

Active Feature-Value Acquisition

Maytal Saar-Tsechansky

McCombs School of Business, University of Texas at Austin, Austin, Texas 78712,
maytal@mail.utexas.edu

Prem Melville

IBM T.J. Watson Research Center, Yorktown Heights, New York 10598,
pmelvil@us.ibm.com

Foster Provost

Stern School of Business, New York University, New York, New York 10012,
fprovost@stern.nyu.edu

Most induction algorithms for building predictive models take as input training data in the form of feature vectors. Acquiring the values of features may be costly, and simply acquiring all values may be wasteful or prohibitively expensive. Active feature-value acquisition (AFA) selects features incrementally in an attempt to improve the predictive model most cost-effectively. This paper presents a framework for AFA based on estimating information value. Although straightforward in principle, estimations and approximations must be made to apply the framework in practice. We present an acquisition policy, sampled expected utility (SEU), that employs particular estimations to enable effective ranking of potential acquisitions in settings where relatively little information is available about the underlying domain. We then present experimental results showing that, compared with the policy of using representative sampling for feature acquisition, SEU reduces the cost of producing a model of a desired accuracy and exhibits consistent performance across domains. We also extend the framework to a more general modeling setting in which feature values as well as class labels are missing and are costly to acquire.

Key words: information acquisition; predictive modeling; active learning; active feature acquisition; data mining; machine learning; business intelligence; imputation; utility-based data mining

History: Received July 18, 2006; accepted April 25, 2008, by Ramayya Krishnan, information systems.

Published online in *Articles in Advance* January 28, 2009.

1. Introduction

...the shift from relying on existing information collected for other purposes to using information collected specifically for research purposes is analogous to primitive man's shifting from food collecting to agriculture... (Siegel and Fouraker 1960, p. 72)

Predictive models play a key role in numerous business intelligence tasks. Models are induced from historical data to predict customer behavior or to detect adversarial acts such as fraud. A critical factor affecting the knowledge captured by such a model is the *quality* of the information from which the model is induced—the “training data.” In the context of predictive modeling, the quality of information pertains to the training sample's composition, the accuracy of the values, and the number of unknown values.

For many predictive modeling tasks, potentially pertinent information is not immediately available but can be acquired at a cost. Traditionally, *information acquisition* and *inductive modeling* are addressed

independently; data are collected irrespective of the modeling objectives. However, information acquisition and predictive modeling in fact are mutually dependent: newly acquired information affects the model induced from the data, and the knowledge captured by the model can help determine what new information would be most useful to acquire (Simon and Lea 1974). We would like to take advantage of this relationship and develop feature value acquisition *policies* for predictive model induction—procedures for evaluating and selecting feature-value acquisitions that will be used for model induction. Mookerjee and Mannino (1997) address a similar problem, where costly feature values of a *test* instance for which *inference* is requested are unknown and are acquired sequentially given an existing knowledge base. Here we study a complementary problem in which values must be acquired for *induction*.

The availability of generic, effective, and computationally feasible information-acquisition policies for model induction can affect business practices by

transforming existing information-acquisition practices and by changing the manner by which firms interact with consumers. As an example, consider the generation of personalized recommendations to customers. Often, a recommender system's underlying predictive model employs customers' ratings of prior purchases as predictors of a customer's preference for a product she has not yet purchased. The availability of many ratings from a large number of customers is critical for successful induction of an accurate model of consumer preferences. However, without costly incentives, most consumers rarely provide this valuable feedback. To improve the model's predictive accuracy, it is infeasible to acquire feedback from all consumers about all products, even those they have already purchased. A better acquisition policy would determine which ratings from which customers would be most cost-effective to acquire via costly incentives, in order to obtain the desired modeling objective for the lowest cost (Huang 2007). Similar scenarios emerge in other modeling tasks where missing feature values can be acquired at a cost. These include modeling of medical treatment effectiveness and diagnosis from medical databases, where patient information, such as details on prior hospitalizations and prior medical tests, is notoriously incomplete. Intelligent information-acquisition policies can also dramatically change already established information-acquisition models: firms acquire bundles of psychographic, consumption, and lifestyle data periodically from third-party suppliers, such as Axiom, to support business intelligence modeling for tasks such as risk scoring, customer retention, and personalized marketing. As with other information goods, a firm should consider how to bundle and price information on consumers. Effective acquisition policies will enable firms to identify and acquire different types of information for different consumers at potentially different prices, to enhance cost-effective modeling. These capabilities may also enable small firms to reap the benefits from business intelligence modeling, allowing them to enrich their potentially limited data by selectively acquiring useful information.

Given training data with missing feature values, an arbitrary classification-model induction algorithm, a set of prospective feature-value acquisitions, and the cost of acquiring each specific feature value (the cost of features may vary feature to feature and instance to instance), the general active feature-value acquisition (AFA) problem is to acquire feature values so as to obtain a desired performance level for minimum cost. In this paper we consider performance to be some function of the model's generalization accuracy.¹

However, because we do not know a priori the population under consideration, the generalization performance cannot be known exactly, and we must estimate it from a sample. Thus, AFA policies cannot be provably optimal, so we employ a heuristic measure of the performance from feature-value acquisition.

Even given a heuristic measure of performance, in principle, identifying the feature-value set that yields the desired performance objective at a minimum cost requires considering all possible sets of prospective acquisitions. Unfortunately, it is not feasible to compute this for most interesting problems. Moreover, given a finite sample, both statistical learning theory and practical experience tell us that more search through possible models often leads to worse performance, because of problems of multiple comparisons (Vapnik 1998, Jensen and Cohen 2000). We will revisit this issue in §5. We therefore revise the objective of AFA for this paper as follows. Given a performance measure, we aim to identify the individual feature value to acquire next, in order to achieve the greatest improvement in the performance measure per unit cost, which implies a greedy, myopic acquisition policy.

The primary contributions of this paper are a general framework for addressing AFA and a specific method for solving AFA problems based on appropriate heuristics. To our knowledge, no prior work² addresses general AFA. We propose an acquisition policy that produces acquisition schedules iteratively based on estimates of the expected utility from different potential acquisitions. In principle this is straightforward, but the AFA setting renders utility estimation particularly challenging: estimations often must be made based on little available information and ought to be sensitive so as to capture the benefit to induction of individual feature values. Further, because the space of possible acquisitions can be immense, estimating the value of each potential acquisition may be computationally infeasible. We develop and study empirically the impact of measures for capturing the value to induction of single feature values in the presence of scarce data and propose different mechanisms to reduce the complexity of the estimation.

This expected-utility approach has several important advantages. The framework is general and can be applied to derive an acquisition schedule for any induction technique. This is important because due to the inherent bias of different modeling techniques and professional regulations in some industries, no single technique is applied across all problems. Another

¹The model's expected accuracy over the population under consideration.

²With the exception of a short paper on our preliminary studies (Melville et al. 2005).

important advantage of the expected-utility approach is that it can be applied to improve any utility function derived from the model's predictive performance, such as estimated generalization accuracy, expected profit in a particular setting, or the expected cost of model error. Finally, this approach can utilize information about the varying cost of information to derive an acquisition schedule, not assuming that the cost of acquiring an unknown value is fixed for all features and/or all instances. Experimental results demonstrate that the resulting method provides significantly better models for a given cost than those obtained with other acquisition policies. Because the method utilizes acquisition cost information, it is particularly advantageous in challenging tasks for which there is significant variance across potential acquisitions with respect to their informativeness and their cost.

Finally, another contribution of this paper is an extension of the policy to a more general acquisition problem. For some modeling tasks, *class labels* (i.e., dependent variables' values) are missing as well as feature values, and either or both may be acquired at cost. We show that because our framework estimates the value to induction of different acquisitions, it allows the dependent variable to be treated as yet another feature. Thus, the AFA framework and method can be extended directly to address this new problem. In practice, the method interleaves the acquisition of class labels and feature values, based on the marginal expected value from each acquisition; we show it to be superior both to uniform acquisitions and to policies that consider the acquisition of only feature values or only class labels.

2. Active Feature-Value Acquisition

Assume a classifier induction problem where each instance is represented with n independent variables plus a discrete, dependent "class variable." The available data set of m instances can be represented by the "incomplete" matrix F , where $F_{i,j}$ corresponds to the value of the j th feature of the i th instance, which may be missing. Missing elements in the matrix F represent missing feature values that can be acquired at a cost. In general, the cost of different feature values may vary, depending on the nature of the particular feature or of the instance for which the information is missing.

Algorithm 1. General Active Feature-Value Acquisition Framework

Given: F : initial (incomplete) instance-feature matrix;
 $Y = \{y_i, i = 1, \dots, m\}$: class labels for all instances;
 T : training set = $\langle F, Y \rangle$; L : classifier induction algorithm; β : size of query batch; C : cost matrix for all instance-feature pairs;

1. Initialize set of possible queries Q to $\{q_{i,j}: i = 1, \dots, m; j = 1, \dots, n; \text{ such that } F_{i,j} \text{ is missing}\}$
2. Repeat until stopping criterion is met
3. Induce a classifier, $M = L(T)$
4. $\forall q_{i,j} \in Q$ compute $\text{score}(q_{i,j}, c_{i,j}, L, T)$
5. Select the subset, S , of β feature value with the highest scores
6. $\forall q_{i,j} \in S$
7. Acquire values for $F_{i,j}$
8. Remove S from Q
9. End Repeat
10. Return $M = L(T)$

We present an iterative, sequential acquisition framework, where at each acquisition phase, alternative acquisitions are evaluated to acquire the value of $F_{i,j}$ at the cost $C_{i,j}$ that provides the largest improvement per unit cost in the performance objective. The iterative framework for AFA is presented in Algorithm 1. The framework is independent of the classification modeling technique; it is given a learner, L , which includes a model induction algorithm and a missing value treatment to allow for induction from the incomplete matrix F .³ At each phase a "score" is estimated and assigned to each potential acquisition, reflecting the estimated added value per unit cost of the acquisition. The acquisition with the highest score is selected, and the corresponding feature value is acquired; a particular approach for assigning scores to potential acquisitions will be described in detail in §2. Once a value is acquired, the training data and the information acquisition cost are appropriately updated, and this process is repeated until some stopping criterion is met, e.g., a desirable model accuracy has been obtained. Often, many values must be acquired to obtain a desired performance level. To reduce the computational burden or based on domain constraints, at each iteration an AFA policy may acquire a "batch" of $\beta \geq 1$ values. As before, computing the set that will result in the greatest improvement in the heuristic measure is computationally complex. In this paper, we select the values with the highest individual scores; the sensitivity to this choice is examined in §3.3.

We now present a method for AFA based on computing the value of the information that may be acquired. The central component of the computation also presents the main difficulties with its implementation: the computation of the value of information prior to acquisition, when only partial knowledge about the

³ Induction algorithms either include an internal mechanism for incorporating instances with missing feature values (Quinlan 1993) or require that missing values be imputed first. Henceforth, we assume that the induction algorithm includes or is coupled with some treatment for instances with missing values.

acquired information is available. We discuss three difficulties with the computation and present approximation techniques to address these difficulties. Together they comprise the proposed AFA method: *sampled expected utility* (SEU).

We estimate the value of a potential acquisition by its expected marginal contribution to predictive performance. Because the true value of the missing feature is unknown prior to its acquisition, it is necessary to estimate the potential impact of an acquisition for different possible acquisition outcomes. The acquisition with the highest information value will be the one that results in the maximum utility in expectation, given a model, a model induction algorithm, and a particular utility function. For the last, the objective may be to maximize the model’s generalization accuracy, to maximize future profit, to minimize the costs incurred due to incorrect predictions, etc. A utility score captures the expected improvement from each potential acquisition. Assuming feature j has K distinct possible values v_1, \dots, v_K , the expected utility of the acquisition $q_{i,j}$, or “query” for short, is given by

$$E(q_{i,j}) = \sum_{k=1}^K \mathcal{U}(F_{i,j} = v_k) P(F_{i,j} = v_k), \quad (1)$$

where $P(F_{i,j} = v_k)$ is the probability that $F_{i,j}$ has the value v_k , and $\mathcal{U}(F_{i,j} = v_k)$ is the utility of knowing (via acquisition) that the feature value $F_{i,j}$ is v_k . The utility $\mathcal{U}(\cdot)$ is the marginal improvement in performance per unit of acquisition cost:

$$\mathcal{U}(F_{i,j} = v_k) = \frac{\mathcal{A}(F, F_{i,j} = v_k)}{C_{i,j}}, \quad (2)$$

where $\mathcal{A}(F, F_{i,j} = v_k)$ is the change in value to induction from augmenting F with $F_{i,j} = v_k$, and $C_{i,j}$ is the cost of acquiring $F_{i,j}$. This *expected utility* policy therefore corresponds to selecting the query that will result in the estimated largest increase in performance per unit cost in expectation. If all feature costs are equal, this corresponds to selecting the query that would result in the classifier with the highest expected performance. Otherwise, *expected utility* allows several low-yield, high-margin acquisitions to be selected instead of one higher-yield acquisition with less expected improvement per unit cost.

In principle, this approach would allow the estimation of the value of each possible acquisition and then the selection of acquisitions by ranking them by their information-value estimates. However, there are significant hurdles to its practical implementation. We introduce the challenges next, and then address each in turn in the following three subsections.

Challenge 1. Estimating contribution to induction. As outlined in Equation (2), for each query

($\forall q_{ij} \in Q$) computing expected utility requires the estimation of the value, $\mathcal{A}(\cdot)$, to induction from the acquisition. Here we assume classification accuracy to be the performance metric of interest; as discussed, the framework applies to other goals such as minimizing misclassification cost or maximizing profit. As we will see, to estimate the expected improvement in classification performance, it is necessary to detect expected changes in the modeling technique’s average class probability estimation that are conducive to improved classification accuracy. It turns out that the obvious measure, classification accuracy itself, is not sensitive enough to such changes.

Challenge 2. Estimating value distributions. For estimating the expected contribution of different acquisitions, a prerequisite is to estimate the conditional distribution $P(F_{i,j} = v_k)$ for each missing value of $F_{i,j}$, as needed in Equation (1). We must identify an estimation mechanism appropriate for the AFA setting; many feature values may be missing, thereby rendering some modeling mechanisms more effective than others. For example, some mechanisms require that missing predictors or their distributions be estimated to produce a prediction. This adds yet another layer of estimation, which may undermine the model’s prediction, sometimes significantly (Saar-Tsechansky and Provost 2007).

Challenge 3. Reducing the consideration set. Even if all unknown values were estimated accurately, selecting the best from *all* potential acquisitions would require estimating the utility of, in the worst case, mn queries. This would be very expensive computationally and is infeasible for most interesting problems.

2.1. Estimating an Acquisition’s Contribution to Performance

Let us consider the measure to be used for estimating the value to induction from an acquisition, $\mathcal{A}(F, F_{i,j} = v_k)$. Let us assume for this discussion that acquiring new information aims to improve the model’s classification accuracy, for a binary classification problem (thus, the decision threshold for maximum a posteriori classification is 0.5). Assuming that the model’s contribution to induction is computed by averaging over a set of hold-out examples, this suggests a simple criterion for identifying an effective utility measure—prefer acquisitions that improve estimated accuracy. Let $A(f)$ denote the value to induction estimated for a single hold-out instance, and where the classifier’s estimated probability that the corresponding instance belongs to the true class is f .

CRITERION 1. For a given hold-out instance, $A(f_1) > A(f_2)$ if $f_1 > \theta$ and $f_2 < \theta$, where f_i refers to the model’s estimated probability that the given instance belongs to the true class, θ is the decision boundary, and $a > b$ denotes that a is better than b .

An obvious measure for this contribution is the model’s classification accuracy itself (i.e., the estimated generalization accuracy) over the augmented sample F . However, as we will show, classification accuracy does not capture fine-grained changes in the models; therefore, we would like a more sensitive measure for evaluating the benefits from prospective acquisitions.

To understand why, it is necessary to examine the dynamics of the modeling setting. Specifically, the training data—and therefore the models induced—continuously change as new information is acquired. Rather than examine the classification performance of a particular model induced from one version of the data, it is useful to examine how new acquisitions affect the *distribution* of estimations induced from different likely variations of the training data. Friedman’s analysis of the relationship between training data and classification error (Friedman 1997) examines how changes in an induction technique’s average estimation of the *probability* of the true class affect the likelihood of classification error. For the sake of discussion, assume binary classification and let $f(y|x)$ and $\hat{f}(y|x)$ denote the actual probability and the model’s estimated probability that an instance belongs to class y , respectively, where x is the input vector of observable attributes. Following Friedman’s analysis, the probability that the predicted class \hat{y} is estimated (erroneously) not to be the most likely class y can be approximated with a standard normal distribution by

$$P(\hat{y} \neq y) = \tilde{\Phi} \left[\text{sign}(f - 1/2) \frac{E\hat{f} - 1/2}{\sqrt{\text{var}\hat{f}}} \right], \quad (3)$$

where $\tilde{\Phi}$ is the upper tail area of the standard normal distribution, $\text{var}\hat{f}$ denotes the estimation variance resulting from variations in the training sample, and $E\hat{f}$ denotes the mean of the probability estimation $\hat{f}(y|x)$ generated by models induced from different variations of the training sample. Henceforth, we refer to $E\hat{f}$ and to $\text{var}\hat{f}$ as the average probability estimation and estimation variance, respectively. When the average probability estimation leads to incorrect class prediction, and given a certain estimation variance, the likelihood of incorrect classification decreases as the average probability estimation of the true class increases. Equation (3) also reveals that the likelihood of correct classification can improve when the average probability estimation *already* leads to a correct prediction. As shown in Equation (3), given a certain estimation variance, $\text{var}\hat{f}$, if the true class probability f and the average probability estimation $E\hat{f}$ lead to the same (correct) classification, then the further from the decision boundary the average probability estimation

$E\hat{f}$ is, the higher the likelihood is of a reduction in classification error. This is because it is less probable for the estimation variance to cause an erroneous classification. This analysis suggests two criteria for an effective measure A of the utility from an acquisition. First, A should favor more extreme (correct) estimates of class membership probability.

CRITERION 1’. For a given hold-out instance, $A(f_1) > A(f_2)$ iff $f_1 > f_2$, where f_i refers to the model’s estimated probability that the given instance belongs to the true class.

This is a specialization of Criterion 1; Criterion 1 will always hold if Criterion 1’ does (but not vice versa).

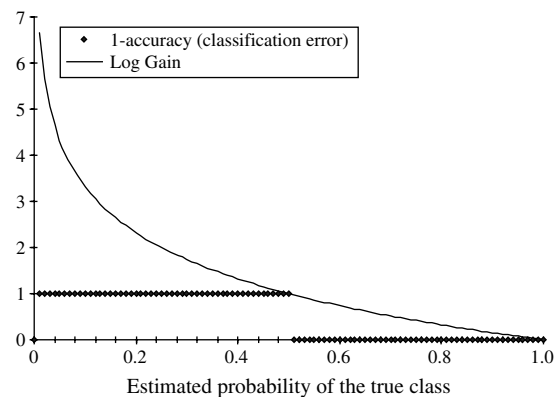
Second, for a correctly classified example, for a fixed-size change in the estimate of class membership probability, A should favor changes to estimates nearer to the decision boundary θ (note that in the analysis above this boundary is assumed to be 0.5):

CRITERION 2. For a given hold-out instance, $A(f_1) - A(f_1 + \Delta) > A(f_2) - A(f_2 + \Delta)$, $\forall \Delta: 0 < \Delta \leq \min(1 - f_1, 1 - f_2)$, and $\theta < f_1 < f_2$.

Based on these two criteria, we can assess the adequacy of different possibilities for the utility measure A , including intuitive alternatives such as estimated accuracy (error rate) or the estimate \hat{f} of the probability of class membership itself. Specifically, classification accuracy is not an adequate AFA utility measure because it does not satisfy either Criterion 1’ or Criterion 2 completely. This is illustrated for binary classification by the diamond line in Figure 1, which shows classification error (1-accuracy) as a function of the model’s estimated probability of the *true* class, assuming maximum a posteriori classification.

Let us consider an alternative AFA utility measure, *Log Gain* (LG) (also known as cross-entropy). For a model induced from a training set F , let $\hat{f}_F(y|x)$ be the probability estimated by the model that instance x belongs to class y , and let $\delta(A)$ be an indicator function such that $\delta = 1$ if A is the correct class and $\delta =$

Figure 1 Log Gain and Classification Error vs. the Probability of the True Class



INFORMS holds copyright to this article and distributed this copy as a courtesy to the author(s). Additional information, including rights and permission policies, is available at http://journals.informs.org/.

0 otherwise. Let $LG(x) = \sum_y -\delta(y) \log_2 \hat{f}_F(y | x)$; Log Gain is “better” as its value decreases. Consider an evaluation data set of t instances; let the value of induction from an acquisition resulting in a training set F be given by the sum of Log Gains over these t instances: $\mathcal{A}(F) = \sum_{e=1}^t LG(x_e)$. Hence, for each value V_k that feature $F_{i,j}$ can take, we would induce a model from the augmented data set and compute this sum of Log Gains. As illustrated in Figure 1, Log Gain satisfies both criteria. It captures important changes in the model’s estimation following an acquisition, which will allow the AFA policy to focus on acquisitions that decrease the likelihood of classification error: using Log Gain will result in higher scores for acquisitions that increase the average probability estimation of the true class when, on average, the model’s class prediction is incorrect. It will also promote acquisitions that lead to more extreme probability estimations when the model’s class prediction is accurate—reducing the risk of erroneous classification as new information is added to the training sample. In principle, other measures that promote these two objectives should benefit the AFA policy as well.

2.2. Estimating Feature-Value Distribution

Let us now address the second term in Equation (1). We need to estimate the conditional probability distribution of a missing feature value given the known values. For each feature j , the probability $P(F_{i,j} = v_k)$ in Equation (1) will be inferred from a model $M_{i,j}$, based on the other information available about instance i (what is known about the other features and the class).

Unfortunately, because of the AFA setting, the instances to which the feature-distribution-estimation model $M_{i,j}$ would be applied may have many missing values. Considering an arbitrary predictive model, predictors whose values are required for inference (rather than for induction) may not be available. The loss in predictive accuracy stemming from the need to estimate missing predictors’ values or their distributions can be avoided if the model incorporates only predictors whose values are known for this instance (Saar-Tsechansky and Provost 2007). However, for AFA this would entail considering a tremendous number of combinations, as at any point in an acquisition schedule different instances may include widely different sets of known feature values.

For the main results of this paper, to estimate the feature-value distribution model $M_{i,j}$ we employ only one predictor: the class variable. This is because for our setting, the class is guaranteed to be known at inference time. In principle, one can employ any set of known predictors to estimate a missing feature-value distribution. In §3 we validate empirically the benefits of relying on known predictors exclusively.

Conditioning on the class variable outperforms a simpler model that does not condition the missing feature distribution at all, because the latter does not take into account any instance-specific information in the estimation. A straightforward application of more complex modeling that captures the interactions among predictors but requires the estimation of unknown predictors for inference does not improve AFA performance.

2.3. Reducing the Consideration Set

Estimating the expectation $\hat{E}(\cdot)$ for each query, $q_{i,j}$, requires training one classifier for each possible value of $F_{i,j}$. Therefore, exhaustively evaluating all possible queries is infeasible for most interesting problems. One way to make this exploration tractable is by applying a fast method to identify a subset of all the possible queries that will subsequently be considered for acquisition. In particular, let the exploration parameter α ($1 \leq \alpha \leq mn/\beta$) control the size of the sample to be considered for acquisition. (Recall that β is the batch size, m the number of examples, and n the number of features.) To acquire a batch of β queries, first a subsample of $\alpha\beta$ queries is selected from the available pool of prospective acquisitions; then the expected utility of each query in this subsample is evaluated using Equation (1). The value of α can be set depending on the amount of time the user is willing to spend on this process and the effectiveness of the selection scheme. We consider two fast methods to identify a subset of queries.

The first approach, *uniform sampling* (US), identifies a representative subset of missing feature values via a uniform random sample of queries. However, when the consideration set is drawn uniformly at random, particularly informative acquisitions may be left out of the consideration set. An alternative approach is to limit the consideration set to a subset of queries that are more likely to be informative for model induction than a query drawn at random. In particular, we propose selecting the consideration set of queries from particularly informative *instances*. This invokes the subproblem of what then constitutes an *informative* instance for model induction. We conjecture that acquired feature values are more likely to have an impact on classification accuracy when the acquired values belong to a *misclassified* example and, as such, embed predictive patterns that are not consistent with the current model. Next, correctly classified instances are more informative if their class prediction is uncertain. The use of uncertainty for active data acquisition originated in work on optimum experimental design (Federov 1972) and has been extensively applied in the active learning literature (Cohn et al. 1994, Saar-Tsechansky and Provost 2004). For a probabilistic model, a lack of discriminative patterns results in

uncertain predictions, for which the model assigns similar likelihoods for class membership of different classes. Formally, for an instance x , let $P_y(x)$ be the estimated probability that x belongs to class y , as predicted by the model. Then the uncertainty score is given by $P_{y_1}(x) - P_{y_2}(x)$, where $P_{y_1}(x)$ and $P_{y_2}(x)$ are the first-highest and second-highest predicted class probability estimates respectively. Motivated by this reasoning, *error sampling* (ES) ranks informative instances higher if they are misclassified by the current model. Next, ES ranks instances in increasing order of the uncertainty score.

We call the approaches in which *uniform sampling* and *error sampling* are used to reduce the set of missing values considered for acquisition, *sampled expected utility-US* (SEU-US) and *sampled expected utility-ES* (SEU-ES), respectively.

3. Experimental Evaluation

We now present a comprehensive set of experiments that demonstrate the efficacy of AFA and the benefits of the measures described in §2.

3.1. Objectives and Methodology

We begin with the key empirical question of whether feature values can be acquired cost-effectively with AFA, compared with a default policy in which a representative set of feature-value acquisitions is drawn uniformly at random. Next, we present extensive empirical results that carefully examine the benefits of the measures we propose for AFA, compared with alternatives. These evaluations provide empirical support to the arguments in §2 regarding measures that are likely to be particularly effective for AFA. We then examine the upper-bound performance that could be obtained with an omniscient “oracle” and how close SEU is to this performance. In addition, we examine which of the two measures that SEU employs is closest to an omniscient measure for the corresponding quantity. These results suggest how improvements in each of SEU’s estimations can contribute to SEU’s overall performance so as to approach the performance of the oracle. Finally, we perform sensitivity analyses, exploring how different settings affect SEU’s performance. These include SEU performance with different feature-value cost structures and parameters that determine the size of the initial sample provided to SEU, the number of examples it considers for acquisition, and the number of examples acquired in each acquisition phase.

3.1.1. Primary Results. To address the first question, we compare, as a function of acquisition cost, the classification performance obtained by the policies SEU-ES, SEU-US, and a policy (uniform) that selects acquisitions uniformly at random. This study also

aims to examine the merits of the two approaches, US and ES, which we propose for reducing the set of prospective acquisitions considered by SEU.

We then demonstrate empirically the properties of the utility measure proposed for AFA in §2. The ability of SEU to rank potential acquisitions accurately will be affected by the accuracy of its estimates of the quantities in Equation (1), namely, the value to induction from a prospective acquisition of a value $F_{i,j} = v_k$, ($\mathcal{A}(F, F_{i,j} = v_k)$) and the estimated distribution of values for each unknown feature ($P(F_{i,j} = v_k)$). Specifically, in §2.1 we showed analytically that Log Gain effectively captures important changes in the estimated class probabilities following an acquisition that affect the likelihood of erroneous classification. Thus, Log Gain improves SEU’s ability to identify acquisitions that are particularly likely to reduce classification error. In contrast, using classification accuracy as the utility measure does not capture changes in probability estimation, except when the estimated class membership also changes. To demonstrate this effect on SEU performance, we compare SEU with a modified policy that employs classification accuracy on the training set to estimate the value of prospective acquisitions. We refer to this policy as *SEU-accuracy*. We also examine our choice for estimating the probability distribution over the values that a prospective acquisition may produce. As discussed in §2.2, prior research has concluded that employing only predictors whose values are known during inference improves prediction significantly. Based on these findings, in this study we conditioned the estimation on the (known) class label. To validate the merits of this approach, we compare it with two alternative methods that reflect two extremes with respect to reliance on predictors and their availability at inference time. The first is a simple approach that computes the *unconditional frequency* of feature values, based simply on their frequency in the training data. We refer to this approach as *SEU-frequency*. In the second approach, we use tree induction, employing all other features and the class label as predictors to estimate the probability distribution of a prospective acquisition. This approach, *SEU-DT*, aims to capitalize on the interactions between predictors for inference. However, as discussed in §2.2, inference is likely to suffer if predictors whose values are unknown must be estimated at inference time. The conditional distribution approach we employ in SEU lies between these two extremes—rather than relying on unknown values or estimating a simple unconditional distribution, SEU conditions the estimation on the known label.

3.1.2. Oracle Policies. To derive an upper bound for SEU’s performance, we employ an omniscient policy (the *oracle*) that knows the true values of missing features to determine the feature-value acquisition

that will lead to the greatest improvement in generalization performance. In addition, we assume the oracle has access to the held-out test data to compute the actual improvement in Log Gain following an acquisition. As with SEU, to render the evaluation feasible, rather than evaluate all possible acquisitions, the oracle selects the best acquisition among a sample of $\alpha\beta$ prospective acquisitions. Both policies select prospective acquisitions from the same set of prospective acquisitions.

We also present experiments that decompose the advantages conferred by the oracle over the imperfect estimations performed by SEU. Specifically, we decompose the relative advantages into the oracle's knowledge of the true model's performance measured over the held-out test data compared with SEU's estimation over the training set, and the oracle's knowledge of the true values of prospective acquisitions compared with SEU's *estimation* of the *expected* benefits from prospective acquisitions over all possible values that a missing feature may have. Recall that in SEU we estimate the distribution of a missing value to compute the benefits from its acquisition *in expectation*. If the actual values of prospective acquisitions were known, one could compute the benefit to model induction from acquiring the corresponding value directly rather than in expectation. To evaluate the upper-bound performance that can be obtained by SEU if the actual values of missing features were known, we constructed a new policy, the *feature oracle*, that has access to the true values of prospective acquisitions for the purpose of evaluating acquisitions. Both SEU and the feature oracle estimate the benefit to induction over the training set and evaluate the same set of prospective acquisitions in each acquisition phase. To evaluate the benefits from assessing the model's accuracy directly over the test data, we compare SEU's performance to that of the *performance oracle*—a policy in which the improvement in Log Gain is computed on the test data. We fix all other components of the policies so that the performance oracle and SEU employ the same measure to estimate feature-value distributions and evaluate the same consideration set of prospective acquisitions at each acquisition phase. We compare SEU to the *performance oracle* and *feature oracle* to estimate how improvements in each of SEU's estimations can contribute to SEU's overall performance so as to approximate the upper-bound performance.

3.1.3. Sensitivity Analyses. Finally, we explore the performance of SEU under different settings. The first of these evaluations considers the robustness of SEU's performance under different feature-value acquisition costs. SEU also incorporates several parameters, such as the size of the sample of feature values, whose contributions to learning are evaluated

by SEU, the number of feature values acquired in each acquisition phase, and the size of the initial sample provided to SEU for evaluating the contributions of prospective acquisitions. We explore in turn how each of these design parameters affects SEU's performance (Langley 2000, Hevner et al. 2004).

3.1.4. Experimental Setup. The empirical evaluations are performed over a set of data sets from a variety of domains. Four data sets⁴—Expedia, eToys, Priceline, and QVC—contain information about Web users and their visits to large retail websites. The target (dependent) variable indicates whether a user made a purchase during a visit. The predictors describe customers' surfing behaviors at the site as well as at other sites over time. We induce models to estimate whether a purchase will occur during a given session and employ the acquisition policies to estimate which unknown feature values are most cost-effective to acquire. These data sets contain both continuous and categorical features; therefore, for simplicity, when estimating value distributions we converted all the continuous features to categorical features using the discretization method of Fayyad and Irani (1993). The remaining data sets are available from the UC Irvine repository (Blake and Merz 1998) and pertain to a variety of domains.

The performance of each acquisition policy is evaluated over 10 independent runs of 10-fold cross-validation as follows. For each cross-validation run, the 10-fold partition was selected at random. In each fold of the cross-validation, all policies were provided with the same subset of initial feature values, drawn uniformly at random from the training portion. All the remaining feature values in the training data constitute the initial pool of potential acquisitions. At each acquisition phase, each policy acquires the values of a set of queries from the pool of prospective acquisitions; then a new model is induced and its classification accuracy is measured on the test data. This process is repeated until a desired number of feature values has been acquired. To reduce computation costs in the experiments, we acquire queries in fixed-size batches at each iteration. For problems for which learning requires more training information, we acquired a larger number of feature values at each phase. For each data set, we selected the initial random sample size to be such that the induced model performed at least better than assigning all instances to the majority class. We later explore the policy's performance for smaller numbers of initial feature values and for different batch sizes. The test data set contains complete instances to allow us to estimate the true generalization accuracy of the constructed model. We set the exploration parameter α

⁴ From the related study by Zheng and Padmanabhan (2006).

Table 1 Error Reductions of SEU Variants Compared to a Uniform Random Acquisition Policy

Data set	β (batch size)	Initial sample (no. of instances)	Consideration set alternatives		Alternative for $A(F_{i,j} = v_k)$ SEU-accuracy	Alternatives for $P(F_{i,j} = v_k)$ SEU-frequency	SEU-DT
			SEU-US	SEU-ES			
Audiology	100	147	14.61*	19.72*	7.81*†	8.38*†	11.31*†
Car	50	1,033	10.93*	11.11*	4.25*†	8.14*†	9.42*
eToys	100	125	19.11*	49.18*	10.17*†	17.99*†	16.14*†
Expedia	100	350	10.04*	16.61*	6.38*†	5.92*†	11.03*
Lymph	20	38	7.20*	3.05*	5.56*†	6.30*	6.70*
Priceline	100	75	10.69*	1.24†	4.46*†	9.82*	9.17*†
QVC	100	225	3.76*	14.82*	−0.20†	2.83*	1.30*†
Vote	10	59	18.30*	8.33*	6.23*†	12.83*†	15.32*†
Average			11.83	15.50	5.58	9.02	10.04

*Policy is better than uniform, $p < 0.05$.† $p < 0.06$.‡SEU-US is better than the alternative policy, $p < 0.05$.

to 10. For model induction we used J48 classification-tree induction, which is the Weka (Witten and Frank 1999) implementation of C4.5 (Quinlan 1993). Integral to this induction algorithm is a missing value treatment, enabling induction from the incomplete data set. In addition, Laplace smoothing was used with J48 to improve class probability estimates.

We compare the performance of any two policies, A and B , by computing the percentage reduction in classification error rate obtained by A over B at each acquisition phase and report the average reduction over all acquisition phases. We refer to this average as the *average percentage error reduction* (Saar-Tsechansky and Provost 2004). The reduction in errors obtained with policy A over the errors of policy B is considered to be significant if the errors produced by policy A are fewer than the corresponding errors (i.e., at the same acquisition phase) produced by policy B , according to a paired t -test ($p < 0.05$) across all the acquisition phases. The learning curves that we present and the average percentage error reduction we report reflect average performance of each policy over the 10 runs of 10-fold cross-validation.

3.2. Results

3.2.1. Primary Results. Table 1 presents the average error reductions obtained by different SEU policies with respect to the uniform sampling policy, which acquires a representative set of feature values drawn uniformly at random. The number of acquisitions acquired at each acquisition phase, β , and the size of the initial sample are also presented in Table 1. In this and subsequent tables each significant value ($p < 0.05$) is marked with an asterisk (*).

Let us first examine whether SEU effectively decreases classification error compared with a uniform sampling policy. In Table 1, the fourth and fifth

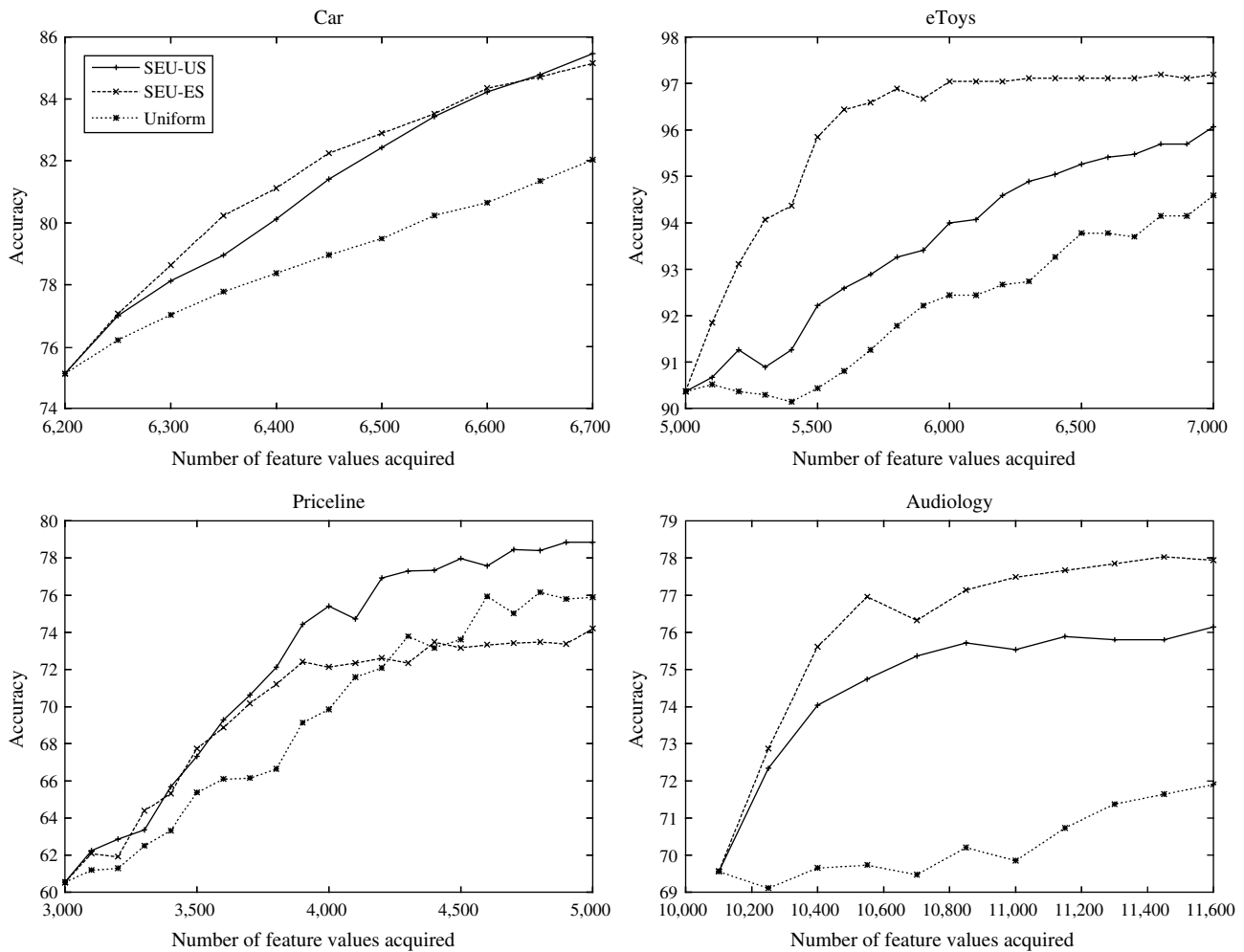
columns present the average error reductions obtained by SEU-US and SEU-ES with respect to uniform query sampling. Figure 2 presents the performance of the three policies on four data sets that exhibit the different patterns of performance we observe. For all data sets SEU builds more accurate models than uniform query sampling. The differences in performance on all data sets, except for SEU-ES on Priceline, are statistically significant.⁵ These results demonstrate that the expected utility framework and the specific methods we employ to estimate the expected improvement in performance are indeed effective for AFA: SEU selects queries that on average are more informative for induction than queries selected uniformly at random.

To underscore the advantage of using SEU, one can observe the cost benefit of using SEU to build a model exhibiting a desired performance level, compared with using a uniform acquisition policy. For example, for the eToys data set, uniform query sampling had to acquire approximately 1,800 feature values to obtain an accuracy of 94%. SEU-ES had to acquire fewer than 400 feature values to achieve the same accuracy. When data-acquisition costs are considerable, this could translate to substantial savings in the cost of building accurate models.

The results also indicate that the method employed to select the queries to be considered for acquisition can have a significant impact on the outcome. Both SEU-US and SEU-ES acquire useful feature values that significantly improve the model's performance. For some data sets, such as eToys and Audiology, SEU-ES selects significantly more informative acquisitions than SEU-US, suggesting that ES identifies a subset of

⁵Note that for SEU-ES on Priceline, the improvement at least is significant at the 0.06 level.

Figure 2 Four Characteristic Patterns of the Improvement of Classification Accuracy as a Function of the Number of Feature Values Acquired, Assuming Uniform Feature Costs



acquisitions to be considered for acquisitions superior to those drawn on average via uniform query sampling. In other data sets—e.g., Priceline—SEU-US is preferable. Recall that ES selects entire instances so all the missing feature values for these instances simultaneously become candidates for acquisition. Fleshing out a smaller set of examples may be more or less preferable in different domains. Also, if examples are incorrectly classified simply because they are outliers, ES will stumble because it will prefer all the unknown features for these examples. Moreover, if only a few features are relevant, selecting all the features will only dilute the candidate set.

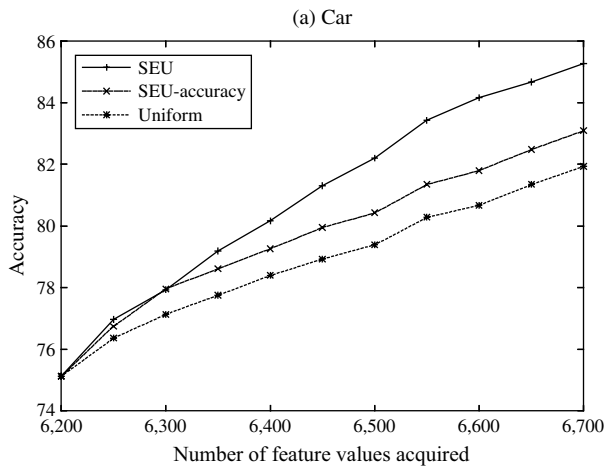
The average error reduction obtained with SEU-US over all acquisition phases ranges between 3.76% and 19.11%. SEU-ES often results in even more substantial savings, but its performance is more varied than that of SEU-US. The average error reduction obtained by SEU-ES ranges between 1.24% and 49.18%. Because it forms the consideration set based on the entire instance, ES may sometimes fail to select instances

with a highly informative feature value if the entire instance seems less informative than another instance.

In sum, both SEU policies provide considerable advantage over uniform query sampling. SEU-ES usually is the better of the two and can sometimes provide very substantial savings. SEU-US is more consistent and thus would be a more conservative choice. In the next section we expand on these results, focusing on the more conservative SEU-US policy. For the remainder of this paper, unless specified otherwise, SEU will refer to SEU-US.

We now compare SEU's performance to its performance using alternative utility measures, providing empirical support for the desired properties outlined in §2. First we examine empirically whether Log Gain is a better measure of prospective utility than is classification accuracy. The sixth column of Table 1 shows the average error reduction obtained by SEU-accuracy over uniform sampling. We mark with an asterisk (*) the data sets where SEU-accuracy is significantly better than uniform acquisition; the data

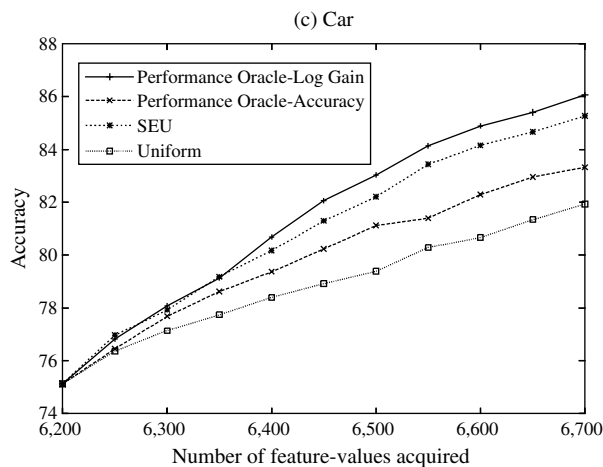
Figure 3 Classification Accuracy vs. Log Gain for Estimating the Value of an Acquisition



(b) Error reductions

Data set	Oracle-Log Gain vs. Oracle-Accuracy
Audiology	2.10*
Car	8.34*
eToys	0.70
Expedia	3.45*
Lymph	6.83*
Priceline	2.45*
QVC	2.94*
Vote	6.52*

* $p < 0.05$.



Note. The three comparisons are described in text.

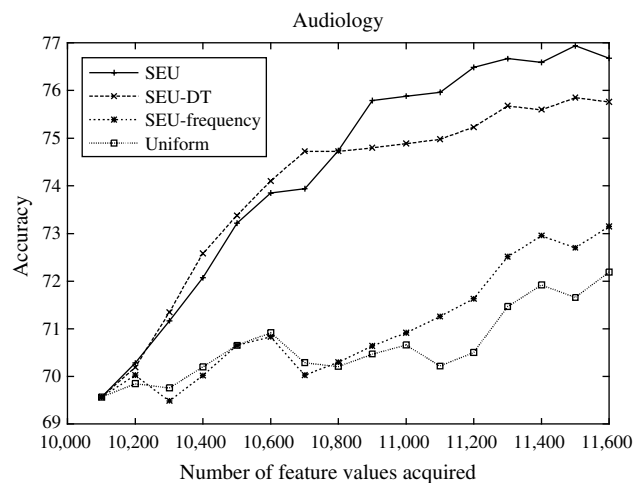
sets for which SEU is significantly better than SEU-accuracy (all of them) are marked with a double-dagger (§). Figure 3(a) shows the performance of each policy as well as of uniform sampling for the Car data set. The improvements obtained by SEU (with Log Gain) can be substantially higher than those obtained by SEU-accuracy, up to more than 12% average error reduction. These results demonstrate that by capturing

changes in the probability estimation, Log Gain indeed is able to select significantly more informative feature values to acquire, leading to better models on average.

A possible reason for the inferior performance obtained by SEU-accuracy may be the difficulty of precisely estimating classification accuracy using only the training data. In practice only the training data are available to the SEU policy, but it is useful to establish whether Log Gain is more informative even when an oracle computes the value of prospective acquisitions directly on the test data or whether classification accuracy is preferable when it is estimated with sufficient precision, in spite of its step-like form. To address this question, we compared SEU's performance with versions of the policy where Log Gain and classification accuracy are measured on the *held-out test data* rather than on the training data. We refer to these policies as *Performance Oracle-Log Gain* and *Performance Oracle-Accuracy*, respectively. Figure 3(b) presents the error reduction obtained with Performance Oracle-Log Gain compared with Performance Oracle-Accuracy. For the Car data set, Figure 3(c) presents the performance obtained by the oracles, SEU, and the uniform acquisition policy. The results confirm the advantage from detecting changes in the model's *probability* estimation via Log Gain over classification accuracy—Log Gain is able to identify more informative acquisitions even when the impact of an acquisition is evaluated with absolute precision over the test data.

Now, let us turn to the value distribution estimation employed by SEU. Table 1 presents average reductions in error when using conditional distributions (SEU, fourth column) compared with using unconditional frequency estimation (SEU-frequency, seventh column) and tree induction (SEU-DT, eighth column). For the Audiology data set, Figure 4 shows the performance of SEU with each estimation approach. These

Figure 4 Three Different Methods for Estimating the Feature-Value Distribution Compared with Uniform Acquisition



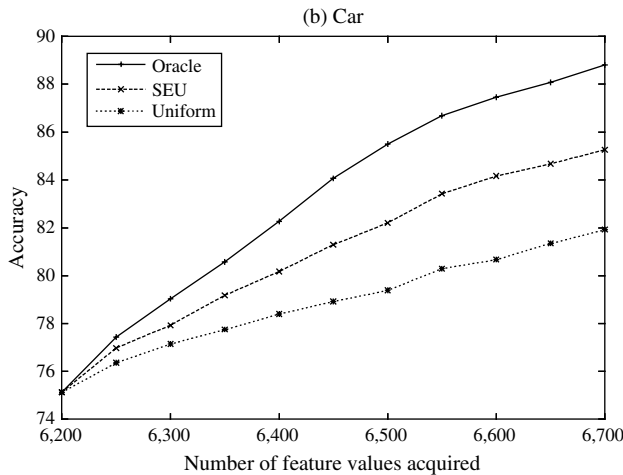
INFORMS holds copyright to this article and distributed this copy as a courtesy to the author(s). Additional information, including rights and permission policies, is available at http://journals.informs.org/.

Figure 5 Error Reduction of the Omniscient Oracle Compared with SEU

(a) Average error reductions

Data set	Oracle vs. SEU	SEU error reduction as % of oracle's
Audiology	9.47*	49.21
Car	13.11*	53.17
eToys	3.11*	40.51
Expedia	4.09*	42.89
Lymph	14.90*	36.38
Priceline	0.63	56.54
QVC	2.42*	47.31
Vote	0.66	62.43

* $p < 0.05$.



Note. SEU achieves about half the total error reduction over uniform acquisition.

results suggest that unconditional distributions provide poorer estimates of the feature-value distribution compared with the conditional distributions; the average improvement over the uniform policy over all data sets is 9.02%, compared with 11.83% obtained by SEU. For individual data sets, SEU's relative advantage with respect to SEU-frequency reached up to 6.23%. We also find that SEU is often better or comparable to SEU-DT, which can capture more complex patterns but relies on predictors that may be unknown at inference time. Errors in estimating missing predictors or their distributions contribute to prediction error cumulatively, and furthermore, because the induction technique implicitly assumes all predictors will be available during inference, it is less likely to capture alternative predictive patterns involving feature values that will be available during inference than relying exclusively on known predictors (Saar-Tsechansky and Provost 2007).

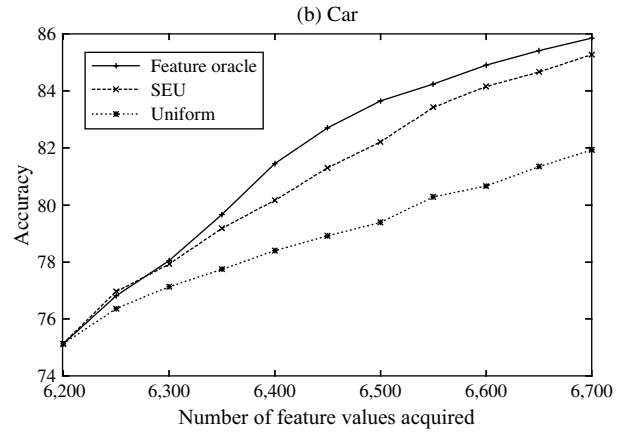
3.2.2. Oracle Policies. Let us now examine the upper-bound performance that can be obtained with SEU. Figure 5(a) shows the average error reduction obtained by the oracle, compared with SEU. Figure 5(b) shows the performance obtained by the

Figure 6 Oracle with Complete Knowledge of Feature Values (Only) Compared with SEU

(a) Average error reductions

Data set	Feature oracle vs. SEU
Audiology	-1.93
Car	3.87*
eToys	-2.60
Expedia	4.00*
Lymph	2.19*
Priceline	-2.24
QVC	-2.08
Vote	0.64

* $p < 0.05$.



oracle, SEU, and the uniform policies for the Car data set. The oracle obtains between 0.63% and 14.9% average error reduction, and its benefit is statistically significant in most cases. We also measured the error reduction obtained by each, the oracle and SEU, compared with uniform sampling to determine what proportion of the improvement obtained by the oracle is obtained by SEU. In Figure 5(a) we denote this measure as "SEU error reduction as % of oracle's." SEU consistently achieves about half the "optimal" error reduction.

We can decompose the advantages conferred by each of the oracle's perfect measures of the quantities in Equation (1) over the corresponding imperfect estimations performed by SEU. Figure 6(a) presents the average error reduction obtained with the feature oracle compared with the SEU policy. For the Car data set, Figure 6(b) shows the performance of the feature oracle, SEU, and uniform sampling. As shown, acquisitions made by the SEU policy often result in models that perform comparably to those induced with acquisitions made by the feature oracle. The feature oracle performs statistically significantly better in only three data sets; in these cases the feature oracle's acquisitions lead to models that are, on average, between 2.19% and 4% more accurate than those produced by SEU. Thus, for the purpose of choosing acquisitions

INFORMS holds copyright to this article and distributed this copy as a courtesy to the author(s). Additional information, including rights and permission policies, is available at <http://journals.informs.org/>.

