D}) - y|x)^2|x]P(x)dx, \quad (1)$$

where $E_T[\cdot]$ denotes expectation over $P(y|x)$ and over training sets \mathcal{D} . The expectation inside the integral may be decomposed as follows (Geman et al., 1992):

$$\begin{aligned} E_T[(f_{\tilde{w}}(x; \mathcal{D}) - y|x)^2|x] &= E[(y|x - E[y|x])^2] \\ &\quad + (E[f_{\tilde{w}}(x; \mathcal{D})] - E[y|x])^2 \\ &\quad + E_{\mathcal{D}}[(f_{\tilde{w}}(x; \mathcal{D}) - E_{\mathcal{D}}[f_{\tilde{w}}(x; \mathcal{D})])^2], \end{aligned} \quad (2)$$

where $E_{\mathcal{D}}[\cdot]$ denotes the expectation over training

⁴ An alternative goal of system identification is discussed briefly in the Appendix, and other interesting goals, such as eigenvalue maximization and entropy minimization, may be found in Fedorov (1972) and MacKay (1992).

⁵ Here, I include the training set to make explicit its role in the process. I also present the equations in the univariate setting. All results in this section apply equally to the multivariate case.

sets \mathcal{D} and the remaining expectations on the right-hand side are expectations with respect to the conditional density $P(y|x)$. It is important to remember here that in the case of active learning, the distribution of \mathcal{D} may differ substantially from the joint distribution $P(x, y)$.

The first term in eqn (2) is the variance of y given x —it is the *noise* in the distribution, and does not depend on the learner or on the training data. The second term is the learner's *squared bias*, and the third is its *variance*; these last two terms comprise the mean squared error of the learner with respect to the regression function $E[y|x]$. When the second term of eqn (2) is zero, we say that the learner is *unbiased*. We shall assume that the neural networks considered in this paper are approximately unbiased; that is, their squared bias is negligible when compared with their overall mean squared error.⁶ Thus it is appropriate to focus on algorithms that minimize the learner's error by minimizing its variance:

$$\sigma_{\tilde{y}}^2 = \int E_{\mathcal{D}} [(f_{\tilde{w}}(x; \mathcal{D}) - E_{\mathcal{D}}[f_{\tilde{w}}(x; \mathcal{D})])^2] P(x) dx. \quad (3)$$

3.1. Estimating Variance

Estimates for $\sigma_{\tilde{y}}^2$ may be obtained by adapting techniques derived for linear systems. We write the network's output sensitivity as $g(x) = \partial \tilde{y}|x / \partial w = \partial f_{\tilde{w}}(x) / \partial w$, and define the Fisher information matrix to be

$$\begin{aligned} A &= \frac{1}{S^2} \frac{\partial^2 S^2}{\partial w^2} \\ &= \frac{1}{S^2} \sum_{i=1}^m \left[\frac{\partial \tilde{y}|x_i}{\partial w} \frac{\partial \tilde{y}|x_i^T}{\partial w} + (\tilde{y}|x_i - y_i) \frac{\partial^2 \tilde{y}|x_i}{\partial w^2} \right] \\ &\approx \frac{1}{S^2} \sum_{i=1}^m g(x_i) g(x_i)^T, \end{aligned} \quad (4)$$

where the training residual S^2 serves as a noise estimate. The approximation in eqn (4) holds when the network fits the data well or the error surface has relatively constant curvature in the vicinity of \tilde{w} . It is worth noting that this approximation does not rely on the *actual* target values y_i ; this is the crucial fact which makes the OED approach tractable in this setting. Given A , we may write the parameter covariance matrix as $\sigma_{\tilde{w}}^2 = A^{-1}$ and the estimated output variance at reference input x_r as

$$\sigma_{\tilde{y}|x_r}^2 \approx g(x_r)^T A^{-1} g(x_r) \quad (5)$$

⁶ The bias term will in fact reappear as a limiting factor in the experimental results described in Section 4.2.

subject to the same approximations (see Thisted (1988) for derivations).⁷ Note that the estimate $\sigma_{\tilde{y}|x_r}^2$ applies only to the variance at a particular reference point. Our interest is in estimating $\sigma_{\tilde{y}}^2$, the average variance over all of X . Although we have assumed that we know the distribution $P(x)$, we do not have a method for directly integrating over it. Instead, we opt for a stochastic estimate based on an average of $\sigma_{\tilde{y}|x_r}^2$, with x_r drawn according to $P(x)$. Writing the second moment of g as $gg^T = \langle g(x)g(x)^T \rangle_X$, this estimate can be computed efficiently as

$$\sigma_{\tilde{y}}^2 = \langle \text{tr}(g^T A^{-1} g) \rangle_X = \langle \text{tr}(A^{-1} gg^T) \rangle_X = \text{Tr}(A^{-1} \overline{gg^T}), \quad (6)$$

where $\text{tr}(\cdot)$ is the matrix trace. Instead of recomputing eqn (5) for each reference point, $\overline{gg^T}$ may be computed over the reference points, and eqn (6) evaluated once.

3.2. Predicting Change in Variance

When an input \tilde{x} is queried, we obtain the resulting output \tilde{y} . When the new example (\tilde{x}, \tilde{y}) is added to the training set, the variance of the model will change. We wish to select \tilde{x} optimally, such that the resulting variance, denoted $\tilde{\sigma}_{\tilde{y}}^2$, is minimized.

The network provides a (possibly inaccurate) estimate of the distribution $P(\tilde{y}|\tilde{x})$, embodied in an estimate of the mean $(\tilde{y}|\tilde{x})$ and variance (S^2). Given infinite computational power then, we could use these estimates to stochastically approximate $\tilde{\sigma}_{\tilde{y}|x}^2$ by drawing examples according to our estimate of $P(\tilde{y}|\tilde{x})$, training on them, and averaging the new variances.⁸ In practice though, we must settle for a coarser approximation. Recall that the approximation in eqn (4) is only indirectly dependent on the actual y_i values of the training set; the dependence is implicit in the choice of \tilde{w} that minimizes S^2 . If $P(\tilde{y}|\tilde{x})$ conforms to our expectations, then \tilde{w} and $g(\cdot)$ will remain essentially unchanged, allowing us to compute the new information matrix \tilde{A} as

$$\tilde{A} \approx A + \frac{1}{S^2} g(\tilde{x}) g(\tilde{x})^T. \quad (7)$$

From the new information matrix, we may compute the new parameter variances, and from there, the new

⁷ If the inverse A^{-1} does not exist, then the parameter covariances are not well defined. In practice, one could use the pseudo-inverse, but the need to do this arose very rarely in our experiments, even with small training sets.

⁸ Recent work by Paass and Kindermann (1995) pursues an approximation of this, using Monte Carlo Markov chain techniques.

output variances. By the matrix inversion lemma

$$\tilde{A}^{-1} = \left(A + \frac{1}{S^2} g(\tilde{x}) g(\tilde{x})^T \right)^{-1} = A^{-1} - \frac{A^{-1} g(\tilde{x}) g(\tilde{x})^T A^{-1}}{S^2 + g(\tilde{x})^T A^{-1} g(\tilde{x})}. \quad (8)$$

The utility of querying at \tilde{x} may be expressed in terms of the expected change in the estimated output variance $\sigma_{\tilde{y}}^2$. The expected new output variance at reference point x_r is

$$\begin{aligned} \langle \sigma_{\tilde{y}|x_r}^2 \rangle_{\tilde{y}} &= g(x_r)^T \tilde{A}^{-1} g(x_r) \\ &= g(x_r)^T \left[A^{-1} - \frac{A^{-1} g(\tilde{x}) g(\tilde{x})^T A^{-1}}{S^2 + g(\tilde{x})^T A^{-1} g(\tilde{x})} \right] g(x_r) \\ &= g(x_r)^T A^{-1} g(x_r) - \frac{[g(x_r)^T A^{-1} g(\tilde{x})]^2}{S^2 + g(\tilde{x})^T A^{-1} g(\tilde{x})} \\ &= \sigma_{\tilde{y}|x_r}^2 - \frac{\sigma_{\tilde{y}|x_r, \tilde{x}}^2}{S^2 + \sigma_{\tilde{y}|\tilde{x}}^2}, \end{aligned}$$

where $\sigma_{\tilde{y}|x_r, \tilde{x}}$ is defined as $g(x_r)^T A^{-1} g(\tilde{x})$. Thus, when \tilde{x} is queried, the expected change in output variance at x_r is

$$\langle \Delta \sigma_{\tilde{y}|x_r}^2 \rangle_{\tilde{y}} | \tilde{x} = \frac{\sigma_{\tilde{y}|x_r, \tilde{x}}^2}{S^2 + \sigma_{\tilde{y}|\tilde{x}}^2}. \quad (9)$$

We compute $\langle \Delta \sigma_{\tilde{y}}^2 \rangle | \tilde{x}$ as a stochastic approximation from $\langle \Delta \sigma_{\tilde{y}|x_r}^2 \rangle | \tilde{x}$ for x_r drawn from $P(x)$. Reusing the estimate $g\bar{g}^T$ from the previous section, we can write the expectation of eqn (9) over X as

$$\langle \Delta \sigma_{\tilde{y}}^2 | \tilde{x} \rangle_{x, \tilde{y}} = \frac{g(\tilde{x})^T A^{-1} \bar{g} \bar{g}^T A^{-1} g(\tilde{x})}{S^2 + g(\tilde{x})^T A^{-1} g(\tilde{x})}. \quad (10)$$

3.3. Choosing an Optimal \tilde{x}

Given eqn (10), the problem remains of how to choose an \tilde{x} that maximizes it. One approach is to use selective sampling: evaluate a number of possible random \tilde{x} , then choose the one with the highest expected gain. This is efficient so long as the dimension of the action space is small. For high-dimensional problems, we may use gradient ascent to efficiently find good \tilde{x} . Defining $\bar{g} \equiv g(\tilde{x})$ and $\Gamma \equiv A^{-1} \bar{g} \bar{g}^T A^{-1}$, we can differentiate eqn (10) with respect to \tilde{x} to get the gradient

$$\nabla_{\tilde{x}} \langle \Delta \sigma_{\tilde{y}}^2 \rangle_{x, \tilde{y}} = \frac{2[(S^2 + \bar{g}^T A^{-1} \bar{g}) \bar{g}^T \Gamma - \bar{g}^T \Gamma \bar{g} \bar{g}^T A^{-1}]}{(S^2 + \bar{g}^T A^{-1} \bar{g})^2} \frac{\partial \bar{g}}{\partial \tilde{x}}. \quad (11)$$

We can “hillclimb” on this gradient to find an \tilde{x} with

a locally optimally expected change in average output variance.

It is worth noting that both of these approaches are applicable in continuous domains, and therefore well-suited to problems with continuous action spaces. Furthermore, the gradient approach is effectively immune to the overabundance of candidate actions in high-dimensional action spaces.

3.4. Two Caveats

3.4.1. Greedy Optimality. I have described a criterion for one-step, or *greedy* optimization. That is, each action/query is chosen to maximize the change in variance on the next step, without regard to how future queries will be chosen. The globally optimal, but computationally expensive approach would involve optimizing over an entire trajectory of m actions/queries. Trajectory optimization entails starting with an initial trajectory, computing the expected gain over it, and iteratively relaxing points on the trajectory towards optimal expected gains (subject to other points along the trajectory being explored). After the iteration has settled, the first point in the trajectory is queried, and the relaxation is repeated on the remaining part of the trajectory. Experiments using this form of optimization did not demonstrate significant improvement, in the average case, over the greedy method, so it appears that trajectory optimization may not be worth the additional computational expense, except in extreme situations (see Sollich (1994) for a theoretical comparison of greedy and globally-optimal querying).

3.4.2. Statistical Assumptions. The interpretation of $\sigma_{\tilde{y}}^2$ as a variance depends on a number of statistical assumptions, some of which may be untenable in real learning problems. The statistical interpretation of the OED framework described here assumes that the target is chosen statistically from a class of admissible functions that depend on the learner’s architecture. In contrast, reality usually dictates a fixed target function which we must try to learn, and we attempt to create an architecture for our learner that includes a reasonable approximation of this function.

This approach also assumes that all misfit of the training data is due to noise, which is assumed to be independent, identically distributed (i.i.d.) and, in the above formulation, Gaussian. In practice, much of the misfit may be nonstatistical in nature—it may be an inherent misfit between the target function and the learner’s architecture. Further, it is likely that what noise there is will not be independent, identically distributed, or Gaussian.

There is one more assumption which is more subtle and affects both the statistical and nonstatistical interpretations of OED: the variance estimates

obtained following eqn (4) are *local* estimates, based on the information matrix given the current (locally optimal) weight vector. It is assumed that this local estimate provides a good approximation of the global variance. In practice, there may be many such locally optimal weight vectors, all giving rather different solutions. If the variance between these solutions is significant (as work by Wolpert (1994), footnote 1, indicates it may be), the local estimate will be a poor approximation of the true variance. Current work by Paass and Kindermann (1995) attempts to compute global variances using Monte Carlo methods.

4. EXPERIMENTAL RESULTS

In this section, I describe two sets of experiments using optimal experiment design for error minimization. The first attempts to confirm that the gains predicted by optimal experiment design may actually be realized in practice, and the second applies OED to learning tasks with static and dynamic constraints. All experiments described in this section were run using feedforward networks with a single hidden layer of 20 units. Hidden and output units used the 0–1 sigmoid as a nonlinearity. No explicit form of regularization (such as weight decay) was used, nor was any noise added to the training sets. All runs were performed on the Xerion simulator (van Camp et al., 1993) using the default weight update rule (“Rudi’s conjugate gradient” with “Ray’s line search”).

4.1. Expected versus Actual Gain

It must be emphasized that the gains predicted by OED are *expected* gains. These expectations are based on the series of approximations detailed in the previous section, which may compromise the realization of any actual gain. In order for the expected gains to materialize, two “bridges” must be

crossed. First, the expected decrease in model variance must be realized as an actual decrease in variance. Second, the actual decrease in model variance must translate into an actual decrease in model MSE.

4.1.1. Expected Decreases in Variances \rightarrow Actual Decreases in Variance. The translation from expected to actual changes in variance requires coordination between the exploration strategy and the learning algorithm: to predict how the variance of a weight will change with a new piece of data, the predictor must know how the weight itself (and other weights in the network) will change. Using a black box routine like backpropagation to update the weights virtually guarantees that there will be some mismatch between expected and actual decreases in variance. Figure 3a plots the correlation between the expected change in output variance and actual change in output variance on the arm kinematics task. The results indicate that in spite of potential mismatch, the correlation between predicted and actual changes in variance is relatively good.

4.1.2. Decreases in Variance \rightarrow Decreases in MSE. A more troubling translation is the one from model variance to model correctness. Given the highly nonlinear nature of a neural network, local minima may leave us in situations where the model is very confident but entirely wrong. Due to high confidence, the learner may reject actions that would reduce its mean squared error and explore areas where the model is correct, but has low confidence. Evidence of this behavior is seen in the lower right corner of Figure 3b, where some actions which produce a large decrease in variance have little effect on E_{MSE} . This behavior appears to be a manifestation of the bias term discussed in Section 3; these queries reduce variance while increasing the learner’s bias, with no net decrease in error. While this demonstrates a weak

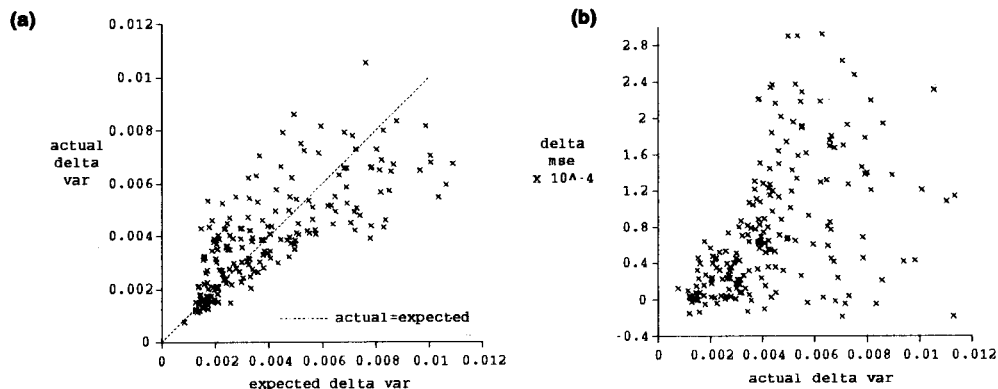


FIGURE 3. (a) Correlations between expected change in output variance and actual change in output variance. (b) Correlations between actual change in output variance and change in mean squared error. Correlations are plotted for a network with a single hidden layer of 20 units trained on 50 examples from the arm kinematics task.

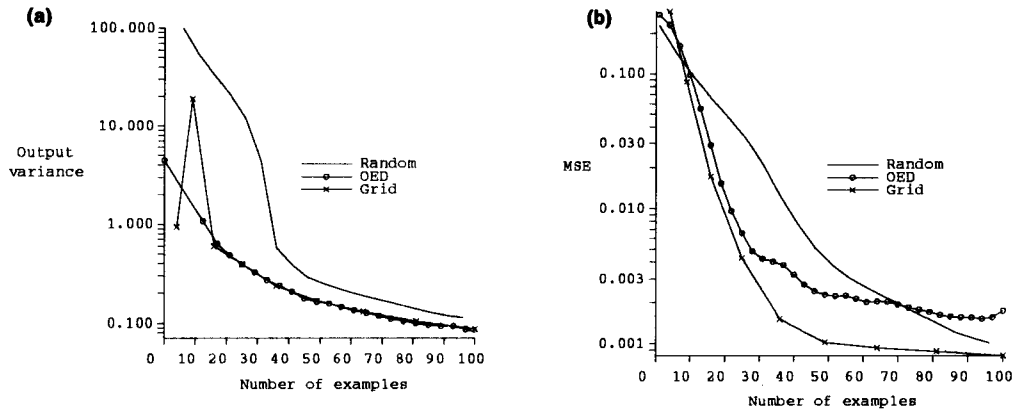


FIGURE 4. Querying with static constraints to learn the kinematics of a planar two-joint arm. (a) Variance using OED-based actions is better than that using random queries, and matches the variance of a uniform grid. (b) MSE using OED-based actions is initially very good, but breaks down at larger training set sizes. Curves are averages over six runs apiece for OED and grid learners, and 12 runs for the random learner.

point in the OED approach (which will be further illustrated below), there appear to be many situations where the approach is still useful.

4.2. Querying with Static Constraints

Here I consider a simple learning problem with static constraints: learning the forward kinematics of a planar arm from examples. The input $X = \{\Theta_1, \Theta_2\}$ specified the arm's joint angles, and the learner attempted to learn a map from these to the Cartesian coordinates $Y = \{C_1, C_2\}$ of the arm's tip. The "shoulder" and "elbow" joints were constrained to the 0–360° and 0–180° respectively; on each time step the learner was allowed to specify an arbitrary $\tilde{x} \in X$ within those limits.

For the greedy OED learner, \tilde{x} was chosen by beginning at a random point in X and hillclimbing the gradient of eqn (9) to a local maximum before querying. This strategy was compared with simply choosing \tilde{x} at random, and choosing \tilde{x} according to a uniform grid over X .⁹

I compared the variance and MSE of the OED-based learner with that of the random and grid learners. The average variance of the OED-based learner was almost identical to that of the grid learner and slightly better than that of the random learner (Figure 4b). In terms of MSE however, the greedy OED learner did not fare as well. Its error was initially comparable to that of the grid strategy, but flattened out at an error approximately twice that of the asymptotic limit (Figure 4b). This flattening appears to be a result of bias. As discussed in

Section 3, the network's error is composed of a variance term and a bias term, and the OED-based approach, while minimizing variance, appears in this case to leave a significant amount of bias (see discussion in Section 6).

4.3. Querying with Dynamic Constraints

For learning with dynamic constraints, I again used the planar arm problem, but this time with a more realistic restriction on new inputs. For the first series of experiments, the learner learned the kinematics by incrementally adjusting Θ_1 and Θ_2 from their values on the previous query. The limits of allowable movement on each step corresponded to constraints with known dynamics. The second set of experiments involved learning the dynamics of the same arm based on torque commands. The unknown next state of the arm corresponded to constraints with unknown dynamics.

4.3.1. Constraints with Known Dynamics. To learn the arm kinematics, the learner hillclimbed to find the Θ_1 and Θ_2 within its limits of movement that would maximize the stochastic approximation of eqn (10). On each time step Θ_1 and Θ_2 were limited to change by no more than $\pm 36^\circ$ and $\pm 18^\circ$ respectively.

I compared variance and MSE of the OED-based learner with that of an identical learner which explored by means of a random walk, and with a learner trained on a series of hand-tuned trajectories.

The greedy OED-based learner found exploration trajectories that, intuitively, appear to give good global coverage of the domain (see Figure 5). In terms of performance, the average variance over the OED-based trajectories was almost as good as that of the best hand-tuned trajectory, and both were far

⁹ Note the uniform grid strategy is not viable for incrementally drawn training sets—the size of the grid must be fixed before any examples are drawn. In these experiments, entirely new training sets of the appropriate size were drawn for each new grid.

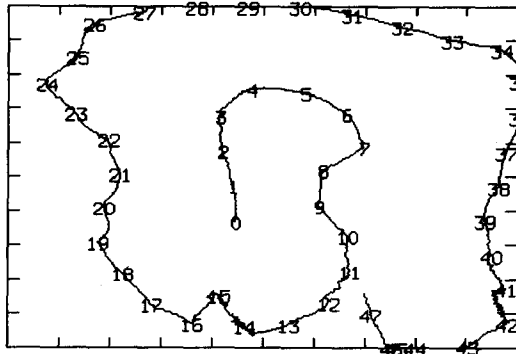


FIGURE 5. Querying with dynamic constraints: learning 2D arm kinematics. Example of OED-based learner's trajectory through angle-space.

better than that of the random exploration trajectories. In terms of MSE, the average error over OED-based trajectories was almost as good as that of the best hand-tuned trajectory, and again, both were far better than the random exploration trajectories (Figure 6). Note that in this case, bias does not seem to play a significant role. I discuss the performance and computational complexity of this task in greater detail in Section 5.

4.3.2. Constraints with Unknown Dynamics. For this set of experiments, I once again used the planar two-jointed arm, but now attempted to learn the arm dynamics. The learner's input $X = \{\theta_1, \theta_2, \dot{\theta}_1, \dot{\theta}_2, \tau_1, \tau_2\}$ specified the joint positions, velocities and torques. Based on these, the learner attempted to learn the arm's next state $Y = \{\theta'_1, \theta'_2, \dot{\theta}'_1, \dot{\theta}'_2\}$. As with the kinematics experiment, I compared random exploration with the greedy OED strategy described in the previous section. Without knowing the dynamics of the input constraints, however, there is no way to specify a preset trajectory.

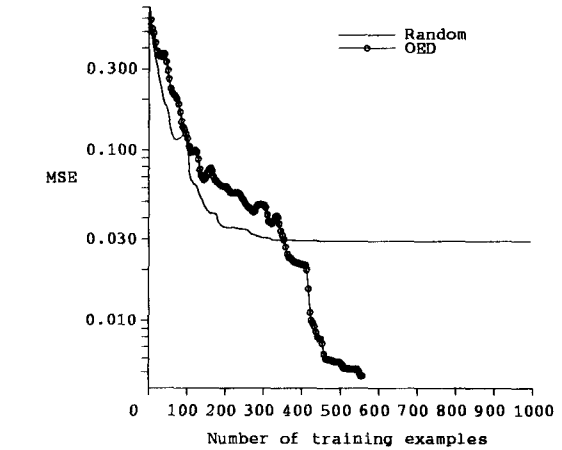
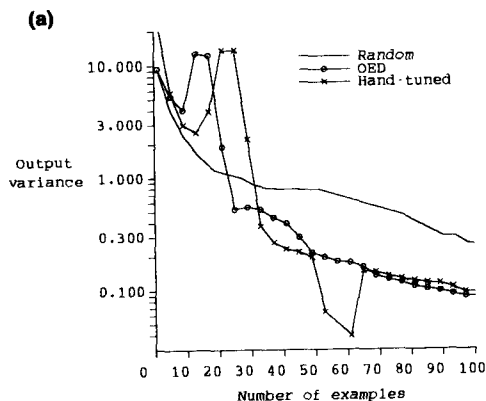


FIGURE 7. MSE of forward dynamic model for two-joint planar arm.

The performance of the learner whose exploration was guided by OED was asymptotically much better than that of the learner following a random search strategy (Figure 7). It is instructive to notice, however, that this improvement is not immediate, but appears only after the learner has taken a number of steps.¹⁰ Intuitively, this may be explainable by the assumptions made in the OED formalism: the network uses its estimate of variance of the current model to determine what data will minimize the variance. Until there are enough data for the model to become reasonably accurate, the estimates will be correspondingly inaccurate, and the search for "optimal" data will be misled. It would be useful to have a way of determining at what point the learner's estimates become reliable, so that one could explore randomly at first, then switch to OED-guided

¹⁰ This behaviour is visible in the other problem domains as well, but is not as pronounced.

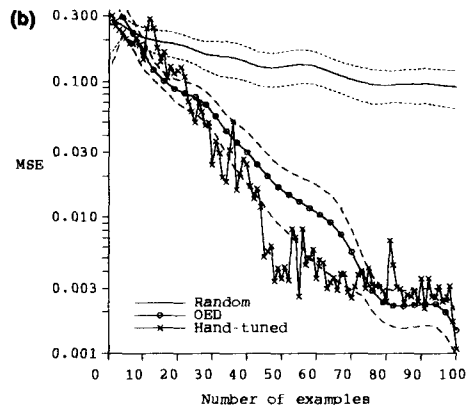


FIGURE 6. Querying with dynamic constraints: learning 2D arm kinematics. (a) Variance using greedy OED actions is better than that using random exploration, and matches the variance of the best hand-tuned trajectory. (b) MSE using greedy OED-based exploration is much better than that of random exploration and almost as good as that of the best hand-tuned trajectory. Curves are averages over five runs apiece for OED-based and random exploration. Dashed lines indicate standard error of estimates.

TABLE 1

Typical Computer Times, in seconds, for Operations Involved in Selecting New Data and Training. Number of Weights in Network = n , Number of Training examples = m , and Number of Reference Points (at which Variance or Gradient is Measured) = r . Time Constants are for Runs Performed on a Sparc 10 Using the Xerion Simulator

Operation	Constant	Order
Batch train	0.029	mn
Incremental train	0.093	n
Compute exact A	3×10^{-1}	mm^3
Compute approx. A	7.2×10^{-6}	mn^2
Invert to get A^{-1}	3.2×10^{-7}	n^3
Compute $\text{var}(x_r)$	5.0×10^{-6}	n^2
Compute $E[\Delta \text{var}(X) x]$	5.4×10^{-6}	mn^2
Compute gradient	1.9×10^{-5}	rn^2

exploration when the learner's model is accurate enough to take advantage of it.

5. COMPUTATIONAL COSTS AND APPROXIMATIONS

The major concern with applying the OED techniques described in this paper is computational cost. In this section I consider the computational complexity of selecting actions via OED techniques, and consider several approximations aimed at reducing the computational costs. These costs are summarized in Table 1, with the time constants observed for runs performed on a Sparc 10.

I divide the learning process into three steps: training, variance estimation, and data selection. I show that, for the case examined, in spite of increased complexity, the improvement in performance more than warrants the use of OED for data selection.

5.1. Cost of Training

I tested two training regimens for the OED-guided learners: batch training reinitialized after each new example was added, and incremental training, reusing the previous network's weights after each new example. While the batch-trained learners' performance was slightly better, their total training time was significantly longer than their incrementally trained counterparts (Figure 8).

5.2. Cost of Variance Estimation (eqn (6))

Variance estimation requires computing and inverting the Hessian. The inverse Hessian may then be used for an arbitrary number of variance estimates and must only be recomputed when the network weights are updated. The approximate Hessian of eqn (4) may be computed in time $O(mn^2)$, but the major cost remains the inversion. I have experimented with diagonal and block diagonal Hessians, which may be inverted quickly, but without the off-diagonal terms, the learner failed to generate reasonable training sets. Recent work by Pearlmutter (1994) offers a way to bring down the cost of computing the first term of eqn (6), but computing the second term remains an $O(n^3)$ operation. Buntine and Weigend (1994) discuss a number of alternative methods for computing the Hessian and its various approximations, and Paass and Kindermann (1995) discuss global methods for variance estimation that are independent of the Hessian.

5.3. Cost of Data Selection (eqns (9)–(11))

Computing eqn (9) is an $O(n^2)$ operation, which must be performed on each of r reference points, and must be repeated for each candidate \tilde{x} . Alternatively, the "moment-based" selection (eqn (10)) and gradient

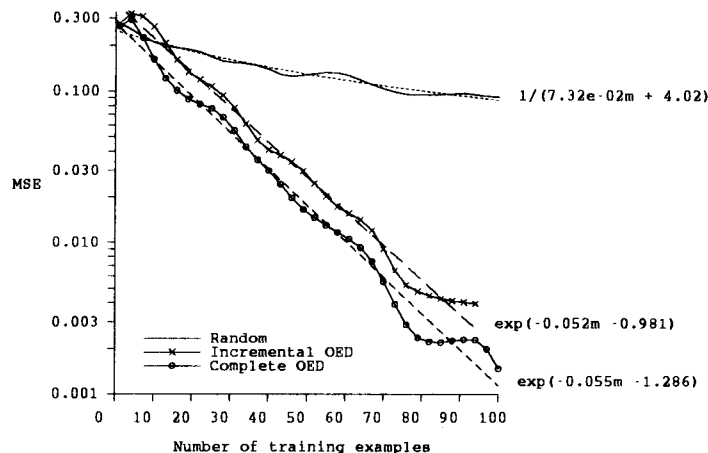


FIGURE 8. Learning curves for the kinematics problem from Section 4.2. Best fit functional forms are plotted for random exploration, incrementally-trained OED and OED completely retrained on new data set.

methods (eqn (11)) both require an $O(n^3)$ matrix multiplication which must be done once, after which any number of iterations may be performed with new \tilde{x} in time $O(n^2)$. Using Pearlmutter's approach to directly approximate $A^{-1}g(\tilde{x})$ would allow an approximation of eqn (10) to be computed in $O(n^2)$ times an "accuracy" constant. I have not determined what effect this time/accuracy tradeoff has on network performance.

5.4. The Payoff: Cost vs Performance

Obviously, the OED-based approach requires significantly more computation time than does learning from random examples. The payoff comes when relative performance is considered. I turn again to the kinematics problem discussed in Section 4.3.1. The approximate total time involved in training a learner on 100 random training examples from this problem (as computed from Table 1) is 170 seconds. For "full-blown" OED, using incremental training, the total time is 790 s. As shown in Figure 8, exploring randomly causes our MSE to decrease roughly as an inverse polynomial, while the various OED strategies decrease MSE roughly exponentially in the number of examples. To achieve the MSE reached by training on OED-selected data, we would need to train on approximately 3380 randomly selected data examples. This would take approximately 7500 seconds, over 2 hours! With this much data, the training time alone is greater than the total OED costs, so regardless of data costs, selecting data via OED is the preferable approach.

With the kinematics example there is the option of hand-tuning a learning trajectory, which requires no more data than the OED approach, and can nominally be learned in less time. This, however, required hours of human intervention to repeatedly re-run the simulations trying different preset exploration trajectories. In the dynamics example and in other cases where the state transitions are unknown, preset exploration strategies are not an option; we must rely on an algorithm for deciding our next action, and the OED-based strategy appears to be a viable, statistically well-founded choice.

6. DISCUSSION AND FUTURE WORK

The experiments described in this paper indicate that, for some tasks, optimal experiment design is a promising tool for guiding active learning in neural networks. It requires no arbitrary discretization of state or action spaces, and is amenable to gradient search techniques. The appropriateness of OED for exploration hinges on the two issues described in the previous two sections: the nature of the input

constraints and the computational load one is able to bear.

For learning problems with static constraints, the advantage of applying OED, or any form of intelligent active learning appears to be problem dependent. Random exploration appears to be reasonably good at decreasing both bias and variance (as seen in Section 4.2), while the variance minimization strategy may leave significant bias, and may in some cases even exacerbate it. For problems where learner bias is likely to be a major factor, then the advantages of the variance-minimizing approach are unclear. Although the variance-minimizing approach does appear to decrease bias in some cases, we have no sound understanding of why it does so. As a very rough approximation, both bias and variance are minimized by selecting examples where there are none; ignoring the finer distinctions, this overlap may account for the observed behavior on some problems.

The real advantage of the OED-based approach appears to lie in problems where the input constraints are dynamic, and where random actions fail to provide good exploration. Compared with arbitrary heuristics, the OED-based approach has the arguable advantage of being the "right thing to do", in spite of its computational costs.

The cost, however, is a major drawback. A decision time on the order of 1–10 seconds may be sufficient for many applications, but is much too long to guide real-time exploration of dynamical systems such as robotic arms. The operations required for Hessian computation and data selection may be efficiently parallelized; the remaining computational expense lies in retraining the network to incorporate each new example. The retraining cost, which is common to all on-line neural exploration algorithms, may be amortized by selecting queries/actions in small batches rather than purely sequentially. This "semi-batched" approach is a promising direction for future work.

Another promising direction, which offers hope of even greater speedups than the semi-batch approach, is switching to an alternative, entirely nonneural learner with which to pursue exploration.

6.1. Improving Performance with Alternative Learners

We may be able to bring down computational costs and improve performance by using a different architecture for the learner. With a standard feedforward neural network, not only is the repeated computation of variances expensive, it sometimes fails to yield estimates suitable for use as confidence intervals (as we saw in Section 4.1.2). A solution to both of these problems may lie in selection of a more amenable architecture and learning algorithm. Two

such architectures, in which output variances have a direct role in estimation, are mixtures of Gaussians (McLachlan & Basford, 1988; Nowlan, 1991; Specht, 1991; Ghahramani & Jordan, 1994) and locally weighted regression (Cleveland et al., 1988; Schaal & Atkeson, 1994). Both have excellent statistical modeling properties, and are computationally more tractable than feedforward neural networks. We are currently pursuing the application of optimal experiment design techniques to these models and have observed encouraging results (Cohn et al., 1995).

6.2. Active Elimination of Bias

Regardless of which learning architecture is used, the results in Section 4.2 make it clear that minimizing variance alone is not enough. For large, data-poor problems, variance will likely be the major source of error, but as variance is removed (via the techniques described in this paper), the bias will constitute a larger and larger portion of the remaining error.

Bias is not as easily estimated as variance; it is usually estimated by expensive cross validation, or by running ensembles of learners in parallel (see, e.g., Geman et al. (1992) and Connor (1993)). Preliminary work at bias minimization has been very encouraging (Cohn, 1995). Future work will need to include methods for efficiently estimating learner bias and variance and taking steps to ensure that both are minimized in an optimal manner.

REFERENCES

- Angluin, D. (1982). A note on the number of queries needed to identify regular languages. *Information and Control*, **51**, 76–87.
- Armstrong, B. (1989). On finding exciting trajectories for identification experiments. *International Journal of Robotics Research*, **8**, 28–48.
- Atkinson, A., & Donev, A. (1992). *Optimum experimental designs*. New York: Clarendon Press.
- Barto, A., Bradtke, S., & Singh, S. (1995). Learning to act using real-time dynamic programming. *Artificial Intelligence*, **72**, 81–138.
- Baum, E., & Haussler, D. (1989). What size net gives valid generalization? In D. Touretzky (Ed.), *Advances in neural information processing systems 1*. San Francisco, CA: Morgan Kaufmann.
- Baum, E., & Lang, K. (1991). Constructing hidden units using examples and queries. In R. Lippmann et al. (Eds.), *Advances in neural information processing systems 3*. San Francisco, CA: Morgan Kaufmann.
- Blumer, A., Ehrenfeucht, A., Haussler, D., & Warmuth, M. (1989). Learnability and the Vapnik–Chervonenkis dimension. *JACM*, **36**, 929–965.
- Box, G., & Draper, N. (1987). *Empirical model-building and response surfaces*. New York: John Wiley.
- Buntine, W., & Weigend, A. (1994). Computing second derivatives in feed-forward networks: a review. *IEEE Transactions on Neural Networks*, **5**(3), 180–188.
- Choueiki, M. H. (1995). *A designed experiment on the use of neural network models in short-term hourly load forecasting*. Doctoral Dissertation, Department of Industrial and Systems Engineering, Ohio State University, Columbus, Ohio.
- Cleveland, W., Devlin, S., & Grosse, E. (1988). Regression by local fitting. *Journal of Econometrics*, **37**, 87–114.
- Cohn, D. (1995). Minimizing statistical bias with queries, AI Lab. Memo 1552, Massachusetts Institute of Technology. Available by anonymous ftp from publications.ai.mit.edu.
- Cohn, D., & Tesauro, G. (1992). How tight are the Vapnik–Chervonenkis bounds? *Neural Computation*, **4**(2), 249–269.
- Cohn, D., Atlas, L., & Ladner, R. (1990). Training connectionist networks with queries and selective sampling. In D. Touretzky (Ed.), *Advances in neural information processing systems 2*, San Francisco, CA: Morgan Kaufmann.
- Cohn, D., Ghahramani, Z., & Jordan, M. (1995). Active learning with statistical models. In G. Tesauro et al. (Eds.), *Advances in neural information processing systems 7*. Cambridge, MA: MIT Press.
- Connor, J. (1993). Bootstrap methods in neural network time series prediction. In J. Alspector et al. (Eds.), *Proceedings of the International Workshop on Application of Neural Networks to Telecommunications*. Hillsdale, NJ: Lawrence Erlbaum.
- Craig, J. (1989). *Introduction to robotics*. New York: Addison-Wesley.
- Eisenberg, B., & Rivest, R. (1990). On the sample complexity of pac-learning using random and chosen examples. In M. Fulk & J. Case (Eds.), *ACM 3rd Annual Workshop on Computational Learning Theory*. San Francisco, CA: Morgan Kaufmann.
- Fedorov, V. (1972). *Theory of optimal experiments*. New York: Academic Press.
- Ford, I., Kitsos, C., & Titterton, D. M. (1989). Recent advances in nonlinear experimental design. *Technometrics*, **31**(1), 49–60.
- Freund, Y., Seung, H. S., Shamir, E., & Tishby, N. (1993). Information, prediction, and query by committee. In S. Hanson et al. (Eds.), *Advances in neural information processing systems 5*. San Francisco, CA: Morgan Kaufmann.
- Geman, S., Bienenstock, E., & Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, **4**, 1–58.
- Ghahramani, Z., & Jordan, NM. (1994). Supervised learning from incomplete data via an EM approach. In J. Cowan et al. (Eds.), *Advances in neural information processing systems 6*. San Francisco, CA: Morgan Kaufmann.
- Haussler, D. (1992). Generalizing the pac model for neural nets and other learning applications. *Information and Computation*, **100**, 78–150.
- Hwang, J. N., Choi, J., Oh, S., & Marks, R. (1990). Query learning based on boundary search and gradient computation of trained multilayer perceptrons. *International Joint Conference on Neural Networks 1990*, San Diego, June 17–21.
- Jacobs, R., Jordan, M., Nowlan, S., & Hinton, G. (1991). Adaptive mixtures of local experts. *Neural Computation*, **3**, 79–87.
- Jordan, M., & Rumelhart, D. (1992). Forward models: supervised learning with a distal teacher. *Cognitive Science*, **16**, 307–354.
- Kuperstein, M. (1988). Neural model of adaptive hand-eye coordination for single postures. *Science*, **239**, 1308–1311.
- Linden, A., & Weber, F. (1993). Implementing inner drive by competence reflection. In H. Roitblat et al. (Eds.), *Proceedings of the 2nd International Conference on Simulation of Adaptive Behavior*. Cambridge, MA: MIT Press.
- MacKay, D. (1992). Information-based objective functions for active data selection. *Neural Computation*, **4**, 590–604.
- McLachlan, G., & Basford, K. (1988). *Mixture models: inference and applications to clustering*. New York: Marcel Dekker.
- Moore, A. (1994). The parti-game algorithm for variable resolution reinforcement learning in multidimensional state-spaces. In J. Cowan et al. (Eds.), *Advances in neural information processing systems 6*. San Francisco, CA: Morgan Kaufmann.

- Nowlan, S. (1991). *Soft competitive adaptation: neural network learning algorithms based on fitting statistical mixtures*. (CMU-CS-91-126), School of Computer Science, Carnegie Mellon University, Pittsburgh, PA.
- Paass, G., & Kindermann, J. (1995). Bayesian query construction for neural network models. In G. Tesauero, et al. (Eds.), *Advances in neural information processing systems 7*. Cambridge, MA: MIT Press.
- Pearlmutter, B. (1994). Fast exact multiplication by the Hessian. *Neural Computation*, 6, 147–160.
- Plutowski, M., & White, H. (1993). Selecting concise training sets from clean data. *IEEE Transactions on Neural Networks*, 4, 305–318.
- Schaal, S., & Atkeson, C. (1994). Robot juggling: an implementation of memory-based learning. *Control Systems*, 14, 57–71.
- Schmidhuber, J., & Storck, J. (1993). *Reinforcement driven information acquisition in nondeterministic environments*. (Technical Report, in preparation), Fakultät für Informatik, Technische Universität München.
- Sollich, P. (1994). Query construction, entropy and generalization in neural network models. *Physical Review E*, 49, 4637–4651.
- Specht, D. (1991). A general regression neural network. *IEEE Transactions on Neural Networks*, 2(6), 568–576.
- Subrahmonia, J., Cooper, D. B., & Keren, D. (1992). *Practical reliable recognition of 2D and 3D objects using implicit polynomials and algebraic invariants*. (Technical Report LEMS-107), Division of Engineering, Brown University, Providence, RI.
- Thisted, R. (1988). *Elements of statistical computing*. New York: Chapman and Hall.
- Thrun, S. (1992). The role of exploration in learning control. In D. White, & D. Sofge (Eds.), *Handbook of intelligent control*. New York: Van Nostrand Reinhold.
- Thrun, S., & Möller, K. (1992). Active exploration in dynamic environments. In J. Moody et al. (Eds.), *Advances in neural information processing systems 4*. San Francisco, CA: Morgan Kaufmann.
- Valiant, L. (1984) A theory of the learnable, *Communications of the ACM* 27, 1134–1142.
- van Camp, D., Plate, T., & Hinton, G. (1993). The Xerion neural network simulator, Department of Computer Science, University of Toronto. For further information, send email to xerion@cs.toronto.edu.
- Whitehead, S. (1991). *A study of cooperative mechanisms for faster reinforcement learning*. (Technical Report 365), Department of Computer Science, University of Rochester, Rochester, NY.
- Wolpert, D. (1994). *Bayesian back-propagation over i-o functions rather than weights*. (Santa Fe Institute Tech Report). Available by anonymous ftp to archive.cis.ohio-state.edu as pub/neuroprose/wolpert.nips.93.ps.Z.

APPENDIX: SYSTEM IDENTIFICATION WITH NEURAL NETWORKS AND OED

Systems identification using OED has been successful on tasks

where the parameters of the unknown system are explicit, but for a neural network model, system identification is problematic. The weights in the network cannot be reasonably considered to represent real parameters of the unknown system being modeled, so there is no good interpretation of their “identity”. A greater problem is the observation that unless the network is fortuitously structured to be exactly the correct size, there will be extra unconstrainable parameters in the form of unused weights, about which it will be impossible to gain information. Distinguishing between unconstrainable parameters (which we wish to delete or ignore) and underconstrained parameters (about which we wish to get more information) is an unsolved problem. Below, I review the derivation of the “*D*-optimality” criterion appropriate for system identification (Fedorov, 1972; MacKay, 1992), and briefly discuss experiments selecting *D*-optimal data.

When doing system identification with a neural network, we are interested in minimizing the covariance of the parameter estimates \hat{w} . For the purposes of optimization, it is convenient to express σ_w^2 as a scalar. The most widely used scalar is the determinant $D = |\sigma_w^2|$, which has an interpretation as the “volume” of parameter space encompassed by the variance (for other approaches see Atkinson and Donev (1992)).

The utility of querying at \bar{x} , from a system identification viewpoint, may be expressed in terms of the expected change in the estimated value of *D*. The expected new value \bar{D} is

$$\bar{D} = |\bar{A}^{-1}| = \frac{S^2 |A^{-1}|}{S^2 + g(\bar{x})^T A^{-1} g(\bar{x})} = \frac{S^2 |\sigma_w^2|}{S^2 + \sigma_{y|\bar{x}}^2}, \quad (\text{A.1})$$

which, by subtraction from the original estimate *D* gives

$$\Delta D|\bar{x} = \frac{D\sigma_{y|\bar{x}}^2}{S^2 + \sigma_{y|\bar{x}}^2}. \quad (\text{A.2})$$

Equation (A.2) is maximized where $\sigma_{y|\bar{x}}^2$ is at a maximum, giving the intuitively pleasing interpretation that for system identification, parameter uncertainty is minimized by querying where our uncertainty is largest. Such queries are, in OED terminology, “*D*-optimal”.

Our experiments using the above criterion to select training data had limited success. For classification problems, the learner selected queries along perceived class boundaries and produced excellent performance. For regression problems such as the arm kinematics however, the learner performed poorly, attempting to select data at $x = \pm\infty$. One may interpret this as an attempt to select data that give the most leverage. While such a strategy is appropriate with a truly linear system, it performs poorly on a neural network, which is at most locally linear. These results are consistent with the comments at the beginning of this section, and with MacKay’s observation that, for learning $X \rightarrow Y$ mappings, the system identification criterion may be the “right solution to the wrong problem” (MacKay, 1992). The criterion addressed in Section 3, also mentioned by MacKay and explored in greater detail in this paper, appears to address the “right” problem.