

SECRET: A Scalable Linear Regression Tree Algorithm

Alin Dobra
Department of Computer Science
Cornell University
Ithaca, NY 14853
dobra@cs.cornell.edu

Johannes Gehrke
Department of Computer Science
Cornell University
Ithaca, NY 14853
johannes@cs.cornell.edu

ABSTRACT

Developing regression models for large datasets that are both accurate and easy to interpret is a very important data mining problem. Regression trees with linear models in the leaves satisfy both these requirements, but thus far, no truly scalable regression tree algorithm is known. This paper proposes a novel regression tree construction algorithm (SECRET) that produces trees of high quality and scales to very large datasets. At every node, SECRET uses the EM algorithm for Gaussian mixtures to find two clusters in the data and to locally transform the regression problem into a classification problem based on closeness to these clusters. Goodness of split measures, like the gini gain, can then be used to determine the split variable and the split point much like in classification tree construction. Scalability of the algorithm can be achieved by employing scalable versions of the EM and classification tree construction algorithms. An experimental evaluation on real and artificial data shows that SECRET has accuracy comparable to other linear regression tree algorithms but takes orders of magnitude less computation time for large datasets.

1. INTRODUCTION

Regression is a very important data mining problem. One very important class of regression models is regression trees. Even though regression trees were introduced in the CART book early in the development of decision trees by Breiman et al. [3], regression tree construction has received far less attention from the research community so far. Quinlan generalized CART regression trees by using a linear model in the leaves to improve the accuracy of the tree [14]; he used as impurity measure the standard deviation of the predictor variable. Karalilc argued that the mean square error of the linear model in a node is a more appropriate impurity measure for linear regression trees since large variance in the data is no indication of the fit of the linear model [9]. Evaluation of the variance is much easier than estimating the error of a linear model (which requires solving a lin-

ear system). Even more, if some of the predictor attributes are discrete, the problem finding the best split attribute for binary trees becomes intractable for linear regression trees since the theorem that justifies a linear algorithm for finding the best split (Theorem 9.4 in [3]) does not seem to apply. To address computational concerns of normal linear regression models, Alexander and Scott proposed the use of simple linear regressors (i.e., the linear model depends on only one predictor attribute), which can be trained more efficiently but are not as accurate [1].

Torgo proposed the use of even more sophisticated functional models in the leaves (i.e., kernel regressors) [17, 16]. For such regression trees both construction and deployment of the model is expensive but they potentially are superior to linear regression trees in terms of accuracy. More recently, Li et al. proposed a linear regression tree algorithm that can produce oblique splits using Principal Hessian Analysis but the algorithm cannot accommodate discrete attributes [10].¹

There are a number of contributions coming from the statistics community. Chaudhuri et al. [4] proposed the use of statistical tests for split variable selection instead of error of fit methods. The main idea is to fit a model (constant, linear or higher order polynomial) for every node in the tree and to partition the data at each node into two classes: data-points with positive residuals² and data-points with negative residuals. In this manner the regression problem is locally reduced to a classification problem, so it becomes much simpler. Statistical tests used in classification tree construction, Student's t-test in this case, can be used from this point on. Unfortunately, it is not clear why differences in the distributions of the signs of the residuals are good criteria on which decisions about splits are made. A further enhancement was proposed recently by Loh [11]. It consists mostly in the use of the χ^2 -test instead of the t-test in order to accommodate discrete attributes, the detection of interactions of pairs of predictor attributes, and a sophisticated calibration mechanism to ensure the unbiasedness of the split attribute selection criterion.

In this paper we introduce SECRET (Scalable EM and Classification based Regression Trees), a new construction algorithm for regression trees with linear models in the leaves, which produces regression trees with accuracy comparable to the trees produced by existing algorithms and at the same time requiring far less computational effort on large

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGKDD '02 Edmonton, Alberta, Canada

Copyright 2002 ACM 1-58113-567-X/02/0007 ...\$5.00.

¹Oblique splits are linear inequalities involving two or more predictor attributes.

²Residuals are the difference between the true value and the value predicted by the regression model.

datasets. Our experiments show that SECRET improves the running time of regression tree construction by up to two orders of magnitude compared to previous work while at the same time constructing trees of comparable quality. Our main idea is to use the EM algorithm on the data partition in an intermediate node to determine two Gaussian clusters, hopefully with shapes close to flat disks. We then use these two Gaussians to locally transform the regression problem into a classification problem by labeling every data-point with class label 1 if the probability of belonging to the first cluster exceeds the probability of belonging to the second cluster, or class label 2 if the converse is true. A split attribute and a corresponding split point to separate the two classes can be determined then using goodness of split measures for classification trees like the gini gain [3]. Least square linear regression can be used to determine the linear regressors in the leaves.

The local reduction to a classification problem allows us to avoid forming and solving the large number of linear systems of equations required for an exhaustive search method such as the method used by RETIS [9]. Even more, scalable versions of the EM algorithm for Gaussian mixtures [2] and classification tree construction [7] can be used to improve the scalability of the proposed solution. An extra benefit of the method is the fact that good oblique splits can be easily obtained.

The rest of the paper is organized as follows. In Section 2 we give short introductions to classification and regression tree construction and to the EM algorithm for Gaussian mixtures. In Section 3 we present in greater detail some of the previously proposed solutions and we comment on their shortcomings. Section 4 contains the description of SECRET, our proposal for a linear regression tree algorithm. We show results of an extensive experimental study of SECRET in Section 5, and we conclude in Section 6.

2. PRELIMINARIES

In this section we give a short introduction to classification and regression trees. More details can be found in the excellent review [13].

2.1 Classification Trees

Classification trees have a tree organization with intermediate nodes labeled by split attributes, the branches starting in an intermediate node labeled by split predicates involving the corresponding split predicate and leaves labeled by class labels. Prediction using classification trees is made by navigating the tree on true predicates until a leaf is reached, when the corresponding label is returned.

The construction of classification trees proceeds in two phases. In the growing phase an oversized tree is build recursively, at every step a split attribute and split predicates involving this attribute are selected in order to maximize the goodness of split criteria. A large number of such criteria have been proposed in the literature [13], here we use Breiman's *gini gain* since is easy to compute and the best split for discrete attributes can be found efficiently [3]. In the second phase the final size of the tree is determined with the goal to minimize the error on unseen examples.

2.2 Regression Trees

Regression models or regressors are functional mappings from the cross product of the domains of predictor attributes

X_1, \dots, X_m to the domain of the continuous predicted attribute, Y . They only differ from classifiers in the fact that the predicted attribute is real valued.

Regression Trees, the particular type of regressors we are interested in, are the natural generalization of decision trees for regression problems. Instead of a class label being associated to every node, a real value or a functional dependency of some of the inputs is used to predict the value of the output.

Regression trees in CART have a constant numerical value in the leaves and use the variance as a measure of impurity [3]. Thus the split selection measure for the node T with children nodes $T_1 \dots T_n$ is:

$$\begin{aligned} \text{Err}(T) &\stackrel{\text{def}}{=} E[(Y - E[Y|T])^2] \\ \Delta\text{Err}(T) &= \text{Err}(T) - \sum_{j=1}^n P[T_j|T] \cdot \text{Err}(T_j) \end{aligned} \quad (1)$$

The use of variance as the impurity measure is justified by the fact that the best constant predictor in a node is the expected value of the predicted variable given that data-points belong to the node, $E[Y|T]$; the variance is thus the mean square error of this best predictor.

As in the case of classification trees, prediction is made by navigating the tree following the branches with predicates that are satisfied until a leaf is reached. The numerical value associated with the leaf is the prediction of the model.

As in the case of classification trees, a top-down induction schema is usually used to build regression trees. Pruning is used to improve the accuracy on unseen examples. Pruning methods for classification trees can be straightforwardly adapted for regression trees [18]. In this paper we use Quinlan's resubstitution error pruning [15]. It consists of eliminating subtrees in order to obtain a tree with the smallest error on the pruning set.

3. PREVIOUS APPROACHES

In this section we analyze some of the previously proposed construction algorithms for linear regression trees and, for each, we point major drawbacks.

3.1 Quinlan's Algorithm

For efficiency reasons, the algorithm proposed by Quinlan [14] *pretends* that a regression tree with constant models in the leaves is constructed until the tree is fully grown, when linear models are fit on the data-points available at leaves. This is equivalent to using the split selection criterion in Equation 1 during the growing phase. Then linear regressors in the leaves are constructed by performing another pass over the data in which the set of data-points from the training examples corresponding to each of the leaves is determined and the least square linear problem for these data-points is formed and solved (using the SVD decomposition [8]).

The same approach was later used by Torgo [16, 17] with more complicated models in the leaves like kernels and local polynomials.

As pointed out by Karalic [9] the variance of the output variable is a poor estimator of the purity of the fit when linear regressors are used in the leaves, since the points can be arranged along a line (so the error of the linear fit is almost zero) but they occupy a significant region (so the

variance is large). To correct this problem, he suggested that the following impurity function should be used:

$$\text{Err}_l(T) \stackrel{\text{def}}{=} E [(Y - E[f(\mathbf{X})|T])^2] \quad (2)$$

where $f(\mathbf{x}) = [1 \ \mathbf{x}^T] \mathbf{c}$ is the linear regressor with the smallest least square error.

To see more clearly that $\text{Err}(T)$ given by Equation 1 is not appropriate for the linear regressor case, consider the situation in Figure 1. The two thick lines represent a large number of points (possibly infinite). The best split for the linear regressor is $x = 0$ and the fit is perfect after the split.

For the case when $\text{Err}(T)$ is used, it can be shown analytically that the the split point is $-(\sqrt{5} - 1)/2 = -0.618034$ or symmetrically 0.618034. Either of these splits is very far from the split obtained using $\text{Err}_l(T)$ (at point 0), thus splitting the points in proportion 19% to 81% instead of the ideal 50% to 50%.

This example suggests that the split point selection based on $\text{Err}(T)$ produces an unnecessary fragmentation of the data that is not related to the natural organization of the data-points for the case of linear regression trees. This fragmentation produces unnecessarily large and unnatural trees, anomalies that are not corrected by the pruning stage.

3.2 Karalic’s Algorithm

Using the split criterion in Equation 2 the problem mentioned above is avoided and much higher quality trees are built. If exhaustive search is used to determine the split point, the computational cost of the algorithm becomes prohibitively expensive for large datasets for two main reasons:

- If the split attribute is continuous, all possible values of this attribute have to be considered as split points. For each of them a linear system has to be formed and solved. Even if the matrix and the vector that form the linear system are maintained incrementally (which can be dangerous from numerical stability point of view), for every level of the tree constructed, a number of linear systems equal to the size of the dataset have to be solved.
- If the split attribute is discrete the situation is worse since Theorem 9.4 in [3] does not seem to apply for this split criterion. This means that an exponential number, in the size of the domain of the split variable, of linear systems has to be formed and solved.

The first problem can be alleviated if a sample of the points available are considered as split points. Even if this simplification is made, the data-points have to be sorted in every intermediate node on all the possible split attributes. Also, it is not clear how these modifications influence the accuracy of the generated regression trees. The second problem seems unavoidable if exhaustive search is used.

3.3 Chaudhuri’s Algorithm

In order to avoid forming and solving so many linear systems, Chaudhuri et al. [4] proposed to locally classify the data-points available at an intermediate node based on the sign of the residual with respect to the least square error linear model. For the dataset shown in Figure 1 this corresponds to points above and below the dashed line. As can be observed, when projected on the X axis, the negative class surrounds the positive class so two split points

are necessary to differentiate between the two groups (the node has to be split into three parts). When the number of predictor attributes is larger than one (the multidimensional case), the separating surface between the class labels + and - is nonlinear in general. Moreover, even if the best regressors are fit to the two classes, the prediction is only slightly improved. The solution adopted by Chaudhuri et al. is to use Quadratic Discriminant Analysis (QDA) to determine the split point.

4. SECRET

For constant regression trees, algorithms for scalable classification trees can be straightforwardly adapted [7]. The main obstacle in doing a similar adaptation for linear regression trees is that the problem of partitioning the domain of a discrete variable in two parts is intractable. Also the amount of sufficient statistics that has to be maintained for the split increases drastically: For constant regressors, we have to maintain two real numbers (the mean and the sum of squares), whereas for linear trees the statistics are quadratic in the number of regression attributes (to maintain the matrix $A^T A$ that defines the linear system) and thus not feasible in practice.

In this work we distinguish as follows between different predictor attributes [11]: (1) *discrete attributes* – used only in the split predicates in intermediate nodes in the regression tree, (2) *split continuous attributes* – continuous attributes used only for splitting, (3) *regression attributes* – continuous attributes used in the linear combination that specifies the linear regressors in the leaves as well as for specifying split predicates. By allowing some continuous attributes to participate in splits but not in regression in the leaves we add greater flexibility to the learning algorithm. The partitioning of the continuous attributes in split and regression is beyond the scope of the paper (and is usually performed by the user [11]).

The main idea behind our algorithm is to locally transform the regression problem into a classification problem by first identifying two general Gaussian distributions in the regressor attributes–output space using the EM algorithm for Gaussian mixtures and then classifying the data-points based on the probability of belonging to these two distributions. Classification tree techniques are then used to select the split attribute and the split point. Our algorithm, called SECRET, is summarized in Figure 4.

The role of the EM Algorithm is to find two natural classes in the data that have an approximately linear organization. The role of the classification is to identify predictor attributes that can differentiate between these two classes in the input space. To see this more clearly, suppose we are in the process of building a linear regression tree and we have to decide on the split attribute and split point for the node T . Suppose the set of training examples available at node T contains tuples with three components: a regressor attribute X_r , a discrete attribute X_d and the predicted attribute Y . The projection of the training data on the X_r, Y space might look like Figure 2. The datapoints are approximately organized in two clusters with Gaussian distributions that are marked as ellipsoids. Differentiating between the two clusters is crucial for prediction, but information in the regression attribute is not sufficient to make this distinction even though within a cluster they can do good predictions. The information in the discrete attribute

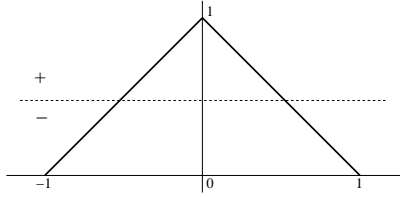


Figure 1: Example where classification on sign of residuals is unintuitive.

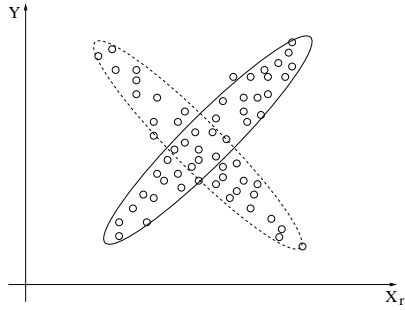


Figure 2: Projection on X_r, Y space of training data.

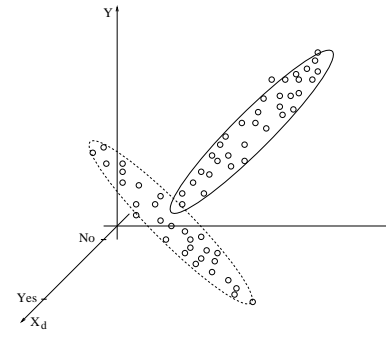


Figure 3: Projection on X_d, X_r, Y space of same training data as in Figure 2

Input: node T , data-partition D

Output: regression tree \mathcal{T} for D rooted at T

Linear regression tree construction algorithm:

BuildTree(node T , data-partition D)

- (1) Normalize data-points to unitary sphere
- (2) Find two Gaussian clusters in regressor-output space (EM)
- (3) Label data-points based on closeness to these clusters
- (4) **foreach** split attribute
- (5) Find the best split point and determine its gini gain
- (6) **endforeach**
- (7) Let X be the attribute with the greatest gini gain and Q the corresponding best split predicate
- (8) **if** (T splits)
- (9) Partition D into D_1, D_2 based on Q and label node T with split attribute X
- (10) Create children nodes T_1, T_2 of T and label the edge (T, T_i) with predicate $q_{(T, T_i)}$
- (11) **BuildTree**(T_1, D_1); **BuildTree**(T_2, D_2)
- (12) **else**
- (13) label T with the least square linear regressor of D
- (14) **endif**

Figure 4: SECRET algorithm

X_d can make this distinction, as can be observed from Figure 3 where the projection is made on the X_d, X_r, Y space. If more split attributes had been present, a split on X_d would have been preferred since the resulting splits are pure.

Observe that the use of the EM Algorithm for Gaussian mixtures is very limited since we have only two mixtures and thus the likelihood function has a simple form (resulting in fewer local minima). Since the EM Algorithm is quite sensitive to distances, we normalize the training data before running the algorithm. Our normalization performs a linear transformation that makes the data look as close as possible to a unitary sphere with the center in the origin. Experimentally, we observed that, with this transformation and in this restricted scenario, the EM algorithm with clusters initialized randomly works well.

Once the two Gaussian mixtures are identified, the datapoints can be labeled based on the closeness to the two clusters (i.e., if a datapoint is closer to cluster 1 than cluster 2 it is labeled with class label 1, otherwise it is labeled with class label 2). With this local labeling of points, attribute and split point selection methods from classification tree construction can be used.

We are using the gini gain as a representative split selection criteria to find the split point; we would like to empha-

size that any split selection criterion from the literature could be used. For each attribute (or collection of attributes for oblique splits) we determine the best split point and compute its gini gain. We then choose the predictor attribute with the largest gini gain as split attribute.

For discrete attributes the algorithm of Breiman et al. finds the split point in time linear in the size of the domain of the discrete attribute (since we only have two class labels) [3]. We directly use this algorithm to find the best split partition for discrete attributes.

4.1 Split Point Selection For Continuous Attributes

Since the EM algorithm for Gaussian mixtures produces two normal distributions, it is reasonable to assume that the projection of the datapoints with the same class label on a continuous attribute X has also a normal distribution. The split point that best separates the two normal distributions can be found using Quadratic Discriminant Analysis (QDA). The reason for preferring QDA to a direct minimization of the gini gain is the fact that it gives qualitatively similar splits but requires less computational effort [12]. The gini of the split can be computed directly from the parameters of the two normal distributions, an extra pass over the data is not necessary.

4.2 Oblique Splits

Ideally, given two Gaussian distributions, we would like to find the separating hyperplane that maximizes the gini gain. Fukunaga showed that the problem of minimizing the expected value of the 0-1 loss (the classification error function) generates an equation involving the normal of the hyperplane that is not solvable algebraically [6]. Following the same treatment, it is easy to see that the problem of minimizing the gini gain generates the same equation. A good solution to the problem of determining a separating hyperplane can be found using Linear Discriminant Analysis (LDA) [6] to determine the projection direction that gives the largest separation between the projections of the two distributions. QDA can then be used on the projection to determine a point contained in the separating hyperplane. This point and the best projection direction completely determine the separating hyperplane. As in the unidimensional case, the gini gain of the split can be determined from the parameters of the two distributions and the separating hyperplane, extra passes being unnecessary.

4.3 Finding Linear Regressors

If the current node is a leaf, or in preparation for the situation that all the descendents of this node are pruned, we have to find the best linear regressor that fits the training data. We identified two ways the LSE linear regressor can be computed.

The first method consist of a traversal of the original dataset and an identification of the subset that falls into this node. The least square linear system that gives the linear regressor is formed with these datapoints and solved. Note that, in the case that all the sufficient statistics can be maintained in main memory, a single traversal of the training dataset per tree level will suffice.

The second method uses the fact that the split selection method tries to find a split attribute and a split point that can differentiate best between the two Gaussian mixtures found in the regressor–output space. The least square problem can be solved at the level of each of these mixtures under the assumption that the distribution of the datapoints is normal with the parameters identified by the EM algorithm. This method is less precise as the data might not fit the Gaussian model, but the method can be used when the number of traversals over the dataset becomes a concern.

5. EMPIRICAL EVALUATION

In this section we present the results of an extensive experimental study of SECRET. The purpose of the study was twofold: (1) to compare the accuracy of SECRET with GUIDE [11], a state-of-the-art linear regression tree algorithm and (2) to compare the scalability properties of SECRET and GUIDE.

The two main findings of our study are:

- **Accuracy of prediction.** SECRET is more accurate than GUIDE on three datasets, as accurate on six datasets and less accurate on three datasets. This suggests that overall the prediction accuracy of SECRET is comparable to the accuracy of GUIDE. On four of the datasets, the use of oblique splits resulted in significant improvement in accuracy.
- **Scalability to large datasets.** For datasets of small to moderate sizes (up to 5000 tuples), GUIDE slightly outperforms SECRET. The behavior for large datasets of the two methods is very different. For datasets with 256000 tuples and 3 attributes, SECRET runs about 200 times faster than GUIDE. Even if GUIDE considers only 1% of the points available as possible split points, SECRET still runs 20 times faster. Also, there is no significant change in running time when SECRET produces oblique splits.

5.1 Experimental Setup

GUIDE is a regression tree construction algorithm that was designed to be both accurate and fast [11]. An extensive experimental study by Loh shows that GUIDE outperforms previous regression tree construction algorithms and compares favorably to MARS [5], a state-of-the-art regression algorithm based on splines [11]. GUIDE uses statistical techniques to pick the split variable and can use exhaustive search or just a sample of the points to find the split point. In our accuracy experiments we set up GUIDE to use exhaustive search. For the scalability experiments we report

Name	Source	# cases	# nominal	# continuous
Abalone	UCI	4177	1	7
Basball	UCI	261	3	17
Kin8nm	DVELVE	8192	0	8
Mpg	UCI	392	3	5
Mumps	SatLib	1523	0	4
Stock	SatLib	950	0	10
TA	UCI	151	4	2
Tecator	SatLib	240	0	11
Cart	Breiman et al.([3])	–	10	1
Fried	Friedman [5]	–	0	11
3DSin	$3 \sin(X_1) \sin(X_2)$	–	0	3

Table 1: Datasets Used In Experiments. Top: Real-life Datasets; Bottom: Synthetic Datasets

running times for both the exhaustive search and split point candidate sampling of size 1%.

For the experimental study we used nine real life and three synthetic datasets. Their characteristics are summarized in Table 1. All datasets except 3DSin have been used before extensively in experimental studies.

We performed all the experiments reported in this paper on a Pentium III 933MHz running Redhat Linux 7.2.

5.2 Experimental Results: Accuracy

For each experiment with real datasets we used a random partitioning into 50% of datapoints for training, 30% for pruning and 20% for testing. For the synthetic datasets we randomly generated 16384 tuples for training, 16384 tuples for pruning and 16384 tuples for testing for each experiment. We repeated each experiment 100 times in order to get accurate estimates. For comparison purposes we built regression trees with both constant (by using all the continuous attributes as split attributes) and linear (by using all continuous attributes as regressor attributes) regression models in the leaves. In all the experiments we used Quinlan’s resubstitution error pruning [15]. For both algorithms we set the minimum number of data-points in a node to be considered for splitting to 1% of the size of the dataset, which resulted in trees at the end of the growth phase with around 75 nodes.

Table 2 contains the average mean square error and its standard deviation for GUIDE, SECRET and SECRET with oblique splits (SECRET(O)) with constant (left part) and linear (right part) regressors in the leaves, on each of the twelve datasets. GUIDE and SECRET with linear regressor in the leaves have equal accuracy (we considered accuracies equal if they were less than three standard deviations away from each other) on six datasets (Abalone, Boston, Mpg, Stock, TA and Tecator), GUIDE wins on three datasets (Baseball, Mumps and Fried) and SECRET wins on the remaining three (Kin8nm, 3DSin and Cart). These findings suggest that the two algorithms are comparable from the accuracy point of view, neither dominating the other. The use of oblique splits in SECRET made a big difference in four datasets (Kin8nm 27%, Stock 24%, Tecator 35% and 3DSin 45%). These datasets usually have less noise and are complicated but smooth (so they offer more opportunities

for *intelligent* splits). At the same time the use of oblique splits resulted in significantly worse performance on two of the datasets (Baseball 13% and Fried 19%).

5.3 Experimental results: Scalability

We chose to use only synthetic datasets for scalability experiments since the sizes of the real datasets are too small. The learning time of both GUIDE and SECRET is mostly dependent on the size of the training set and on the number of attributes, as is confirmed by some other experiments we are not reporting here. As in the case of accuracy experiments, we set the minimum number of data-points in a node to be considered for further splits to 1% of the size of the training set. We measured only the time to grow the trees, ignoring the time necessary for pruning and testing. The reason for this is the fact that pruning and testing can be implemented efficiently and for large datasets do not make a significant contribution to the running time. For GUIDE we report running times for both exhaustive search and sample split point (only 1% of the points available in a node are considered as possible split points), denoted by GUIDE(S).

Results of experiments with the 3DSin dataset and Fried dataset are depicted in Figures 5 and 6 respectively. A number of observations are apparent from these two sets of results: (1) the performance of the two versions of SECRET (with and without oblique splits) is virtually indistinguishable, (2) the running time of both versions of GUIDE is quadratic in size for large datasets, (3) as the number of attributes went up from 3 (3DSin) to 11 (Fried) the computation time for GUIDE(S), SECRET and SECRET(O) went up about 3.5 times but went slightly down for GUIDE, and (4) for large datasets (256000 tuples) SECRET is two orders of magnitude faster than GUIDE and one order of magnitude faster than GUIDE(S). It is also worth pointing out that, for SECRET, most of the time is spent in the EM algorithm. If used, sampling would not decrease the precision of EM much and at the same time would considerably decrease the computation time. For this reason the comparison with GUIDE(S) is not fair, nevertheless starting from medium sized datasets SECRET outperforms significantly the sampled version of GUIDE.

6. CONCLUSIONS

In this paper we introduced SECRET, a new linear regression tree construction algorithm designed to overcome the scalability problems of previous approaches. SECRET's scalability is based on a novel approach to tree construction that reduces the regression problem to a classification problem using the EM clustering algorithm. Our experiments show the efficacy and quality of SECRET.

We believe that this algorithm is only the first step towards a family of truly scalable regression tree algorithms. We see two immediate extensions. First, the pruning of classification trees has significant impact on their accuracy; we plan to investigate this problem for regression trees. Second, we will look into other possible approaches of split variable and split point selection.

Acknowledgements. This work was sponsored by NSF grants IIS-0121175 and IIS-0084762, by the Cornell Information Assurance Institute, and by generous gifts from Microsoft and Intel.

7. REFERENCES

- [1] W. P. Alexander and S. D. Grimshaw. Treed regression. *Journal of Computational and Graphical Statistics*, (5):156–175, 1996.
- [2] P. S. Bradley, U. M. Fayyad, and C. Reina. Scaling clustering algorithms to large databases. In *Knowledge Discovery and Data Mining*, pages 9–15, 1998.
- [3] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth, Belmont, 1984.
- [4] P. Chaudhuri, M.-C. Huang, W.-Y. Loh, and R. Yao. Piecewise-polynomial regression trees. *Statistica Sinica*, 4:143–167, 1994.
- [5] J. H. Friedman. Multivariate adaptive regression splines. *The Annals of Statistics*, 19:1–141 (with discussion), 1991.
- [6] K. Fukunaga. *Introduction to Statistical Pattern Recognition, Second edition*. Academic Press, 1990.
- [7] J. Gehrke, R. Ramakrishnan, and V. Ganti. Rainforest – a framework for fast decision tree construction of large datasets. In *Proceedings of the 24th International Conference on Very Large Databases*, pages 416–427. Morgan Kaufmann, August 1998.
- [8] G. H. Golub and C. F. V. Loan. *Matrix Computations*. Johns Hopkins, 1996.
- [9] A. Karalic. Linear regression in regression tree leaves. In *International School for Synthesis of Expert Knowledge, Bled, Slovenia*, 1992.
- [10] K.-C. Li, H.-H. Lue, and C.-H. Chen. Interactive tree-structured regression via principal hessian directions. *Journal of the American Statistical Association*, (95):547–560, 2000.
- [11] W.-Y. Loh. Regression trees with unbiased variable selection and interaction detection. *Statistica Sinica*, 2002. in press.
- [12] W.-Y. Loh and Y.-S. Shih. Split selection methods for classification trees. *Statistica Sinica*, 7(4), October 1997.
- [13] S. K. Murthy. Automatic construction of decision trees from data: A multi-disciplinary survey. *Data Mining and Knowledge Discovery*, 1997.
- [14] J. R. Quinlan. Learning with Continuous Classes. In *5th Australian Joint Conference on Artificial Intelligence*, pages 343–348, 1992.
- [15] J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufman, 1993.
- [16] L. Torgo. Functional models for regression tree leaves. In *Proc. 14th International Conference on Machine Learning*, pages 385–393. Morgan Kaufmann, 1997.
- [17] L. Torgo. Kernel regression trees. In *European Conference on Machine Learning*, 1997. Poster paper.
- [18] L. Torgo. A comparative study of reliable error estimators for pruning regression trees. *Iberoamerican Conf. on Artificial Intelligence*. Springer-Verlag, 1998.

	Constant Regressors			Linear Regressors		
	GUIDE	SECRET	SECRET(O)	GUIDE	SECRET	SECRET(O)
Abalone	5.32±0.05	5.50±0.10	5.41±0.10	<i>4.63±0.04</i>	4.67±0.04	4.76±0.05
Baseball	0.224±0.009	0.200±0.008	0.289±0.012	0.173±0.005	0.243±0.011	0.280±0.009
Boston	23.34±0.72	28.00±0.92	30.91±0.94	40.63±6.63	24.01±0.69	26.11±0.66
Kin8nm	0.0419±0.0002	0.0437±0.0002	0.0301±0.0003	0.0235±0.0002	0.0222±0.0002	0.0162±0.0001
Mpg	12.94±0.33	30.09±2.28	26.26±2.45	34.92±21.92	15.88±0.68	16.76±0.74
Mumps	1.34±0.02	1.59±0.02	1.56±0.02	1.02±0.02	1.23±0.02	1.32±0.04
Stock	2.23±0.06	2.20±0.06	2.18±0.07	1.49±0.09	1.35±0.05	1.03±0.03
TA	0.74±0.02	0.69±0.01	<i>0.69±0.01</i>	0.81±0.04	0.72±0.01	0.79±0.08
Tecator	57.59±2.40	49.72±1.72	28.21±1.75	13.46±0.72	12.08±0.53	7.80±0.53
3DSin	0.1435±0.0020	0.4110±0.0006	0.2864±0.0077	0.0448±0.0018	0.0384±0.0026	0.0209±0.0004
Cart	1.506±0.005	1.171±0.001	N/A	N/A	N/A	N/A
Fried	7.29±0.01	7.45±0.01	6.43±0.03	1.21±0.00	1.26±0.01	1.50±0.01

Table 2: Accuracy on real (upper part) and synthetic (lower part) datasets of GUIDE and SECRET. In parenthesis we indicate O for oblique splits. The winner is in bold face if it is statistically significant and in italics otherwise.

Size	GUIDE	GUIDE(S)	SECRET	SECRET(O)
250	0.07	0.05	0.21	0.21
500	0.13	0.07	0.33	0.34
1000	0.30	0.12	0.55	0.58
2000	0.94	0.24	1.08	1.12
4000	3.28	0.66	2.11	2.07
8000	12.58	2.40	4.07	4.12
16000	48.93	9.48	8.16	8.37
32000	264.50	43.25	16.71	16.19
64000	1389.88	184.50	35.62	35.91
128000	6369.94	708.73	73.35	71.67
256000	25224.02	2637.94	129.95	131.70

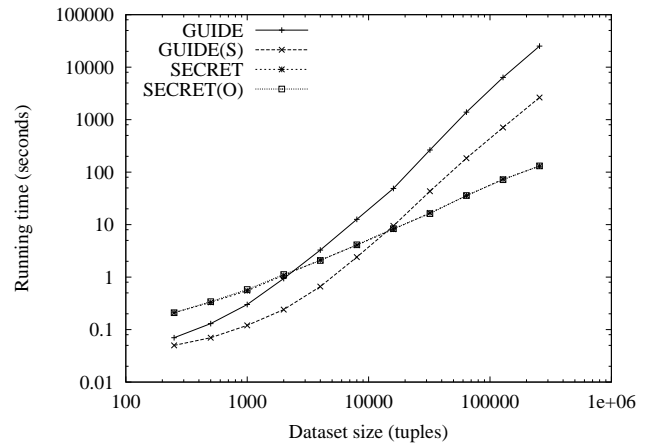


Figure 5: Running time (in seconds) of GUIDE, GUIDE with 0.01 of point as split points, SECRET and SECRET with oblique splits for synthetic dataset 3DSin (3 continuous attributes).

Size	GUIDE	GUIDE(S)	SECRET	SECRET(O)
250	0.09	0.07	0.47	0.43
500	0.17	0.14	0.87	0.92
1000	0.36	0.28	1.85	1.83
2000	1.12	0.80	3.58	3.69
4000	2.90	2.38	7.33	7.36
8000	10.46	8.43	13.77	14.05
16000	42.16	33.09	27.80	28.68
32000	194.63	123.63	56.87	58.01
64000	1082.70	533.16	122.26	124.60
128000	4464.88	1937.94	223.42	222.75
256000	18052.16	8434.33	460.12	470.68

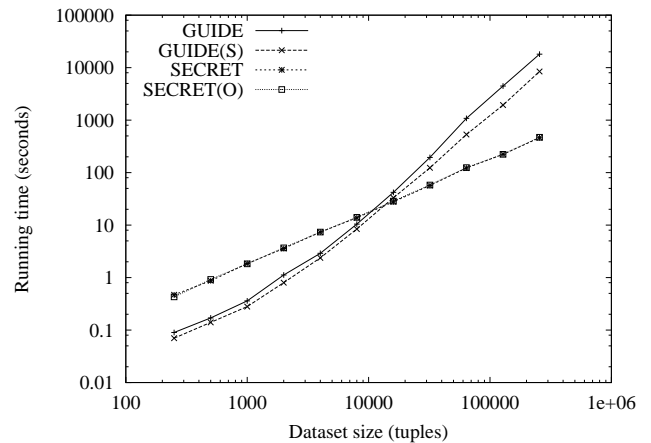


Figure 6: Running time (in seconds) of GUIDE, GUIDE with 0.01 of point as split points, SECRET and SECRET with oblique splits for synthetic dataset Fried (11 continuous attributes).