This is a draft of a paper that is to appear in the IEEE Transactions on Knowledge and Data Engineering

# Towards Intelligent Assistance for a Data Mining Process: An Ontology-based Approach for Cost-sensitive Classification[1]

**Abraham Bernstein[1,2], Foster Provost[2], and Shawndra Hill[2]**
[1]Department of Information Technology   [2]Stern School of Business
University of Zurich                       New York University
Switzerland                               U.S.A.

## Abstract

A data mining (DM) process involves multiple stages. A simple, but typical, process might include preprocessing data, applying a data-mining algorithm, and postprocessing the mining results. There are many possible choices for each stage, and only some combinations are valid. Because of the large space and non-trivial interactions, both novices and data-mining specialists need assistance in composing and selecting DM processes. Extending notions developed for statistical expert systems we present a prototype Intelligent Discovery Assistant (IDA), which provides users with (i) systematic enumerations of valid DM processes, in order that important, potentially fruitful options are not overlooked, and (ii) effective rankings of these valid processes by different criteria, to facilitate the choice of DM processes to execute. We use the prototype to show that an IDA can indeed provide useful enumerations and effective rankings in the context of simple classification processes. We discuss how an IDA could be an important tool for knowledge sharing among a team of data miners. Finally, we illustrate the claims with a comprehensive demonstration of cost-sensitive classification using a more involved process and data from the 1998 KDDCUP competition.

## Index Terms

Data mining, data-mining process, intelligent assistants, knowledge discovery

---

# 1  Introduction

Knowledge discovery from data is the result of an exploratory process involving the application of various algorithmic procedures for manipulating data, building models from data, and manipulating the models. The Knowledge Discovery (KD) process [Fayyad, Piatetsky-Shapiro & Smyth, 1996] is one of the central notions of the field of Knowledge Discovery and Data mining (KDD). The KD process deserves more attention from the research community; processes comprise multiple algorithmic components, which interact in non-trivial ways. Even data-mining specialists are not familiar with the full range of components, let alone the vast design space of possible processes. Therefore, both novices and data-mining specialists are apt to overlook useful instances of the KD process. We consider tools that will help data miners to navigate the space of KD processes systematically, and more effectively. In particular, this paper focuses on a subset of stages of the KD process—those stages for which there are multiple algorithm components that can apply; we will call this a data mining (DM) process (to distinguish it from the larger knowledge discovery process). For most of this paper, we consider a prototypical DM process template, similar to the one described by Fayyad et al. [1996] and [Chapman et al., 2000], which is shown in Figure 1. We concentrate our work here on three DM-process stages: automated preprocessing of data, application of induction algorithms, and automated postprocessing of models. We have chosen this set of steps because, individually, they are relatively well understood—and they can be applied to a wide variety of benchmark data sets.[3] In the final case study, we expand our view to a more involved DM process.
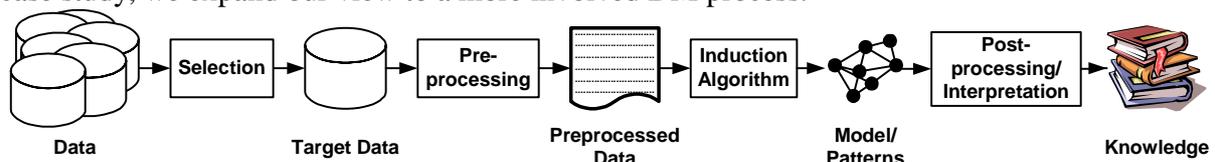


**Figure 1: The KD process (adapted from Fayad et al. [1996])**

Figure 2 shows three simple, example DM processes.[4] Process 1 comprises simply the application of a decision-tree inducer. Process 2 preprocesses the data by discretizing numeric attrib-

---

[3] More generally, because we will assemble these components automatically into complete processes that can be executed by a user, the scope of our investigation is necessarily limited to KD-process stages for which there exist automated components, and for which their requirements and functions can be specified. Important but ill-understood stages such as "business process analysis" or "management of discovered knowledge" are not included [Senator, 2000]. We also do not consider intelligent support for more open-ended, statistical/exploratory data analysis, as has been addressed by St. Amant and Cohen [1998].

[4] Descriptions of all of the techniques can be found in a data mining textbook [Witten & Frank, 2000].

utes, and then builds a naïve Bayesian classifier. Process 3 preprocesses the data first by taking a random subsample, then applies discretization, and then builds a naïve Bayesian classifier.

Consider an intelligent assistant that helps a data miner with the exploration of the space of valid DM processes. A *valid* DM process violates no fundamental constraints of its constituent techniques. For example, consider an implementation of a naïve Bayesian classifier that applies only to categorical attributes (as do many implementations). If an input data set contains numeric attributes, simply applying this classifier is not a valid DM process. However, Process 2 is valid, because it preprocesses the data with a discretization routine, transforming the numeric attributes to categorical ones. An automated system can take advantage of an explicit ontology of data-mining techniques, which defines the various techniques and their properties.
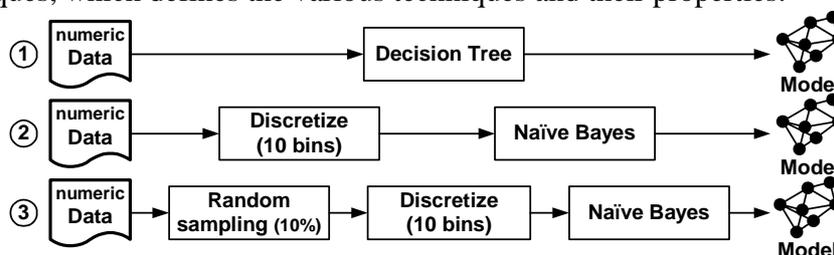


**Figure 2: Three valid DM processes**



**Figure 3: Simplified elements of a DM ontology[5]**

The Intelligent Discovery Assistant (IDA) determines characteristics of the data and of the desired mining result, and uses the ontology to search for and enumerate the DM processes that are valid for producing the desired result from the given data. Each search operator corresponds to the inclusion in the DM process of a different data-mining technique; preconditions constrain its

---

[5] The naïve Bayesian approach generally allows induction from data containing continuous attributes. The implementations used in the data mining community often have the restriction of not handling continuous data.

applicability and there are effects of applying it. Figure 3 shows some (simplified) ontology entries. Then the IDA assists the user in choosing processes to execute, for example, by ranking the process (heuristically) according to what is important to the user. The ranking shown in Figure 2 (based on the number of techniques that form the plan) would be useful if the user were interested in minimizing fuss. Another user may want to minimize run time. In that case the reverse of the ranking shown in Figure 2 would be better. Other ranking criteria are accuracy, cost sensitivity, comprehensibility, etc., and combinations thereof.

We claim that such a system can provide users with two main benefits:

1.  a systematic enumeration of valid DM processes, so they do not miss important, potentially fruitful options;

2.  effective rankings of these valid processes by different criteria, to help them choose between the options.

We also assert that an ontology-based IDA provides an infrastructure for sharing knowledge about data-mining processes, which can lead to what economists call network externalities. We do not provide experimental proof for this third hypothesis, but argue that behavioral research in the area of knowledge sharing has shown such effects in analogous applications.

We support the first claim by presenting in detail the design of an effective IDA for cost-sensitive classification, including a working prototype, describing how valid plans are enumerated based on an ontology that specifies the characteristics of the various component techniques. We show plans that the prototype produces, and argue that they would be useful not only to novices, but even to expert data miners. We provide support for the second claim with an experimental study, using ranking heuristics. Although we do not claim to give an in-depth treatment of ranking methods, we demonstrate the ability of the IDA prototype to rank potential processes by speed and by accuracy (both of which can be assessed objectively) and by combinations of the two, in the context of a classification task. Finally, we provide additional support for all the claims with an empirical demonstration, using the KDDCUP 1998 data-mining problem, showing how an IDA can take advantage of knowledge about a problem-specific DM process. For most of the paper we use simple processes, such as those presented in Figure 2, to provide sup-

port for our claims. The final demonstration goes into more depth (but less breadth) with a particular, more complex process.

## 2 Motivation and General Procedure

When engaged in design activities, people rarely explore the entire design space [Ulrich and Eppinger, 1995, p. 79]. When confronted with a new problem, data miners, even data-mining experts, often do not explore the design space of DM processes thoroughly. For example, the ACM SIGKDD Conference holds an annual competition, in which a never-before-seen data set is released to the community and teams of researchers and practitioners compete to discover knowledge from the data. KDDCUP-2000 received 30 entrants (teams) attempting to mine knowledge from electronic-commerce data. As reported by Brodley and Kohavi [2000], most types of data-mining algorithm were tried by only a small fraction of participants.

Expert data miners may ignore many data-mining approaches because they do not have access to the tools; however, readily and freely available data-mining toolkits make this reason suspect. More likely, experts do not use that many data-mining tools—especially tools that require additional pre- and post-processing or those entailing complicated installation or execution procedures (e.g., complicated parameter tweaking). Indeed, the only algorithm that was tried by more than 20% of the KDDCUP-2000 participants was decision-tree induction, which often performs reasonably well with no tweaking and with little pre-/post-processing.

The overall meta-process followed by our IDA is shown in Figure 4. The user provides data, metadata, goals, and desiderata. Then the IDA composes the set of valid DM processes, according to the constraints imposed by the user inputs, the data, and/or the ontology. This composition involves choosing induction algorithm(s), and appropriate pre- and post-processing modules (as well as other aspects of the process, not considered in this paper). Next, the IDA ranks the suitable processes into a suggested order based on the user's desiderata. The user can select plans from the suggestion list, hopefully aided by the ranking. Finally, the IDA will produce code for and can execute (automatically) the suggested processes on the selected data.
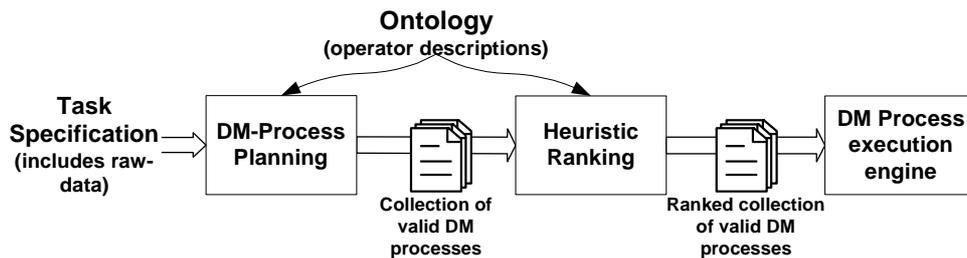
**Figure 4: The overall process followed by an IDA**

# 3 Enumerating Valid Data Mining Processes

Our first claim is that an ontology-based IDA can enumerate DM processes useful to a data miner. We support our claim in two ways. First, we describe how the ontology can enable the composition of only valid plans. Second, we describe process instances produced by our prototype (called IDEA), in order to provide evidence that they can be non-trivial. Later we will describe how problem-specific elements can be incorporated into an IDA; for clarity and generality first we concentrate on domain-independent elements of the DM process. For example, when presented with a data set to mine, a knowledge-discovery worker (researcher or practitioner) generally is faced with a confusing array of choices [Witten & Frank, 2000]: should I use C4.5 or naive Bayes or a neural network? Should I use discretization? If so, what method? Should I subsample? Should I prune? How do I take into account costs of misclassification?

## 3.1 An Ontology-based Intelligent Discovery Assistant

Consider a straightforward example: a user presents a large data set, including both numeric and categorical data, and specifies *classification* as the learning task, along with the appropriate dependent variable. The IDA asks the user to specify his/her desired tradeoffs between accuracy and speed of learning. Then the IDA determines which DM processes are appropriate. For our example task, decision-tree learning alone might be appropriate. Or, a decision-tree program plus subsampling as a pre-process, or plus pruning as a post-process, or plus both. Are naive Bayes or neural networks appropriate for this example? Perhaps not by themselves. Not if naive Bayes only takes categorical attributes. Typically, neural networks only take numeric attributes. However, appropriate pre-processing (transforming the data type) may enable their use.

The IDA uses the ontology to assist the user in composing valid and useful DM processes. Basing our design on AI planning [Ghallab et al., 1998] and semantic web services [Ankolekar et al., 2001], the prototype's ontology contains for each operator:

- A specification of the conditions under which the operator can be applied, including a pre-condition on the state of the DM process, its compatibility with preceding operators, and the inputs necessary for the execution of the algorithm.

- A specification of the operator's effects on the DM process's state and on the data.

- Estimations of the operator's effects on attributes such as speed, accuracy, model comprehensibility, etc. (shown as heuristic indicators in Figure 3).

- A help function to obtain comprehensible information about each of the operators.

In addition the ontology contains schemata for generic problems such as target marketing. The schemata are represented internally as complex, decomposable operators with the same parameters as their simple counterparts. The only difference is that some of the steps within the complex operators might not be completely specified, opening a design space of sub-solutions (section 5 provides an example). The collection of all schemata is a case-base of proven useful processes.



**Figure 5: The data-mining ontology**
(partial view, the italicized leaf nodes were used in the ranking experiments)

Figure 5 shows a structural view of the prototype ontology, which at the highest level groups the DM operators into: pre-processing, induction, post-processing. Each group is further sub-divided. At the leaves of this tree are the actual operators (selected examples are shown in italics). We constructed this prototype ontology by first considering the types of operators provided

by typical data-mining toolkits. We chose the three stages of the data-mining process that have received the most automation—resulting in a non-trivial selection of operators. We then divided these into subcategories, focusing on operators that would be useful for cost-sensitive classification. Finally we chose implementations that were available in the Weka toolkit [Witten and Frank, 2000]. This ontology clearly is not complete, even for cost-sensitive learning; however, it is sufficient for the experimental demonstrations that follow. After the ontological structure was in place, we called on our own expertise to provide the operator specifics (preconditions, postconditions, etc., discussed above). Obviously, the ontology is limited by our knowledge. All of this was done before the experiments below were conducted, except as mentioned below.

Based on the prototype ontology, we built a prototype IDA, which we call IDEA (**I**ntelligent **D**iscovery **E**lectronic **A**ssistant). Following our general framework for IDAs (see Figure 4), IDEA first gathers a *task specification* for the DM process, analyzes the data that the user wishes to mine and extracts the relevant meta-data, such as the types of attributes included (e.g., continuous, categorical). Using a GUI, the user then can complement the gathered information with additional knowledge about the data, and can specify the type of information/model he/she wishes to mine. IDEA's first core component, the *DM-process planner* is described in Section 3.2. A collection of valid DM processes typically will contain processes that are undesirable for certain user goals—e.g., sacrificing too much accuracy to obtain a model fast. IDEA's second core component, the *heuristic ranker*, ranks the valid DM processes using a combination of several heuristic functions. The GUI also allows the user: to specify tradeoffs (weights) between ranking functions, to sort the list of plans using any (weighted combination) of the rankings, to examine the details of any process plan, and to generate code for and to run the processes.

A final function of IDEA is to supply the users with an interface to the ontology. It allows users to browse the ontology entries with a tree-like hierarchy browser. To add new operators to the ontology requires adding a new element in the ontology tree and specifying its parameters. When adding the Weka ID3 tree learner operator, for example, the user would first choose Decision Tree as an appropriate parent in the ontology. Then the user would proceed to add the appropriate parameters. As a child of Decision tree, the new ID3 operator would inherit some parameter values. The user must complement these parameters with the actual algorithm implementation and

the call interface (for Weka's ID3: java-class weka.classifiers.trees.Id3 with the relevant method calls or a reference to a WSDL[6] file), and the heuristic parameters (such as `Speed = + 25`).[7]

## 3.2 Enumerating Valid DM Processes: IDEA's procedure

Our first claim is that an IDA can produce a systematic enumeration of DM processes that will be useful to data miners, to help them keep from overlooking important processes. To enumerate (only) valid DM processes, IDEA performs a straightforward search of the space of processes defined by the ontology, constrained by the restrictions on operator application defined in the ontology. The structure of the search problem is amenable to more complex, AI-style planning, but so far the straightforward search has been sufficient. IDEA constructs specification of the sequence of DM operators (i.e., the DM process) that moves from the start state—the meta-data description of the data set—to the goal state—typically a prediction model with some desired properties. Starting with an empty process at the start state, at every state it finds the applicable operators using the compatibilities, adds each operator (separately) to the partial process that brought it to the current state, and transforms the state using the operator's effects. Using the example above, in order to apply naïve Bayes, the current state must not contain numeric attributes; which would be the case after discretization. The planner would not apply discretization twice, because after the first application, the state no longer would contain numeric attributes, and thus the preconditions of discretization no longer would apply. The planner stops pursuing a given process when it has reached either the goal state or some "dead-end" state that will not lead to the goal state. The planner also can add complex operator schemata to any solution. Akin to hierarchical planning, it then must revisit all the non-specified steps and treat them as planning problems themselves (see section 5 for an example). The central difference from traditional, AI planning is that execution does not stop when a first viable solution is found. Instead the search

---

[6] WSDL (Web Services Description Language) is an XML format for describing services (or remote procedures), providing all information necessary to locate and call the service [Christensen et al., 2001].

[7] One method of ranking by predicted speed follows a simple "composition" approach that our results below show is effective. Each DM operator in the ontology contains a heuristic entry of the form <operator> <number>, where the operator is either +,- ,*,/ and number is an integer. Starting with a total score of 0 the algorithm visits each DM operator in the plan in the order of its appearance then applies the <operator> and <number> to compute the new total score.

returns as many valid processes as possible,[8] as users of KD tools are often not able to express their preferences precisely or completely before seeing possible available alternatives.

The constraints in the ontology are essential. If we use the ontology whose overall structure is shown in Figure 5, give the goal of *classification*, and constrain the search only with the ordering of the logical groupings imposed by the prototype ontology (i.e., pre-processing precedes induction which precedes post-processing), IDEA generates 163,840 DM processes. Adding the constraints imposed by the pre- and post-conditions of the operators,[9] IDEA produces 597 valid process instances—less than one-half of one percent of the size of the unconstrained enumeration. Adding metadata (e.g., the data set contains numeric attributes) or user desiderata (e.g., the user wants cost-sensitive classification) allows the enumeration to be constrained even further.

### 3.3 Enumerating Valid DM Processes: example enumerations from IDEA

The enumerations of processes produced by IDEA are not trivial. In many cases they would be valuable not only to novice data miners, but even to experts.

**Example 1)** When IDEA is given the goal of producing a *cost-sensitive classifier for a two-class problem*, it produces an enumeration comprising 189 DM processes. The enumeration includes building a class-probability estimator and setting a cost-specific threshold on the output probability. It includes building a regression model and determining (empirically) an effective threshold on the output score. The enumeration also includes using class-stratified sampling with any classification algorithm (which transforms an error-minimizing classifier into a cost-minimizing classifier). Novice data miners certainly do not consider all these options when approaching a cost-sensitive problem. In fact, we are aware of no single published research paper on cost-sensitive learning that considers one of each of these types of option [Turney, 1996].

**Example 2)** When we give IDEA the goal of producing *comprehensible classifiers*, the top-ranked DM process is: `subsample the instances` → `feature selection` → `use a rule learner` → `prune the resultant rule set`. Although comprehensibility is a goal of much machine-learning research, we are not aware of this process being used or suggested. This process is interesting because each component individually has been

---

[8] As long as the number of DM operators that will be available to an IDA is not huge, the speed of planning is unlikely to be problematic. For example, with the prototype DM ontology (currently incorporating a few dozen operators), the current DM-process planner can generate all valid processes (up to several hundred for problems with few constraints) in less than a second.
[9] For example: neural networks require numeric attributes; decision-tree pruning can only apply to decision trees, etc.

shown to yield more comprehensible models; why shouldn't the composition yield even more comprehensible models? As another DM process highly ranked by comprehensibility, IDEA suggests: `build a decision tree` ➔ `convert tree to rules` ➔ `prune rule set`. This also is a non-trivial suggestion: it is the process introduced by Quinlan [1987] and shown to produce a combination of comprehensibility and high accuracy. Although the addition to the ontology of `convert tree to rules` was influenced by Quinlan's work, we did not "program" the system to produce this process. IDEA composed and ranked processes based on knowledge of individual operators. This is particularly valuable, because the addition of a new operator to the ontology can have far-reaching effects (e.g., adding "convert trees to rules" results in this plan being suggested strongly for comprehensible classification).

**Example 3)** Consider the case where the user is interested in classification, but wants to get results fast. Does IDEA's enumeration contain particularly useful (fast) processes? Indeed it suggests processes that use fast induction algorithms, such as C4.5 (shown to be very fast for memory-resident data, as compared to a wide variety of other induction algorithms [Lim et al., 2000]). It also produces suggestions not commonly considered [Provost & Kolluri, 1999]. For example, the enumeration contains plans that use discretization as a preprocess. Research has shown that discretization as a preprocess can produce classifiers with comparable accuracy to induction without the preprocess [Kohavi & Sahami, 1996]; but with discretization, many induction algorithms run much faster. For example, as described by Provost and Kolluri, most decision tree inducers repeatedly sort numeric attributes, increasing the computational complexity considerably; discretization eliminates the sorting. IDEA's suggestions of fast plans also include plans that use subsampling as a preprocess. Most researchers studying scaling up have not considered subsampling explicitly, but of course it produces classifiers much faster—and for large data sets it has been shown to produce classifiers with comparable accuracies [Oates & Jensen, 1997].

## 4  IDAs can produce effective rankings

Large enumerations of DM processes can be unwieldy. It is important to help the user choose from among the candidate processes. Rankings of DM processes can be produced in a variety of ways. For example, static rankings of processes for different criteria could be stored in the system. Flexible rankings also are important—so that as new ontological knowledge is added, the

system can take advantage of it immediately. IDEA produces rankings dynamically by composing the effects of individual operators. The ontology contains estimations of the effects of each operator on each goal. For example, an induction algorithm may be estimated to have a particular speed (relative to the other algorithms). Taking a 10% random sample of the data as a pre-process might be specified to reduce the run time by a factor of 10 (which would be appropriate for algorithms with linear run times). Correspondingly, sampling might be specified to reduce the accuracy by a certain factor (on average), and to increase the comprehensibility by a different factor (cf., the study by Oates and Jensen [1997]). For a given DM process plan, an overall score is produced as the composition of the functions of the component operators.

## 4.1  Details of ranking experiments

In order to provide a demonstration to support our claim that IDAs can produce useful rankings, we coupled IDEA with a code generator that generates code for the Weka data-mining toolkit[10] [Witten and Frank, 2000]. The system generates Java code for executing the plans, as well as code for evaluating the resulting models based on accuracy and speed of learning. We assess IDEA's ability to rank processes by speed and by accuracy, because these are criteria of general interest to users and for which there are well-accepted evaluation metrics. Furthermore, one expects a rough tradeoff between speed and accuracy [Lim, *et al.*, 2000], and a user of an IDA may be interested in points between the extremes.

For the experiments in this section, we restricted the ontology to a subset for which it is feasible to study an entire enumeration of plans thoroughly. The ontology subset uses seven common pre-processing, post-processing, and induction techniques (for which there were appropriate functions in Weka, see below). The experimental task is to build a *classifier*, and has as its start state a data set *containing at least one numeric attribute* (which renders some inducers inapplicable without preprocessing). Table 1 shows on the left the list of 16 valid process plans IDEA created for this problem; on the right is a legend describing the 7 operators used.[11] Even this small ontology produces an interesting variety of DM-process plans. For example, the ontology

---

[10] The choice of Weka was driven by the availability of a large number of suitable machine learning operators. Weka does have the drawback that it mostly operates on in-memory structures making it unsuitable for exploration of some realistic large-scale data sets. In particular, the pre-processing steps, which oftentimes entail accessing large databases, should be handled a suitable database environment or within a full-scale data mining pre-processing environment like Mining Mart [Morik and Scholz, 2003].
[11] The last operator in Table 1, cpe, which places an appropriate threshold on a class-probability estimator, becomes a no-op for Naïve Bayes (nb) in the Weka implementation, because Weka's implementation of nb thresholds automatically.

specifies that naïve Bayes only considers categorical attributes, so the planner needs[12] to include a preprocessor that transforms the data. Indeed, although the ontology for the experiments is very small, the diversity of plans is greater than in many research papers.

In Table 1, the first column ranks the plans by the number of operators in the plan. This may be interesting to users who will be executing plans manually, who may be interested in minimizing fuss. We will not consider this ranking further except to reference plans by number. The *heuristic rank* columns of Table 1 show a pairs of speed rankings computed by heuristics. The "credit-g" ranking is a static ranking created by running all the plans on one, randomly selected data set (viz., credit-g[13]). A static ranking makes practical sense if the flexibility to add new operators is not of primary importance. Adding new operators (or otherwise changing the ontology) changes the space of plans, in which case a static ranking would have to be updated or recomputed. The "composition" ranking was generated by a functional composition; the ontology specifies a base accuracy and speed for each learner, and specifies that all the preprocessing operators will reduce accuracy and will increase speed, by different amounts (see Footnote 7). The heuristic functions are subjective, based on our experience with the different data-mining techniques and on our reading of the literature (e.g., [Lim *et al.*, 2000]). The ranking functions were fixed before we began using Weka's particular implementations, with one exception: because speed ratings differ markedly by implementation, we ran Weka on one data set (credit-g) to instantiate the base speed for the learning algorithms and the improvement factors for sampling and for discretization.

---

[12] This is not strictly true for the Weka implementation, for which naïve Bayes is augmented with a density estimator for processing numeric attributes. The Weka implementation would be considered naïve Bayes plus a particular numeric preprocessor.
[13] We did not use credit-g as a testing data set in our experiments.

| | steps | heuristic rank | |
|---|---|---|---|
| | | credit-g speed | composition speed |
| Plan # 1 | c4.5 | 13 | 13 |
| Plan # 2 | part | 16 | 16 |
| Plan # 3 | rs, c4.5 | 2 | 5 |
| Plan # 4 | rs, part | 8.5 | 10 |
| Plan # 5 | fbd, c4.5 | 12 | 11 |
| Plan # 6 | fbd, part | 15 | 14 |
| Plan # 7 | cbd, c4.5 | 11 | 12 |
| Plan # 8 | cbd, part | 14 | 15 |
| Plan # 9 | rs, fbd, c4.5 | 4 | 3 |
| Plan # 10 | rs, fbd, part | 6.5 | 8 |
| Plan # 11 | rs, cbd, c4.5 | 5 | 4 |
| Plan # 12 | rs, cbd, part | 6.5 | 9 |
| Plan # 13 | fbd, nb, cpe | 8.5 | 6 |
| Plan # 14 | cbd, nb, cpe | 10 | 7 |
| Plan # 15 | rs, fbd, nb, cpe | 1 | 1 |
| Plan # 16 | rs, cbd, nb, cpe | 3 | 2 |

**Table 1: 16 process plans and rankings**

| Legend for operators used in plans | |
|---|---|
| acronym | name/algorithm |
| rs | Random sampling (result instances = 10% of input inst.) |
| fbd | Fixed-bin discretization (10 bins) |
| cbd | Class-based discretization (Fayyad & Irani's [1993] MDL |
| c4.5 | C4.5 (using Witten & Frank's [2000] J48 implementation) |
| part | Rule Learner (PART, Frank & Witten [1998]) |
| nb | Naïve Byes (John & Langley [1995]) |
| cpe | CPE-thresholding post-processor |

| Dataset name | Size |
|---|---|
| heart-h | 294 |
| heart-c | 303 |
| ionosphere | 351 |
| balance-scale | 625 |
| credit-a | 690 |
| diabetes | 768 |
| vehicle | 846 |
| anneal | 898 |
| vowel | 990 |
| credit-g | 1000 |
| segment | 2310 |
| move | 3029 |
| dna | 3186 |
| gene | 3190 |
| adult10 | 3256 |
| hypothyroid | 3772 |
| sick | 3772 |
| waveform-5000 | 5000 |
| page | 5473 |
| optdigits | 5620 |
| insurance | 9822 |
| letter | 20000 |
| adult | 32561 |

**Table 2: Data set names and sizes**

Our experiments—to assess the feasibility of using an IDA to provide rankings by speed and by accuracy—compare the proposed rankings to rankings generated by actually running the plans on the data sets. For the experiments, we used 23 data sets from the UCI Repository [Blake & Merz, 2001], each containing at least one numeric attribute. The data sets and their total sizes are listed in Table 2. Unless otherwise specified, for each experiment we partitioned each data set randomly into halves (we will refer to these subsets as $D_1$ and $D_2$). We used ten-fold cross-validation within $D_2$ to compute average classification accuracy and average speed—which then are used to assess the quality of the ex-ante rankings, and to construct the "actual" (ex-post) rankings for all comparisons. (We will use the $D_1$s, later, to construct auto-experimentation rankings; the $\{D_1, D_2\}$ partitioning ensures that all results are comparable.)

| Plan Name | credit-g ranking | composition ranking | D2 ("actual") ranking |
|---|---|---|---|
| Plan # 2 | 16 | 16 | 16 |
| Plan # 6 | 15 | 14 | 15 |
| Plan # 8 | 14 | 15 | 14 |
| Plan # 1 | 13 | 13 | 13 |
| Plan # 7 | 11 | 12 | 12 |
| Plan # 4 | 9 | 10 | 11 |
| Plan # 5 | 12 | 11 | 10 |
| Plan # 14 | 10 | 7 | 9 |
| Plan # 10 | 7 | 8 | 8 |
| Plan # 12 | 7 | 9 | 7 |
| Plan # 3 | 2 | 5 | 6 |
| Plan # 13 | 9 | 6 | 5 |
| Plan # 11 | 5 | 4 | 4 |
| Plan # 9 | 4 | 3 | 3 |
| Plan # 16 | 3 | 2 | 2 |
| Plan # 15 | 1 | 1 | 1 |

**Table 3: Adult data set rankings by speed**

| | speed | | accuracy |
|---|---|---|---|
| | credit-g ranking | composition ranking | auto-experiment. |
| heart-h | 0.39 | 0.30 | -0.06 |
| heart-c | **0.62** | **0.59** | 0.06 |
| ionosphere | **0.80** | **0.70** | 0.20 |
| balance-scale | **0.82** | **0.81** | **0.55** |
| credit-a | **0.94** | **0.91** | **0.71** |
| diabetes | **0.55** | **0.64** | 0.49 |
| vehicle | **0.94** | **0.95** | **0.91** |
| anneal | **0.98** | **0.92** | **0.90** |
| vowel | **0.90** | **0.93** | **0.90** |
| segment | **0.89** | **0.92** | **0.92** |
| move | **0.90** | **0.95** | **0.87** |
| dna | **0.98** | **0.94** | **0.91** |
| gene | **0.92** | **0.95** | **0.88** |
| adult10 | **0.97** | **0.97** | **0.86** |
| hypothyroid | **0.95** | **0.91** | **0.96** |
| sick | **0.95** | **0.89** | 0.18 |
| waveform-5000 | **0.90** | **0.94** | **0.94** |
| page | **0.86** | **0.85** | **0.74** |
| optdigits | **0.89** | **0.87** | **0.84** |
| insurance | **0.95** | **0.93** | **0.84** |
| letter | **0.90** | **0.96** | **0.96** |
| adult | **0.93** | **0.98** | **0.86** |
| mean | **0.86** | **0.85** | **0.70** |
| median | **0.90** | **0.92** | **0.86** |

**Table 4: Spearman ranks for ranking heuristics for speed and accuracy**

## 4.2 Ranking by Speed

Our first experiments examine whether the heuristics can be effective for ranking DM processes by speed. Since being able to rank well by speed is most important for larger data sets, consider the largest of our data sets: adult. Table 3 shows the two rankings and the actual (ex-post) ranking based on the average run times for all the plans. The table is sorted by the actual ranking, and the table entries are the positions of each plan in each ranking (i.e., 1 is the first plan in a ranking, 2 the next, etc.). Both heuristics rank very well. For the credit-g ranking (on the adult data set) Spearman's rank-correlation $r_s = 0.93$ and for the composition ranking, $r_s = 0.98$ (recall that perfect rank correlation is 1, no correlation is 0, and a perfectly inverted ranking is -1).

Table 4 shows for all the domains the correlations between the rankings produced by the heuristics and the rankings based on the actual speeds. Here the data sets are presented in order of increasing size (large ones toward the bottom). Highlighted in bold are the cases where $r_s > 0.5$ (all but the smallest data set).[14]  Neither heuristic is superior, but both are effective; for both ranking heuristics, the average is approximately $r_s = 0.85$.

---

[14] The choice of 0.5 was ad hoc, but was chosen before running the experiment. Examining hand-crafted rankings with various $r_s$ values seemed to indicate that 0.5 gave rankings that looked good.

### 4.3 Ranking by Accuracy – Using Auto-Experimentation

Our next set of experiments examines whether the IDA can be effective for ranking DM processes by accuracy. Note that one would not expect to be able to do nearly as well at this task as for ranking by speed. Nevertheless, it would be helpful to be able to give users guidance in this regard, especially when a system proposes a process containing a component with which the user is not familiar—if the process were ranked highly by accuracy, it would justify learning about this new component. However, our attempt to use heuristic scores, similar to those that were successful for ranking by speed, did not produce particularly good accuracy rankings. Fortunately, an IDA can perform *auto-experimentation*, composing process plans and running its own experiments to produce a ranking of the plans by accuracy.[15] Although this may be the best possible ranking method (albeit time consuming), careful experimental evaluations of the accuracies of predictive models produce only estimations of the accuracy of the models on unseen data. The quality of the rankings of DM processes produced by such estimation will vary (e.g., by data-set size), and for any particular domain must be determined empirically.

We now present an experiment to assess the effectiveness of such a procedure. For each domain, IDEA composed DM process plans and generated Weka code for the plans and for their evaluations via cross-validation. For each data set, the cross-validation was performed on data subset $D_1$ to produce an estimation of the accuracy that would result from running the plan on a data set from the domain. These accuracies were used to construct a ranking of the DM-process plans by accuracy for each data set. These rankings then were compared to the ranking produced on data set $D_2$ (identically to all previous experiments). Table 4 lists the resulting rank correlations in the rightmost column. As expected, the empirically determined rankings are considerably better for the larger data sets: averaged $r_s = 0.86$ for the data sets with $>= 5000$ records.

### 4.4 Trading off Speed and Accuracy

For large data sets auto-experimentation provides good accuracy rankings, but one pays a considerable run-time price as the data-set size grows. What if a user is willing to trade off some speed for a better accuracy ranking, but does not have the time for full-blown auto-experimentation (i.e., running all the plans on all the data)? An alternative is to perform auto-experimentation on

---

[15] This is not an option for speed rankings, because the auto-experimentation process itself may be (very) time consuming.

subsamples of the data to estimate the accuracy ranking for the full data set. We now demonstrate that our IDA can allow users to trade off quality of ranking for timeliness.

IDEA ran the process plans for the six largest data sets (each having 5000 or more total records) on increasingly larger subsets of the data. Specifically, for each domain's $D_1$, we selected random subsets of 10%, 20%, …, 100% of the data. For each subset, IDEA performed cross-validation to determine empirically an expected accuracy ranking, identically to the previous experiment. For this experiment, we consider only the eight DM-process plans that do not (already) contain random sampling. Figure 6 plots the rank correlations as the size of the sample grows, and in bold shows the average rank correlation as size grows. As expected, the largest samples give better rankings than the smallest ones. For the 100% sample, all are above 0.5, and all but optdigits are above 0.8. On the other hand, for several of the data sets (page, adult, letter) the rankings with the 10% sample are not much better than random.
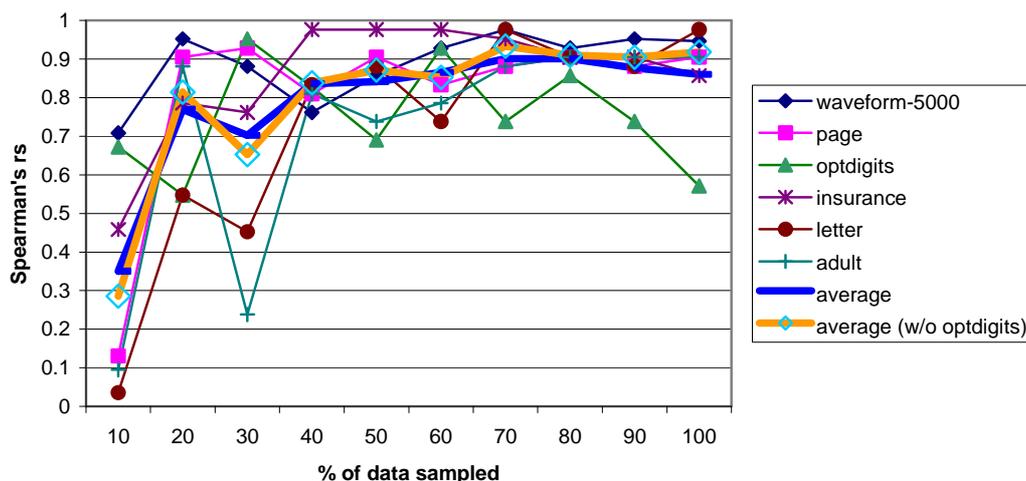


**Figure 6: Rank correlations and sample size**

With one notable exception, the rank correlations become relatively stable when about half of the data have been seen. The optdigits curve is unusual: the rank correlations do not increase and do not become more stable as more data are used. Further investigation shows that optdigits is, in an important sense, too easy. Specifically, all methods perform extremely well, even with small training sets, so it is not possible to rank them meaningfully beyond a certain level. The last series in Figure 6 (bold and marked with a ◊) graphs the average without the optdigits data, show-

ing that the average performance is as desired (generally increasing, but with decreasing marginal benefits).

In sum, the results in this section demonstrate that it is possible for an IDA to produce effective rankings of generated processes by different desiderata (speed and accuracy), and to produce rankings that make tradeoffs between the two.

## 5  Demonstration with a more complex DM process

We now present the results of a final set of experiments, to demonstrate further the power of IDAs. The prototypical DM-process template that we used for the discussions and experiments above was straightforward—as was necessary to introduce our notion of an IDA and to run experiments on a suite of benchmark data sets. However, in real-world situations the DM process can be more complex [Agrawal, 1998]. We assert that in such cases the potential value of an IDA is even greater, because there is greater need for expertise in technique and process.

The data we use for our demonstration were the subject of the 1998 KDDCUP. The rationale for choosing the KDDCUP 1998 dataset was threefold. First, the data set highlights the strengths of the planning-and-ranking approach: the combination of human insight about the problem and machine support for the systematic exploration of the design space. Second, it allows us to show the applicability of IDEA in the context of a more complex, cost-sensitive learning problem, rather than the straightforward classification problem used for the previous demonstrations. Finally, the data set has already been preprocessed extensively, making it suitable for our prototype, which concentrates on the building of the classification model, not on feature construction and selection. Even with the extensive preprocessing, the KDDCUP 1998 problem is not trivial.

The KDDCUP 1998 problem was to select a subset of customers to whom to mail solicitations, in order to maximize profit (revenues minus the cost of mailing). Participants built models from the training data, using a wide variety of different methods. To determine the winners, the organizers evaluated (on a separate test set for which the true answers were hidden) how much profit each team's model would have garnered. More specifically, KDDCUP 1998 was based on data from a fund-raising campaign undertaken by a national veterans association. The customer base was a set of individuals who donated in prior campaigns, and the goal was to select those from whom to solicit donations in the current campaign. Each observation in the data set is an indi-

vidual, and includes (for example) the response to the prior campaign. The training set from the competition consists of 95412 records and the test set consists of 96367 records. The mailing cost is $0.68 and the average donation is $15.60 with a range of $1-$200. The donation frequency is about 5% of the population. Using the default strategy of mailing to everyone, the average profit over the test set is $10,560. The (actual) results of KDDCUP 1998 are presented in Table 5. For our experiment, we use the variables used by Zadrozny and Elkan [2001].[16]

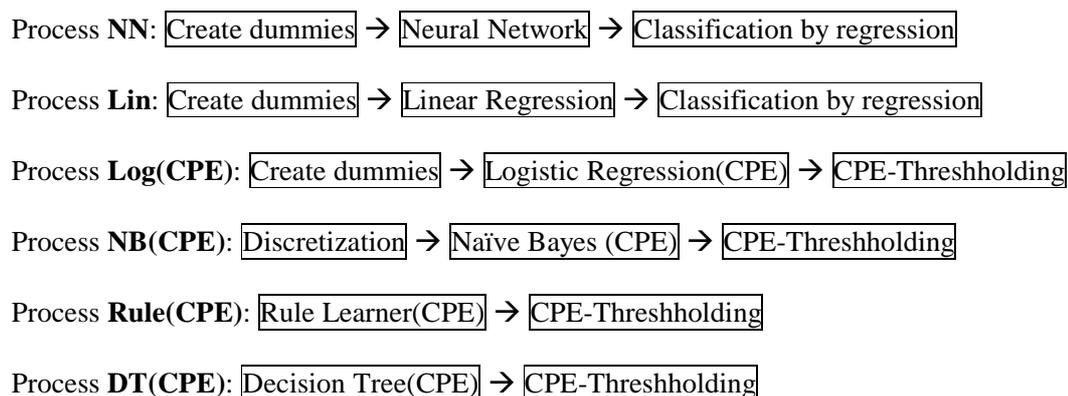| Participants | Profit | %Gain |
|---|---|---|
| Urban Science | $14,712 | 39.32 |
| SAS | $14,662 | 38.84 |
| #3 | $13,954 | 32.14 |
| #4 | $13,825 | 30.92 |
| #5 | $13,794 | 30.63 |
| #6 | $13,598 | 28.77 |
| #7 | $13,040 | 23.48 |
| #8 | $12,298 | 16.46 |
| #9 | $11,423 | 8.17 |
| #10 | $11,276 | 6.78 |
| #11 | $10,720 | 1.52 |
| #12 | $10,706 | 1.38 |
| #13 | $10,112 | -4.24 |
| #14 | $10,049 | -4.84 |
| #15 | $9,741 | -7.76 |
| #16 | $9,464 | -10.38 |
| #17 | $5,683 | -46.18 |
| #18 | $5,484 | -48.07 |
| #19 | $1,925 | -81.77 |
| #20 | $1,706 | -83.84 |
| #21 | ($54) | -100.51 |

**Table 5: Results of 1998 KDDCUP**

This was a challenging competition: the spread between the different competitors is quite large. Notice that 9 of 21 entries produced lower profits than did the default strategy of mailing to everyone. In fact, the last-place entry actually lost money. The winners achieved a 39% increase in profit over the default strategy. Notice also that the winners are experts in this sort of data mining: Urban Science specializes in building models for target marketing (and in fact, they also won the 1997 KDDCUP). In second place is SAS, who also have extensive experience with this sort of modeling. The competitors with the lower scores most likely applied data mining tools in the manner typical of data-mining/machine-learning research. As we will demonstrate, the *straightforward* application of existing tools is insufficient for high-level performance on these data. However, the inclusion of application-specific, DM-process-related knowledge is. As we will

---

[16] Note that selection and construction of features also is part of the KD process. We do not treat them in this paper, except in Limitations, below.

show, it is essential for IDEA to incorporate application-specific process knowledge. First let us consider how IDEA performs without doing so.

We followed a methodology intended to mimic the algorithmic portion of the process that KDDCUP competitors would have taken (i.e., not including feature construction and selection). Specifically, we create rankings of DM processes considering only the training set (estimating the profit that would be obtained). To assess the quality of a ranking, we calculate the "actual" profits on the test set. The 1998 KDDCUP focused on a problem of cost-sensitive classification: classify into one of two categories, solicit or do not, taking into account the cost of false positives (the mailing costs) and the cost of false negatives (the lost revenue). We use a larger set of induction algorithms than in the experiments above, but for clarity, for this experiment we do not consider pre- and post-processing explicitly.

Process **NN**: Create dummies → Neural Network → Classification by regression

Process **Lin**: Create dummies → Linear Regression → Classification by regression

Process **Log(CPE)**: Create dummies → Logistic Regression(CPE) → CPE-Threshholding

Process **NB(CPE)**: Discretization → Naïve Bayes (CPE) → CPE-Threshholding

Process **Rule(CPE)**: Rule Learner(CPE) → CPE-Threshholding

Process **DT(CPE)**: Decision Tree(CPE) → CPE-Threshholding

**Figure 7: DM processes generated for cost-sensitive classification**

Figure 7 shows 6 DM processes generated for cost-sensitive classification. As mentioned above, a wider variety of learning algorithms (from Weka) is used here, and only one process with each algorithm is generated. Specifically, the first two processes produce regression models: process "NN" is the application of a neural network learner and process "Lin" is the application of linear regression. As mentioned in section 3.3, regression models can be converted to cost-sensitive classification models by a postprocessor that chooses (by experimenting with the training data) an appropriate threshold on the predicted (output) value ("classification by regression"). Both of these algorithms require categorical variables to be preprocessed into a set of binary "dummy" variables. The last four processes use algorithms that create "class probability estimators," which give an estimation of the probability that a new example belongs to the class

in question (here, "will donate"). Such a model can be converted to a cost-sensitive classifier with a postprocessor that chooses a threshold decision-theoretically, taking into account the misclassification costs. Process Log(CPE) uses logistic regression, which also requires preprocessing of categorical variables into dummies. Process NB(CPE) uses naïve Bayes, for which discretization is used as a preprocess. Processes Rule(CPE) and DT(CPE) build rule-based and decision-tree models, respectively; these do not require the preprocessing of numeric or categorical variables.

Table 6 shows the ranking of these processes by estimated profit, the actual profit calculated on the test set, and the resulting percentage gain over the default strategy of mailing to everyone. The profit was estimated by auto-experimentation (using cross-validation, as above) on the training data. Note that except for the neural network classifier, the ranking by estimated profit is perfect. Unfortunately, even without the error, the procedure would have placed only 9[th] (of 21) in the competition. What's worse, only one of the processes actually beats the default strategy of mailing to everyone. To be fair, this was a very difficult problem for data miners not intimate with modeling for problems such as target marketing. Indeed, the participants in the contest were serious data-mining researchers and tool vendors, and only half were able to do significantly better than the default strategy.

| Plan | Rank | Profit | %Gain |
|------|------|--------|-------|
| NN | 1 | $6,919 | -34.48 |
| Lin | 2 | $11,968 | 13.33 |
| Log(CPE) | 3 | $10,520 | -0.37 |
| Rule(CPE) | 4 | $9,924 | -6.02 |
| NB(CPE) | 5 | $9,538 | -9.68 |
| DT(CPE) | 6 | $8,496 | -19.54 |

| Legend for Operators Used in Plans | |
|------|------|
| acronym | name/algorithm |
| j48 | Decision Tree |
| Log | Logistic Regression |
| NB | Naive Bayes |
| Rule | Rule Learner |
| Lin | Linear Regression |
| NN | Neural Network |

**Table 6: Process plans ranked by estimated profit, showing actual profit and gain over default strategy**

What did the winner(s) do differently? They did not use more complicated mining algorithms. Rather, *they used a different DM process*, one that is known by specialists to be particularly effective for target marketing. Specifically, as shown in Figure 8, a class probability estimator (CPE) is built to estimate the probability of donation; separately, a regression model is built (from the donors in the training set) to estimate the amount to be donated conditioned on the presence of a donation. These two models are used in combination: the product of the two, for any individual, estimates his/her expected donation. If the expected donation is greater than the

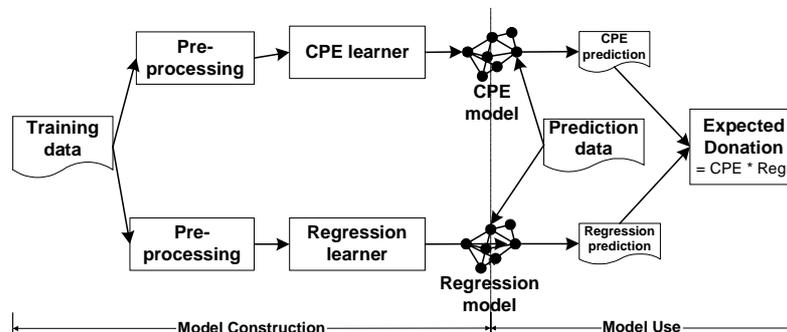cost of the promotion to that individual, in this case $0.68, then a mailing should be sent. Otherwise it should not.



**Figure 8: Target Marketing Process**

We claim that such process knowledge, in this case about how to combine techniques to form effective special-purpose DM processes, can be added to an IDA's ontology by specialists, subsequently to be brought to bear by others. The specialists can simply add the target marketing process as a problem-solving schema to the ontology. Note that there still is a large degree of freedom, even given such a process template. What type of learner should be used for class-probability estimation? What type of regressor? Given the learner, what type of pre-/post-processing is required? Using hierarchical planning, the IDA constructs DM processes within the constraints imposed by this template, in addition to the simpler, default template (which we used in previous sections).

For our final experiment, we considered the cost-sensitive plans built with both the default template and the plans built with the target-marketing template. In order not to bias the ranking with our prior knowledge (we know what the winners did), we use only auto-experimentation (cross-validation) to rank processes. In addition to the six process plans produced with the default (linear) DM process template, using the target-marketing template produces eight additional plans: the cross product of the available CPE learners (four) and the available regression learners (two). All the plans then are ranked by their estimated profit, produced via cross-validation on the training set. If one plan were to be submitted to a contest such as the KDDCUP competition, it would be the highest-ranking plan. Of course, we have the luxury of examining the entire list.

The fourteen process plans, ranked by cross-validated estimated profit, are listed in Table 7 along with their test-set profits and the percentage gain (loss) over the default mailing strategy. The estimated ranking reflects the actual profit ranking quite well (with a couple notable glitches;

Spearman's $r_s = 0.798$).  Indeed, the range of gains is remarkably similar to the actual ranking of submissions to the contest (note that we excluded processes such as: `(just) build a simple decision tree`, which produces zero profit).  The top-ranked plans indeed are competitive with the winners' submissions.  The penultimate plan is the one used by the winning submission, and performs comparably in terms of profit.  We did not expect the IDA to perform this well, because we figured SAS and Urban Science must have left some tricks up their sleeves (e.g., proprietary twists on the modeling algorithms).  The top-ranked process actually would have beaten the winning submission.

| Plan | Rank | Actual Profit | %Gain |
|---|---|---|---|
| Log(CPE) + NN | 1 | $14,914 | 41.23 |
| Log(CPE) + Lin | 2 | $14,778 | 39.95 |
| Rule(CPE) + NN | 3 | $13,672 | 29.47 |
| Rule(CPE) + Lin | 4 | $13,456 | 27.42 |
| DT(CPE) + NN | 5 | $11,055 | 4.69 |
| NN | 6 | $6,919 | -34.48 |
| DT(CPE) + Lin | 7 | $10,843 | 2.68 |
| Lin | 8 | $11,968 | 13.33 |
| Log(CPE) | 9 | $10,520 | -0.37 |
| NB(CPE) + NN | 10 | $10,070 | -4.64 |
| RULE(CPE) | 11 | $9,924 | -6.02 |
| NB(CPE) | 12 | $9,538 | -9.68 |
| NB(CPE) + Lin | 13 | $10,113 | -4.23 |
| DT(CPE) | 14 | $8,496 | -19.54 |

| Legend for Operators Used in Plans | |
|---|---|
| acronym | name/algorithm |
| DT | Decision Tree |
| Log | Logistic Regression |
| NB | Naive Bayes |
| Rule | Rule Learner |
| Lin | Linear Regression |
| NN | Neural Network |
| CPE | Class Prob. Estimator |

**Table 7:  Process plans ranked by estimated profit, showing actual profit and gain over default strategy**

These results illustrate not only the power of the IDA generally to enumerate and to rank processes effectively, but also the power of the IDA as a knowledge-sharing device.  If one specialist includes knowledge about the target-marketing process, and another includes knowledge about neural networks, and yet another includes knowledge about logistic regression, other users would benefit from the IDA's composition of these to form a top-performing DM process. Furthermore, the example illustrates how the specialist knowledge about the target marketing process together with the systematic exploration of solutions by an IDA can lead to surprisingly useful results.

## 6  Related Work

An IDA provides users with non-trivial, personalized "catalogs" of valid DM-processes, tailored to their task at hand, and helps them to choose among these processes in order to analyze their data. We know of little work that directly studies IDAs for the overall DM-process, although some have argued that they are important [Brazdil, 1998; Morik 2000]. There is, however, quite a long tradition of work that addresses some of the same goals (such as

recommending and ranking) or using similar techniques (e.g., planning, auto-experimentation, and the use of ontologies) for recommending and for ranking individual induction algorithms.

## 6.1   The Use of IDAs

Especially in the European community, researchers have argued for the importance of IDAs. Morik [2000], for example proposes to use a case-based repository to store successful chains of pre-processing operators.[17]   As pre-processing chains are partial DM processes, the insights gained should complement our work, and ideally could be integrated with a system such as IDEA.  The European MetaL project[18] has as one of its foci IDA-like systems; we are not aware of any existing system that uses background knowledge and/or experimentation to compose and rank DM *processes*, although Brazdil argues that it is important to do so [Brazdil, 1998].

   The only implemented IDA-like system we are aware of was presented by Engels *et al.*, who describe a user-guidance module for DM processes called CITRUS ([Engels, 1996], [Engels *et al.*, 1997], [Wirth *et al.*, 1997], and [Verdenius and Engels, 1997]).  In particular, the user-guidance module uses a task/method decomposition [Chandrasekaran *et al.,* 1992] to guide the user through a stepwise refinement of a high-level DM process, in order to help the user to construct the best plan using a limited model of operations. Finished plans are compiled into scripts for execution. The system is implemented by extending SPSS Inc.'s Clementine® system, which provides a visual interface to construct DM-processes manually.

   This work is similar to our approach as it provides the user with assistance when constructing DM processes, and uses AI planning techniques. In contrast, our approach is based on two notions that have led us in a different direction.  First, even with a well-specified goal it is very difficult to discern the one best plan, because the results of running data-mining methods are difficult to predict.  Secondly, users' goals and desired tradeoffs often cannot be specified easily or completely at the onset of an investigation.  This is because many desiderata are tacit and difficult to specify precisely (e.g., one may have an aversion to certain representations, based on experience with the domain experts).  Moreover, knowledge discovery is an exploratory process; users must be given as much flexibility as possible.  An IDA presents the user with many valid

---

[17] see http://www-ai.cs.uni-dortmund.de/FORSCHUNG/PROJEKTE/MININGMART/index.eng.html
[18] MetaL stands for "Meta-Learning," the process of learning models of the performance of learning algorithms as a function of characteristics of data sets; see http://www.cs.bris.ac.uk/Research/MachineLearning/metal.html.

plans to choose from and helps him/her to choose among them, via rankings based on different criteria (and on combinations thereof). The user has no obligation to choose the highest-ranked plan in any given ranking—all of the plans in the ranking will be valid.

Akin to the approach of Engels and colleagues, Buntine et al. [1999] introduce a method that generates data-analysis programs using program synthesis based on a declarative specification of the data-analysis problem. The declarative specification is a generalized Bayesian network. This approach is similar to ours, in that it composes new KD programs from a declarative problem specification. It differs from our approach in that it attempts to synthesize the one best program for data analysis (based on a optimization specification) rather then providing the user with a series of options and help to trade off the attributes of the different programs. The use of a deductive reasoning system for process synthesis is attractive in that it allows guiding the planning process using new declarations rather than changing the planner.

## 6.2  Projects with Related Goals: Recommending and Ranking

A variety of research projects address issues regarding recommending/selecting optimal induction algorithms (rather than processes) and ranking induction algorithms. The knowledge generated from such projects could help to populate an IDA's ontology, as well as to inform the construction of more advanced functions for ranking processes. The MLT-Consultant [Craw 1992] was one of the first such systems. It used a MYCIN-style knowledge base [Davis 1984] with a Hypertext-based GUI to recommend to a user an algorithm to choose (from a machine-learning library). Several projects have since studied the selection of individual induction algorithms or subcomponents of algorithms based on certain forms of background knowledge. For example, Brodley [1995] chooses subcomponents to form a hybrid decision tree, based on expert knowledge of algorithm applicability. In Europe the StatLog project[19] [Michie et al., 1994] has investigated what induction algorithms to use given particular circumstances. Brazdil et al. [1994], Gama & Brazdil [1995], and others, use meta-rules drawn from experimental studies, to help predict which algorithms will be better; the rules consider measurable characteristics of the data (e.g., number of cases, number of attributes, kurtosis). This notion of "meta-learning" is the basis for the MetaL project, mentioned above. Finally, Hilario & Kalousis [2001] use a case-

---

[19] see http://www.ncc.up.pt/liacc/ML/statlog/index.html

based system to advise users regarding which induction algorithm (and its respective parameter settings) to choose given a particular data-mining task.

A different tradition of meta-level systems for data mining [Buchanan et al., 1978], sometimes called "automatic bias selection," involves the selection of one of the following, based in part on feedback from the performance of the learner: vocabulary terms, the induction algorithm itself, components of the induction algorithm, parameters to the induction algorithm [desJardins and Gordon, 1995]. Bias-selection work generally assumes the goal is accuracy maximization, but also applies to other desiderata [Tcheng et al., 1989; Provost and Buchanan, 1995].

Addressing the need for improved ranking methods, several research projects have studied the use of experimental comparison to rank individual induction algorithms. Brazdil [1998] summarizes some prior methods. This work is closely related to our ranking of DM processes (especially since one may put a conceptual box around a DM process and call it an induction algorithm, although this obscures important issues regarding the composition of processes). More recently, Brazdil and Soares have studied the ranking of individual induction algorithms, based on (functions of) their performances on previously seen data sets [Brazdil & Soares, 2000; Soares and Brazdil, 2000]. They compare various methods for ranking, which perform comparably, and they consider ranking combining accuracy and speed.

## 6.3 Projects using similar techniques: Landmarking, Planning, Knowledge Management, and Ontologies

As we have seen, many of the component methods necessary for building IDAs have been the subject of recent study, especially in the European community. Several researchers have studied the notion of using fast processes (of different sorts) to help estimate the performance of less efficient ones. Pfahringer *et al*. [2000] and Fürnkranz and Petrak [2001] provide overviews of such "landmarking" techniques. In particular, Petrak [2000] shows the effectiveness of using subsamples from the data set in question to predict which learning algorithm will yield the lowest error on the entire data set; the technique works remarkably well—although it should be noted that for large data sets often one can achieve high accuracy with a surprisingly small subset of the data (cf., progressive sampling [Provost *et al*., 1999]). On the other hand, the relative performance of algorithms can change markedly with the amount of data [Perlich, *et al*., 2001].

One important related research area is the domain of statistical expert systems (cf. [Gale 1985], [Hand, 1994]), which are generally concerned with providing statistical advice. Most of the systems we are aware of based their advice on some statistical strategy, which is defined by Olford and Peters [1985] as "… the reasoning used by the experienced statistician in the course of analysis of some aspect of a substantive statistical problem" (p. 337). Typically, those strategies are hand-coded to contain the multiple analysis alternatives of different problems such as regression analysis [Oldford, 1997]. They help to guide the analysis of data by a human, to inform about which steps are likely to work next, and to allow direct execution. In contrast to our approach they do not offer support for a systematic exploration of the design space of possible processes (beyond the hand-coded strategies) nor for their relative rankings.

St. Amant and Cohen [1998] study intelligent, computer-based support for open-ended, statistical/exploratory data analysis. While focusing on somewhat different application areas—their approach on statistical, exploratory data analysis and ours on the DM process—both approaches employ mixed-initiative planning, where an AI-planner proposes different courses of action. The two approaches differ, however, in how the human and the machine share control of the process. Statistical/exploratory analysis necessitates step-by-step guidance, where the user can evaluate each step and get advice on what to do next. Our approach, on the other hand, presents the user with all possible plans and forecasts of their (relative) performance, allowing the user to choose one (or more) of the plans, run it, and then re-run the system based on insights gained. This latter approach is better suited in a domain (like knowledge discovery) where algorithms may run for extended periods of time. It may be worthwhile to create a hybrid approach that combines step-by-step guidance with overall planning allowing for the support of both types of data analysis.

The European Mining Mart project [Morik and Scholz, 2003] stores best-practice cases of pre-processing chains that were developed by experienced users. The project developed a data-mining workbench, which allows users to draw on a case base to develop new data-mining processes. Given its close integration with a database it focuses on very large, real-world data sets. Mining Mart is similar to our approach, in that it provides process-oriented discovery assistance to users based on an operator meta-model [Scholz and Euler, 2002] and a case-base. It differs in that it does not provide users any planning facility or active support while choosing among the

cases stored in the database. Finally, even though the Mining Mart meta model is richer then the one we chose (it describes not only the operators but also captures meta-data about the dataset), it does not seem to take advantage of inheritance features, which can vastly simplify the implementation and engineering of an ontology. Our approach is complimentary to Mining Mart and much could be gained from combining the strengths of both approaches.

Kerber *et al.* [1998] document the DM process using active links to DM processes (that have been visually programmed) and to the rationale for major design choices. They collect these descriptions in a repository. This approach facilitates the reuse of DM processes, resulting in a knowledge management system for DM processes. It is complementary to our approach, as it emphasizes the documentation and retrieval of past knowledge, which could be integrated well with our notion of active support as represented by IDAs.

The only work of which we are aware that uses an explicit ontology within a meta-level machine-learning system is described by Suyama & Yamaguchi [1998]. As far as we can tell, this system uses the ontology to guide the composition, by genetic programming, of fine-grained induction algorithm components.

## 7  Discussion, Limitations and Future Work

We have argued for a systematic exploration of the design space of DM processes, without which users (even experts) seldom are systematic in their search of the DM-process space and therefore may overlook important, useful DM processes. Our IDA does not mimic the behavior of experts, who often use heuristics to pre-prune the hypothesis space to a "consideration set." This often leads experts to overlook excellent solutions, which lie outside of their consideration set [Ulrich and Eppinger 1995].

For emphasis we have discussed novice users and expert users. However, this is not a true dichotomy—there is a spectrum of expertise along which users reside. For the most novice, any help with DM process planning will be helpful. For the most expert, an IDA could be useful for double-checking, and for automating previously manual tasks, as well as for suggesting additional processes. For others along the expertise spectrum, IDAs will have both types of benefits. Furthermore, even among experts, different users have different expertise: a data miner trained in the statistics community and a data miner from the machine-learning community can be experts and novices with respect to different methods. An IDA may help to educate any user. For ex-

ample, when the system produces a highly ranked plan that a user had not considered previously, the user can examine the ontology, and become educated on some new aspect of the DM process.

A unique benefit of an explicit ontology is the synergy it can provide between teams of users. If users contribute to the ontology, other users instantly receive the benefit of their contributions. Thus, an IDA may exhibit what economists call *network externalities* or *network effects*: the value to each user increases as the "network" grows. An IDA becomes more valuable to each user as the number of contributing users grows. All users get the benefit of each contributor's work automatically. No single member must be expert in all data-mining technology.

Consider the following example of network effects in action. Jill is a member of a large team of data miners, with several on-going projects. While reading the statistics literature she discovers a technique called *dual scaling* [Nishisato, 1994], a preprocessing operator that transforms categorical data into (scaled) numeric data, in a manner particularly useful for classification. Jill codes up a new preprocessor (call it DS) and uses it in her work. Such discoveries normally are isolated; they do not benefit a team's other projects. However, consider what happens if Jill simply adds DS into the IDA. When another team member, Jack, uses the system, DM-process plans may be generated that use DS (when appropriate). In some cases, these plans will be highly ranked (when DS is likely to do a good job satisfying Jack's criteria). In such cases, Jack could experiment with DS immediately, or could read about it (using the documentation Jill added), or could follow pointers to the literature, or could call Jill directly and talk to her about it. Thereby, the tool brings to bear shared knowledge in the context of a particular need.

While we have provided no true experimental support for this assertion (adding the target-marketing template did greatly improved the performance of IDEA for the KDDCUP-1998 problem), empirical studies of the social aspects of knowledge sharing provide support for analogous claims in different application domains. Pentland [1992], for example, shows how workers in a software hotline use a shared database as a central knowledge-sharing tool to become more effective as a team. Ackerman and Mandel [1999] show how a knowledge-sharing tool helps scientists (astrophysicists) to learn from each other how to perform specific data analysis tasks.

We are not suggesting automating the DM process totally; in contrast, intricate user interaction is critical to successful discovery. We have shown that it is possible to provide automated,

knowledge-based *assistance* for certain aspects of DM process design. We only have covered a few aspects so far, and for the most part only in the prototypical linear process. For example, our current prototype does not produce cyclic processes and our code generator does not yet produce code for more-complicated components, such as iterative feature selection (e.g., around a sub-process), wrappers for parameter tweaking, or progressive sampling. This is the subject of on-going work—we do not believe that there are fundamental roadblocks. However, it should be clear that the space of DM processes will grow, and more knowledge or interaction may need to be brought to bear than is evident in the demonstrations we have provided here. On the other hand, this difficulty faces human data miners as well as IDAs, and the result seems to be that even expert humans end up using only a small set of tools, those with which they are familiar. Even a moderately effective IDA would expand this set.

Our experiments with rankings serve to demonstrate that valid processes can be ranked effectively. As stated above, we have not yet studied the production of rankings in depth. Our IDA ranks the enumerations by characteristics such as speed, accuracy, and model comprehensibility. Some of those desiderata, such as speed and accuracy, have clear objective measures. Others are highly subjective. A statistician, for example, might rate a logistic regression equation as being very comprehensible, whereas a physician might not. A decision tree, on the other hand, might have the opposite result. Such preferences could be entered directly into a user-specific ontology, or could be discovered using relevance feedback methods by the IDA.

The related work on ranking induction *algorithms* should be very helpful for designing IDAs, but also provides important caveats. For example, our use of the Spearman rank-correlation co-efficient in effect weights equally the positions throughout a ranking. However, for our purposes, the processes near the top of the ranking probably would be the critical ones (especially for large number of generated process plans). Soares et al. [Soares, Costa & Brazdil, 2001] introduce a weighted modification to Spearman's coefficient, that takes into account position in the ranking. This same group of researchers also point out other challenges in comparing rankings, stemming from the fact that the "ideal" ranking typically is based only on estimates of the true error rates [Brazdil and Soares, 2000] [Soares, Brazdil & Costa, 2000].

Ranking DM processes using auto-experimentation raises the concern of multiple comparisons problems that are likely to make overfitting an issue. We shouldn't forget, though, that the exploration might also lead to better solutions. Furthermore, this problem is not only a problem of an IDA, but also of a manual exploration of multiple plans. One could argue, for example, that the data mining community is overfitting the UCI datasets. Hence, once the ontology is large there will be a chance that overfitting becomes an issue. Nevertheless, one needs to be aware of and deal with appropriately (if possible) in order to find better solutions through process planning.

An additional issue is that our current prototype does not rely on detailed meta-data (beyond attribute type). Exploring detailed meta-data to inform the ranking by extending the findings of the Statlog and METAL projects mentioned above to DM processes could lead to better ranking and we plan to explore this further.

Furthermore, we only have considered here parts of the process that are relatively well understood. Preprocessing existing variables, induction algorithms, and post-processing learned models have received considerable attention in the literature. Other parts of the process are not as well understood or documented. For example, although feature construction has received research attention for years, our understanding of when and how to use it effectively pales in comparison with our understanding of these other parts of the process. Consider the KDDCUP 1998 problem we presented above. We ignored the issue of feature construction, which (we assume) was crucial to success in the competition. Does enough knowledge exist to provide an IDA an ontology that will be effective to assist a user with feature construction? To our knowledge, this has yet to be shown convincingly. However, if generally effective methods or problem-specific heuristics exist, an IDA should be able to incorporate them. We also have assumed that the user will perform the selection of the discovery task(s) to perform. A separate task is intelligent assistance for the selection of discovery tasks. This typically is ignored in discussions of the knowledge discovery process, but was addressed in early knowledge discovery work by Lenat [1982] and recently by Livingston, Rosenberg, and Buchanan [2001a,b].

Finally, although studies such as this are necessary for the development of useful IDAs, we also need well-designed (and executed) user studies to assess whether IDAs actually are effective in helping real data miners. Such studies could also provide indication of which features of IDAs

are most effective in supporting the knowledge discovery process and, therefore, provide guidance for further improvements of IDAs.

# 8 Conclusion

Both novices and specialists need assistance in navigating the space of possible DM processes. We presented an ontology-based IDA, arguing that it can generate valid, non-trivial, and sometimes surprisingly interesting DM-process instances. Further, we have given empirical evidence that it is possible for IDAs to rank process instances effectively by speed and by accuracy, and have argued that they could rank by model comprehensibility. Finally, we have argued that IDAs can be particularly useful as a knowledge-sharing environment for teams of data miners, creating network effects wherein the tool becomes increasingly valuable as it gets more and more users.

The knowledge discovery process has been a key concept in the field of KDD for a decade, but very little research addresses it explicitly. After having undertaken this work, we understand better why. Treating the DM process requires a tremendous breadth of knowledge of research and practical technique. Even most researchers know only a fraction of what is necessary to do a comprehensive job of building an ontology (and we certainly have mistreated certain topics, although we have been careful). In retrospect, we believe even more strongly that in order for research on the knowledge discovery process to advance, systems like IDAs are essential—they document and automate parts of the process that are better understood, in order for research to concentrate on the large areas that are not well understood.

# References

[Ackerman & Mandel, 1999] M. Ackerman and E. Mandel. 1999. Memory in the Small: Combining Collective Memory and Task Support for a Scientific Community. Journal of Organizational Computing and Electronic Commerce, 9(2-3), 105-127.

[Agrawal, 1998] Agrawal, N., 1998, http://www.kdnuggets.com/meetings/kdd98/gain-kddcup98-release.html

[Ankolekar et al., 2001] Ankolekar, A., Burstein, M., Hobbs, J. R., Lassila, O., Martin, D. L., et al. 2001. DAML-S: Semantic Markup For Web Services. Presented at the First Semantic Web Working Symposium, Stanford University, California, USA.

[Blake & Merz, 2001] Blake, C.L. & Merz, C.J., 2000. UCI Repository of machine learning databases [http://www.ics.uci.edu/~mlearn/MLRepository.html].

[Brazdil et al., 1994] P. Brazdil, J. Gama, & B. Henery. Characterizing the Applicability of Classification Algorithms Using Meta-Level Learning. In *Proceeedings of the European Conference on Machine Learning (ECML-94),* pp. 83-102.

[Brazdil, 1998] Pavel. B. Brazdil. Data Transformation and Model Selection by Experimentation and Meta-learning. Proc. of ECML-98 Workshop on Upgrading Learning to the Meta-Level: Model Selection and Data Transformation, pp. 11-17.

[Brazdil & Soares, 2000] P. Brazdil & C. Soares. A Comparison of Ranking Methods for Classification Algorithm Selection. Proc. *11th European Conference on Machine Learning (ECML-2000)*: 63-74.

[Brodley, 1995] C. Brodley. Recursive Automatic Bias Selection for Classifier Construction. *Machine Learning* 20(1-2): 63-94.

[Brodley and Kohavi, 2000] Brodley, C. and R. Kohavi. KDD-Cup 2000 Peeling the Onion. Talk at KDD-2000, Boston, August 2000.

[Buchanan et al., 1978] Buchanan, B., C. Johnson, T. Mitchell, and R. Smith, 1978. Models of Learning Systems. In J. Beltzer, (ed.), *Encyclopedia of Computer Science and Technology*, pp. 24-51.

[Buntine et al., 1999] W. Buntine, B. Fischer, T. Pressburger: Towards Automated Synthesis of Data Mining Programs. *Proc. of the 5$^{th}$ Intl. Conf. on Knowledge Discovery in Databases*, 372-376.

[Chandrasekaran *et al.*, 1992] Chandrasekaran, B., Johnson, T. R., and Smith, J. W. 1992. Task-Structure Analysis for Knowledge Modeling. CACM 35, 9, 124-137

[Chapman et al. 2000] Pete Chapman, Julian Clinton, Randy Kerber, Thomas Khabaza, Thomas Reinartz, Colin Shearer, and Rüdiger Wirth, CRISP-DM 1.0: Step-by-step data mining guide, *SPSS White paper – technical report CRISPWP-0800,* SPSS Inc., 2000

[Christensen et al., 2001] E. Christensen, F. Curbera, G. Meredith, and S. Weerawarana. 2001. Web Services Description Language (WSDL) 1.1, W3C Note, 15 March 2001, http://www.w3.org/TR/wsdl

[Craw *et al.*, 1992] Craw, S., Sleeman, D., Graner, N., Rissakis, M., and Sharma, S. 1992. CONSULTANT: Providing Advice for the Machine Learning Toolbox. *Proc. of the Research and Development in Expert Systems IX* CUP, Cambridge, 5 - 23

[Davis 1984] R. Davis. Interactive Transfer of Expertise. In B. Buchanan and E. H. Shortliffe, editors. Rule-based Expert Systems, pages 171-205. Addison-Wesley, Reading, MA. 1984

[desJardins and Gordon, 1995] desJardin, M., and D. Gordon (1995) Special Issue on Bias Evaluation and Selection. Machine Learning 20, 1/2.

[Engels, 1996] Engels, R. 1996. Planning Tasks for Knowledge Discovery in Databases; Performing Task-Oriented User-Guidance. Proc. of the International Conference on Knowledge Discovery & Data Mining AAAI-Press, Portland, OR, 170-175

[Engels et al., 1997] Engels, R., Lindner, G., and Studer, R. 1997. A Guided Tour through the Data Mining Jungle. Proc. of the 3rd Intl. Conf. on Knowledge Discovery in Databases.

[Fayyad & Irani, 1993] Fayyad, U.M., and K. B. Irani. Multi-Interval discretization of continuous-valued attributes for classification learning. Proc. of the Thirteenth International Joint Conference on Artificial Intelligence. San Francisco: Morgan Kaufman, pp. 1022-1027.

[Fayyad, Piatetsky-Shapiro & Smyth, 1996] Fayyad, U., Piatetsky-Shapiro, G. and P. Smyth. The KDD process for extracting useful knowledge from volumes of data. CACM 39 (11), pp. 27-34.

[Frank & Witten, 1998] Frank, E., and Witten, I. H. Generating Accurate Rule Sets Without Global Optimization, 1998. Proc. of the Fifteenth International Conference on Machine Learning, 144-151.

[Fürnkranz & Petrak, 2001] J. Furnkranz and J. Petrak. An Evaluation of Landmarking Variants. In C. Giraud-Carrier, N. Lavrac, and S. Moyle (eds.), ECML/PKDD 2001 Workshop on Integrating Aspects of Data Mining, Decision Support, and Meta-Learning, pp. 57-68.

[Gale, 1985] W. Gale. 1985. Artificial Intelligence and Statistics. Reading, MA: Addison-Wesley.

[Gama & Brazdil, 1995] J. Gama and P. Brazdil: Characterization of Classification Algorithms. 7th Portuguese Conference on Artificial Intelligence, EPIA '95, pp. 189-200.

[Ghallab et al., 1998] M. Ghallab, C. Howe, D. McDermott, A. Ram, M. Veloso, D. Weld., et al. (1998). PDDL - The Planning Domain Definition Language (Technical Report No. TR-98-003/DCS TR-1165). New Haven, CT: Yale Center for Computational Vision and Control.

[Hand, 1994] D. Hand. 1994. Statistical Expert Systems. Chance, 7(1), 28-34.

[Hilario & Kalousis 2001] M. Hilario and A. Kalousis, "Fusion of Meta-Knowledge and Meta-Data for Case-Based Model Selection," University of Geneva, Geneva UNIGE-AI-01-01.

[John & Langley, 1995] George H. John and Pat Langley. Estimating Continuous Distributions in Bayesian Classifiers. Proc. of the Eleventh Conference on Uncertainty in Artificial Intelligence. pp. 338-345. Morgan Kaufmann, San Mateo.

[Kerber et al., 1998] Kerber, R., Beck, H., Anand, T., and Smart, B. 1998. Active Templates: Comprehensive Support for the Knowledge Discovery Process. Intl. Conf. on Knowledge Discovery and Data Mining, 244-248

[Kohavi & Sahami, 1996] Ron Kohavi and Mehran Sahami, 1996. Error-Based and Entropy-Based Discretization of Continuous Features. Proc. of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), pp. 114-119.

[Lenat, 1982] D. Lenat, 1982. AM: Discovery in Mathematics as Heuristic Search. In R. Davis and D. Lenat (eds.), *Knowledge-based Systems in Artificial Intelligence*, pp. 3-255. New York: McGraw-Hill.

[Lim et al., 2000] Lim, T.-S., Loh, W.-Y., and Shih, Y.-S. A Comparison of Prediction Accuracy, Complexity, and Training Time of Thirty-Three Old and New Classification Algorithms. Machine Learning 40 (3), 203-228.

[Livingston, Rosenberg & Buchanan, 2001a] G. Livingston, J. Rosenberg, and B. Buchanan, 2001. Closing the Loop: an Agenda- and Justification-Based Framework for Selecting the Next Discovery Task to Perform. Proc. *of the 2001 IEEE International Conference on Data Mining*, pp. 385-392. IEEE Computer Society Press.

[Livingston, Rosenberg & Buchanan, 2001b] G. Livingston, J. Rosenberg, and B. Buchanan, 2001. Closing the Loop: Heuristics for Autonomous Discovery. Proc. *of the 2001 IEEE International Conference on Data Mining*, pp. 393-400. IEEE Computer Society Press.

[Michie et al., 1994] D. Michie, D. Spiegelhalter, and C. Taylor (eds.). *Machine Learning, Neural and Statistical Classification*. Ellis Horwood.

[Morik, 2000] Morik, K., *The Representation Race - Preprocessing for Handling Time Phenomena*. Proc. of the European Conference on Machine Learning 2000 (ECML 2000), Ramon López de Mántaras and Enric Plaza (ed.), Lecture Notes in Artificial Intelligence, Vol. 1810, Springer Verlag Berlin, 2000.

[Morik and Scholz, 2003] K. Morik and M. Scholz, The MiningMart Approach to Knowledge Discovery in Databases. In Handbook of Intelligent IT, Ning Zhong and Jiming Liu (ed.), IOS Press, 2003, to appear.

[Nishisato, 1994] Nishisato, S. *Elements of Dual Scaling: An Introduction to Practical Data Analysis.* Hillsdale: Lawrence Erlbaum Associates.

[Oates & Jensen, 1997]. Oates, T. and D. Jensen, 1997. The effects of training set size on decision tree complexity. Proc. of the Fourteenth International Conference on Machine Learning, pp. 254-262.

[Oldford & Peters, 1985] R. Oldford, & S. Peters. 1985. Implementation and Study of Statistical Strategy. In W. A. Gale (Ed.), Artificial Intelligence and Statistics (pp. 335-353). Reading, MA: Addison-Wesley.

[Oldford, 1997] R. Oldford. 1997. Computational Thinking for Statisticians: Training by Implementing Statistical Strategy. *Presented at the 29th Symposium on the Interface*, Houston TX.

[Pentland, 1992] B. Pentland. 1992. Organizing Moves in Software Support Hot Lines. Administrative Science Quarterly, 37(4), 527-548.

[Perlich, et al., 2001] C. Perlich, F. Provost, and J. Simonoff. "Tree Induction vs. Logistic Regression: A Learning-curve Analysis." CeDER Working Paper #IS-01-02, Stern School of Business, New York University, NY, NY 10012.

[Petrak, 2000]. J. Petrak. Fast Subsampling Performance Estimates for Classification Algorithm Selection. Tech. Report TR-2000-07. Austrian Research Institute for Artificial Intelligence.

[Pfahringer, et al., 2000]. Meta-Learning by Landmarking Various Learning Algorithms. In P. Langley (ed.), Proc. of the Seventeenth International Conference on Machine Learning (ICML-2000), pp. 743-750. San Francisco: Morgan Kaufmann.

[Provost & Buchanan, 1995] F. Provost & B.. Buchanan. Inductive Policy: The Pragmatics of Bias Selection. Machine Learning 20(1-2): 35-61 (1995).

[Provost et al. 1999] F. Provost, D. Jensen, and T. Oates. Efficient Progressive Sampling. *Proc. of the 5th ACM SIGKDD Intl. Conf. on Knowledge Discovery and Data Mining*, pp. 23-32.

[Provost & Kolluri, 1999]. Provost, F. & V. Kolluri. A Survey of Methods for Scaling Up Inductive Algorithms. *Data Mining and Knowledge Discovery* 3 (2):131-169, 1999.

[Quinlan, 1987] J. R. Quinlan: Simplifying Decision Trees. *International Journal of Man-Machine Studies* 27(3): 221-234 (1987)

[Quinlan, 1993] J. R. Quinlan (1993). *C4.5: Programs for Machine Learning*. San Mateo, CA: Morgan Kaufmann.

[Scholz and Euler, 2002] M. and T. Euler, Documentation of the MiningMart Meta Model (M4). TechnicalReport, TR12-05, IST Project MiningMart, IST-11993, 2002.

[Senator, 2000] T. Senator (2000). Ongoing Management and Application of Discovered Knowledge in a Large Regulatory Organization: A Case Study of the Use and Impact of NASD Regulation's Advanced Detection System (ADS). Proc. *of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 44-53.

[Soares & Brazdil, 2000] C. Soares & P. Brazdil. Zoomed Ranking: Selection of Classification Algorithms Based on Relevant Performance Information. Proc. *of Principles of Data Mining and Knowledge Discovery, 4th European Conference (PKDD-2000),* 126-135.

[Soares, Brazdil & Costa, 2000] C. Soares, P. Brazdil, & J. Costa. Measures to Evaluate Rankings of Classification Algorithms. In H. Kiers, J. Rasson, P. Groenen and M. Schader (eds.), *Data Analysis, Classification, and Related Methods, Proc. of the Seventh Conference of the International Federation of Classification Societies (IFCS),* pp. 119-124. Springer.

[Soares, Costa & Brazdil, 2001] C. Soares, P. Brazdil, & J. Costa. Improved Statistical Support for Matchmaking: Rank Correlation Taking Rank Importance into Account. In *JOCLAD 2001: VII Jornadas de Classificao e Analise de Dados*, pp. 72-75.

[St. Amant & Cohen, 1998] St. Amant, R., and Paul R. Cohen. Intelligent Support for Exploratory Data Analysis. Journal of Computational and Graphical Statistics, Volume 7, Number 4, Pages 545–558

[Suyama & Yamaguchi 1998] A. Suyama & T. Yamaguchi. Specifying and Learning Inductive Learning Systems using Ontologies. In *Working Notes from the 1998 AAAI Workshop on the Methodology of Applying Machine Learning: Problem Definition, Task Decomposition and Technique Selection,* pp. 29-36.

[Tcheng, et al., 1989] D. Tcheng, B. Lambert, S. Lu, L. Rendell: Building Robust Learning Systems by Combining Induction and Optimization. IJCAI 1989: 806-812.

[Turney, 1996] Turney, P., 1996. Cost-sensitive Learning Bibliography. http://extractor.iit.nrc.ca/bibliographies/cost-sensitive.html.

[Ulrich and Eppinger, 1995] Ulrich, K. T., and Eppinger, S. D. 1995. Product Design and Development. McGraw-Hill, New York, NY,

[Verdenius and Engels, 1997] Verdenius, F., and Engels, R. 1997. A Process Model for Developing Inductive Applications. Proc. of the Seventh Belgian-Dutch Conference on Machine Learning Tilburg, NL, 119-128

[Wirth et al., 1997] Wirth, R., et al. Towards Process-Oriented Tool Support for KDD. Proc. of the 1st European Symposium on Principles of Data Mining and Knowledge Discovery Trondheim, Norway

[Witten and Frank, 2000] Witten, I., and E. Frank (2000). Data Mininig: Practical machine learning tools and techniques with Java implementations. San Francisco: Morgan Kaufmann.

[Zadrozny and Elkan, 2001] Zadrozny, B. and C. Elkan (2001). Learning and Making Decisions When Costs and Probabilities are Both Unknown. *Proc. of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.