

Early Prediction of Coronary Artery Calcification Levels Using Statistical Relational Learning

Sriraam Natarajan, Kristian Kersting*, Saket Joshi⁺,
Santiago Saldana, Edward Ip, David R. Jacobs, Jr**, Jeffrey Carr

Wake Forest University School of Medicine, USA

*Fraunhofer IAIS, Germany

⁺ Oregon State University, USA

**University of Minnesota, USA

Abstract

Coronary heart disease (CHD) is a major cause of death worldwide. Although a multitude of cardiovascular risks factors have been identified, CHD most likely reflects actually complex interactions of these factors even over time. Today's datasets from longitudinal studies offer great promise to uncover these interactions but also pose enormous analytical problems due to typically large amount of both discrete and continuous measurements and risk factors with potential long-range interactions over time. Our investigation demonstrates that a statistical relational analysis of longitudinal data can easily uncover complex interactions of risks factors and actually predict future coronary artery calcified (CAC) plaque levels — an indicator of the amount of CHD present in an individual — significantly better than traditional non-relational machine learning approaches. The uncovered long-range interactions between risk factors conform to existing clinical knowledge and are successful in identifying risk factors at early adult stage. This makes it possible to design patient-specific treatments in young adults to mitigate the risk later.

1. Introduction

Heart disease and stroke – cardiovascular diseases, generally – encumber society with enormous costs. According to the World Heart Federation ¹, cardiovascular disease costs the European Union €169 billion in 2003 and the USA about €310.23 billion in direct and indirect annual costs.

One major cardiovascular disease is coronary heart disease (CHD). It is reported to be a major cause of morbidity and death in adults through heart attacks or acute myocardial infarctions (AMI). CHD is a condition which includes plaque build up inside the coronary arteries, i.e., atherosclerosis. Atherosclerosis is a disease process that begins in childhood, eventually resulting in clinical events later in life. The factors that determine development and progression of coronary artery disease are in large part established; however, the causes are very closely related with risk factors present in youth. Early detection of risks will help in designing effective treatments targeted at youth in order to prevent cardiovascular events in adulthood and to dramatically reduce the costs associated with cardiovascular diseases.

Our major contribution is to demonstrate the impact of machine learning on CHD research. We show that relationships between the measured risk factors and the development of advanced CAD lesions and overall plaque burden can be automatically extracted and understood. As the cohort ages and sufficient clinical events occur, this work will allow us to apply these methods to clinical events such as AMI and heart failure. Specifically, we propose to use the longitudinal

¹See <http://www.world-heart-federation.org/cardiovascular-health/global-facts-map/economic-impact/> and references in there.

data collected from the Coronary Artery Risk Developments in Young Adults (CARDIA)² study over different years to automatically estimate models using machine learning techniques for predicting the Coronary Artery Calcification (CAC) amounts, a measure of subclinical CAD, at year 20 given the measurements from the previous years. This longitudinal study began in 1985 – 86 and measured risk factors in different years (5,10,15,20) respectively. Several vital factors such as *body mass index (bmi)*, *cholesterol*-levels, *blood pressure* and *exercise* level are measured along with *family history*, *medical history*, *nutrient intake*, *obesity questions*, *pyschosocial*, *pulmonary function* etc. Using the predictions of the CAC levels we can predict cardiovascular events such as heart attacks. This in turn allows us to enable pro-active treatment planning for the high-risk patients i.e., identify young adult patients who are potentially at high-risk to cardiovascular events and design patient-specific treatments that would mitigate the risks.

We use Statistical Relational Learning (SRL) (Getoor & Taskar, 2007) algorithms for predicting CAC-levels in year 20 (corresponding to year 2005 when the patients were between 38 and 50 years old) given the measurements from all the previous years. SRL approaches, unlike what is traditionally done in statistical learning, seek to avoid explicit state enumeration as, in principle, is traditionally done in statistical learning through a symbolic representation of states. The advantage of these models is that they can succinctly represent probabilistic dependencies among the attributes of different related objects leading to a compact representation of learned models that allow for sharing of parameters between similar objects. Given that the CARDIA data is highly relational (multiple measurements are performed over examinations for each participant) and temporal, we use SRL methods to learn to predict CAC-levels. We use a de-identified version of the data set for methodological development.

More precisely, we use two kinds of SRL algorithms for this task – *Relational Probability Trees* (RPT) (Neville et al., 2003) and a more recently popular *Relational Functional Gradient Boosting* (RFGB) (Kersting & Driessens, 2008; Natarajan et al., 2011) approach. RPTs upgrade the attribute-value representation used within classical classification trees. The RFGB approach, on the other hand, involves learning a set of regression trees, each of which represents a potential function. The functional gradient approach has been found to give state of the art results in many relational

problems and we employ the same for CAC prediction. We use a sub-set of measurements from the CARDIA data set and predict the CAC-levels. We compared the SRL algorithms against propositional machine learning algorithms and demonstrated the superiority of the SRL algorithms. The learned models were also verified by a domain expert and the results conform to known medical risks. The results also provided a few insights about the relationships between risk factors and age of the individual. Identifying risk factors such as cholesterol level in young adulthood has potential to enable both the physician and the subject to devise a personalized plan to optimize it. Keeping track of these risk factors in young adulthood will prevent serious cardio-vascular events in their late adulthood.

The introduction of this domain to the SRL community itself is the second major contribution of the present paper. So far, significant performance gains have been reported for classical machine learning applications such as entity resolution, link prediction, and social network analysis. Discovering and understanding relationships between the measured risk factors and the development of advance CAD lesions and overall plaque burden has not been studied so far. Indeed, there has been some work on using machine learning for CAC prediction. For instance, in the work by Sun et al. (Sun et al., 2008), the authors use SNP data in a sub-set of population and employ ensemble methods such as random forests and RuleFit to predict CAC levels. This approach mainly relies on genetic information which are harder to obtain than clinical measurements. Others have employed machine learning techniques for identifying the presence of coronary artery disease. For instance the work by Hung et al. (Hung et al., 1985) predicts CAD on men using a logistic regression method. Lewis et al. (Lewis et al., 2006) use multivariate logistic regression models to identify whether chronic exposure to everyday discrimination affects the CAC levels in a certain section of the society. As far as we are aware, none of these methods take into account the clinical measurements from such a rich data set such as CARDIA spanning over 20 years to predict CAC-levels. Our current work is the first attempt at using ML techniques for this very significant problem.

2. Methodology

Before explaining how to adapt the CARDIA data to the relational setting, we will justify and detail our relational methodology.

²<http://www.cardia.dopm.uab.edu/>

2.1. The Need for Relational Models

Are relational models really beneficial? Could we also use propositional models? As we show, relational approaches are able to comprehensively outperform standard machine learning and data mining approaches. Beyond this, there are several justifications for adopting statistical relational analyses. First, the data consists of several diverse features (e.g., demographics, psychosocial, family history, dietary habits) that interact with each other in many complex ways making it *relational*. Second, the data was collected as part of a longitudinal study, i.e., over many different time periods such as 0, 5, 10, years etc., making it *temporal*. Third, like most data sets from biomedical applications, it contains missing values i.e., all data are not collected for all individuals. Fourth, the nature of SRL algorithms allow for more complex interactions between features. Finally, the learned models can be generalized across different sub-groups of participants and across different studies themselves. This data poses the following challenges for SRL methods creating the need for some preprocessing.

(1) Since the data are longitudinal, there are multiple measurements of many of the risk factors over different years of study. Hence time has to be incorporated into the model. To do so, the features are treated as fluents with time being the last argument. For instance, $weight(X, W, T)$ would refer to person X having weight W at time T . (2) CAC-levels of the participants are negligible (and often actually unobserved) in early years. This prevents us from using standard Dynamic Bayesian Network or HMMs; the values are nearly always zero in the initial years, being non-zero only in 10% at year 15 and 18% at year 20. (3) The input data consists of mainly continuous values. SRL methods use predicate logic where attributes are binary. In the case of features such as *cholesterol* level, *ldl*, *bmi*, we discretized them into bins based on domain knowledge. This is one of the key lessons learned: *using the domain expert’s knowledge (for discretization in our case) makes it possible to learn very highly predictive models in real problems.* (4) The cohort decreased over the years. There were a number of participants who did not appear for a certain number of years and returned for others. We did not try to normalize the data set by removing all the missing participants or replacing them with the most commonly observed value. Instead, we allowed the values to be missing. The only case where we dropped the participants from the data base was when they were not present in year 20 where we predict the CAC-levels. This is to say that we are not considering the problem to be a semi-supervised learning problem but

treat it as a strictly supervised learning one. (5) Recall the goal of the study is to identify the effect of the factors in early adulthood on cardiovascular risks in middle-aged adults. The algorithm should be allowed to search through all the risk factors in all the years for predicting CAC-levels. This implies that the data must not be altered or tailored in any form. In this work, we did not make any modifications to the data except for the discretizations mentioned earlier. As we show, our methods are very successful in identifying long-range correlations. One of the biggest lessons learned from this work is that risk factors between the ages of 25 through 40 are very significant for CAC-level prediction at age 38 to 50.

However, which SRL approach should we use?

2.2. Relational Gradient Boosting

One of the most important challenges in SRL is learning the structure of the models, i.e., the weighted relational rules. This problem has received much attention lately. Most approaches (Kok & Domingos, 2010) follow a traditional greedy hill-climbing search: first obtain the candidate rules/clauses, score them, i.e., learn the weights, and select the best candidate. The temporal nature of our task at hand makes it difficult to use these approaches. Therefore, we use a boosting approach based on functional gradients recently proposed that learns the structure and parameters simultaneously (Natarajan et al., 2011). It was proven successful in several classical SRL domains and achieves state-of-the art performances. Also, it easily allows — as we will show — to account for the temporal aspects of CAC-level prediction.

Functional gradient methods have been used previously to train conditional random fields (CRF) (Dietterich et al.(2004)) and their relational extensions (TILDE-CRF) (Gutmann & Kersting, 2006). Assume that the training examples are of the form (\mathbf{x}_i, y_i) for $i = 1, \dots, N$ and $y_i \in \{1, \dots, K\}$. We use \mathbf{x} to denote the vector of features. The goal is to fit a model $P(y|\mathbf{x}) \propto e^{\psi(y, \mathbf{x})}$. Dietterich et al. used an approach to train the potential functions based on Friedman’s(2001) gradient-tree boosting algorithm where the potential functions are represented by sums of regression trees that are grown stage-wise. Since the stage-wise growth of these regression trees are similar to the Adaboost algorithm (Freund & Schapire, 1996), this is called as *gradient-tree boosting*. More formally, functional gradient ascent starts with an initial potential ψ_0 and iteratively adds gradients Δ_i . After m iterations, the potential is given by $\psi_m = \psi_0 + \Delta_1 + \dots + \Delta_m$. Here, Δ_m is the functional gradient at episode m and

is

$$\Delta_m = \eta_m \times E_{x,y}[\partial/\partial\psi_{m-1}\log P(y|x; \psi_{m-1})] \quad (1)$$

where η_m is the learning rate. Dietterich *et al.* suggested evaluating the gradient at every position in every training example and fitting a regression tree to these derived examples i.e., fit a regression tree h_m on the training examples $[(x_i, y_i), \Delta_m(y_i; x_i)]$. They point out that although the fitted function h_m is not exactly the same as the desired Δ_m , it will point in the same direction (assuming that there are enough training examples). So ascent in the direction of h_m will approximate the true functional gradient.

Let us denote the CAC-level as y and for ease of explanation assume that it is binary valued (i.e., present vs absent). Let us denote all the other variables measured over the different years as \mathbf{x} . Our aim is to learn $P(y|\mathbf{x})$ where, $P(y|\mathbf{x}) = e^{\psi(y;\mathbf{x})} / \sum_y e^{\psi(y;\mathbf{x})}$. Note that in the functional gradient presented in Equation 1, the expectation $E_{x,y}[\dots]$ cannot be computed as the joint distribution $P(\mathbf{x}, \mathbf{y})$ is unknown. Instead of computing the functional gradients over the potential function, they are instead computed for each training example i given as $\langle \mathbf{x}_i, y_i \rangle$. Now this set of local gradients form a set of training examples for the gradient at stage m . The main idea in the gradient-tree boosting is to fit a regression-tree on the training examples at each gradient step. In this work, we replace the propositional regression trees with relational regression trees.

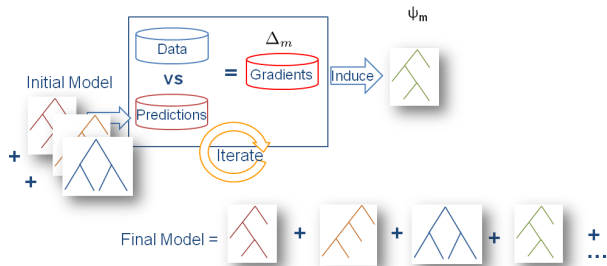


Figure 1. Relational Functional Gradient Boosting. This is similar to the standard FGB where trees are induced in stage-wise manner the key difference being that the trees are relational regression trees.

The functional gradient with respect to $\psi(y_i = 1; \mathbf{x}_i)$ of the likelihood for each example $\langle y_i, \mathbf{x}_i \rangle$ can be shown to be: $\frac{\partial \log P(y_i; \mathbf{x}_i)}{\partial \psi(y_i = 1; \mathbf{x}_i)} = I(y_i = 1; \mathbf{x}_i) - P(y_i = 1; \mathbf{x}_i)$, where I is the indicator function that is 1 if $y_i = 1$ and 0 otherwise. The expression is simply the adjustment required to match the predicted probability with the true label of the example. If the example is positive (i.e., if the participant has significant CAC-level

in year 20) and the predicted probability is less than 1, this gradient is positive indicating that the predicted probability should move towards 1. Conversely, if the example is negative and the predicted probability is greater than 0, the gradient is negative driving the value the other way.

We use *Relational Regression Trees* (RRTs) (Blokkeel & Raedt, 1998) to fit the gradient function for every training example. Each RRT can be viewed as defining several new feature combinations, one corresponding to each path from the root to a leaf. The resulting potential functions from all these different RRTs still have the form of a linear combination of features but the features can be quite complex.

This idea is illustrated in Figure 1. First a tree is learned from the training examples and this tree is used to determine the weights (i.e., functional gradients) of the examples for the next iteration (which in this case is the difference between the true probability of being true and the predicted probability). Once the examples are weighted, a new tree is induced from the examples. The trees are then considered together and the regression values are added when weighing the examples and the process is repeated.

At a fairly high level, the learning of RRT proceeds as follows: The learning algorithm starts with an empty tree and repeatedly searches for the best test for a node according to some splitting criterion such as weighted variance. Next, the examples in the node are split into *success* and *failure* according to the test. For each split, the procedure is recursively applied further obtaining subtrees for the splits. We use weighted variance on the examples as the test criterion. In our method, we use a small depth limit (of at most 3) to terminate the search. In the leaves, the average regression values are computed.

The *key* idea underlying the present work is to represent the distribution over CAC-levels as a set of RRTs on the features. When learning to predict the CAC-levels in year 20, we use the data collected from all the previous years. We ignore the CAC-levels that are present for some individuals at year 15 since we are interested in planning preventive treatments in early adulthood based on other risk factors. We bear in mind that CAC rarely regresses from present to absent or from a higher level to a lower level.

These trees are learned such that at each iteration the new set of RRTs aim to maximize the likelihood of the distributions w.r.t ψ . When computing $P(cac(X)|\mathbf{f}(X))$ for a particular patient X , given the feature set \mathbf{f} , each branch in each tree is considered

to determine the branches that are satisfied for that particular grounding (x) and their corresponding regression values are added to the potential ψ .

To investigate the usefulness of other relational learners, we also considered Relational Probability Trees (Neville et al., 2003). We modified the RPT learning algorithm to learn a regression tree similar to TILDE to predict positive examples and turn the regression values in the leaves into probabilities by exponentiating the regression value and normalizing them. We modified TILDE to automatically include aggregate functions such as count, mode, max, mean etc. while searching for the next node to add to the tree. Also, the regression tree learner can use conjunctions of predicates in the inner nodes as against a single test by the traditional RPT learner. This modification has been shown to have better performance than RPTs (Natarajan et al., 2011) and hence we employ this modified RPT learner in our experiments.

3. Adapting the CARDIA Data

The CARDIA Study examines the development and determinants of clinical and subclinical cardiovascular disease and its risk factors. It began in 1985 – 6 (Year 0) with a group of 5115 men and women whose age were between 18-30 years from 4 centers. The same participants were asked to participate in follow-up examinations during 87 – 88 (Year 2), 90 – 91 (Year 5), 92 – 93 (Year 7), 95 – 96 (Year 10), 2000 – 2001 (Year 15), and 05 – 06 (Year 20). A majority of the group has been examined at each of the follow-up examinations (90%, 86%, 81%, 79%, 74%, and 72%, respectively). Data has been collected on a variety of factors believed to be related to heart disease. This rich data set provides a valuable opportunity to identify risk factors in young adults that could cause serious cardiovascular issues in their adulthood. This in turn, will allow physicians to develop patient-specific preventive treatments that can improve the quality of life in later years.

We used known risk factors such as *age*, *sex*, *cholesterol*, *bmi*, *glucose*, *hdl level*, *ldl level*, *exercise*, *trig level*, *systolic bp* and *diastolic bp* that are measured between years 0 and 20 over the patients. Our goal is to predict if the CAC-levels of the patients are above 0 for year 20 given the above mentioned factors. Any CAC-level over 0 indicates the presence of advanced coronary atheroma and elevated risk for future CHD and need to be monitored. So, we are in a binary classification setting of predicting 0 vs non-0 CAC levels. In our data set, most of the population had CAC-level of 0 (2981 out of 3043 subjects) in year 20. Hence

Feature	Thresholds
cholesterol	70, 100, 150, 200, 250, 300, 400
dbp	0, 30, 50, 70, 90, 100, 150
glucose	0, 50, 100, 200, 300, 400
hdl	10, 30, 50, 70, 100,120,200
ldl	0,50,100,150,200,400
trig	0,25,50,100,300,1000,3000
bmi	0,16,18.5,25,30,35,40,100

Classifier	Parameters
J48	C 0.25 M 2
SVM	C 1.0, L 0.01, P 1E12, N 0, V 1, W 1 Poly Kernel
AdaBoost	P 100, S 1, l 10
Bagging	P 100, S 1, l 10, M 2, V 0.001, N 3, S 1, L -1
Logistic	R 1.0E-8, M -1

Table 1. (Top) Domain expert’s discretization of some of the input features. (Bottom) Parameters of the propositional classifiers

there is a huge skew in the data set where there is a very small number of positive examples (less than 20% of subjects had significant CAC-levels). It must be mentioned that the data for year 25 is still being compiled and hence our prediction models have been designed for year 20 data.

We converted the data set into predicate logic, see e.g. (De Raedt, 2008) for an introduction. The first argument of every predicate is the ID of the person and the last argument is the year of measurement. It is possible for our algorithm to search at the level of the variables or ground the variable to a constant while searching for the next predicate to add to the tree. For example, we could use some values such as “never smoked”, “quit smoking” etc. directly in the learned model and in other cases, use variables in the node. This is yet another reason for employing a relational learning algorithm.

The risk factors, however, are continuous variables. For instance, *ldl*, *hdl*, *glucose*, *bmi*, *dbp*, *sdp* etc. all take real numbered values with different ranges. While many methods exist that can discretize the data and/or directly operate on the continuous data and automatically discretizing the features based on the data is preferred, some of these risk factors have been analyzed by the medical community for decades and the thresholds have been identified. For instance a *bmi* of less than 16 is *severely underweight*, greater than 40 is *extremely obese* etc. Hence, we used inputs from a domain expert to identify the thresholds for discretizing the numeric features and these are presented in Table 1.

We also included the difference between the two successive measurements as input features. This repre-

sents the “change” in risk factor for the subject. For the boosting algorithm (RFGB), we used the preset number of parameters of trees, namely 20. The tree-depth of the trees were set to 3 and hence we preferred a reasonably large set of small trees. As mentioned above, we allowed the algorithm to construct the aggregators on the fly. We compare against learning a single tree(RPT) of depth 10. This is due to the fact that every path from root to leaf indicates an interaction between the risk factors and our domain expert indicated that 10 should be the upper limit of the interactions. We also compared our algorithms against the standard ML algorithms using the weka package. Hence, we propositionalized our features by creating one feature for every measurement at every year. We included the change (difference between measurements in successive years) features for the propositional data set as well. The default parameters for the propositional classifiers are presented in Table 1. We tried a few different parameter settings for the algorithms in Weka and report the best results.

We performed 5-fold cross-validation. Since there is an unequal distribution of the pos and neg examples measuring accuracy can be misleading, and hence we present the area under PR and ROC curves in the next section.

4. Predicting CAC Levels

Comparison with Propositional Learners: We present the results of learning to predict CAC-levels using our algorithms and the standard ML techniques. A full test set has a very large skew towards one class. Hence in the test set (since accuracies in relational data can be very inflated), we sampled twice the number of negatives as positives. Recall that the positive class would mean that the CAC-level of the subject in year 20 is significant (i.e., greater than 0). Table 2 compares the results of our techniques – boosting (*RFGB*) and *RPT* — against decision-trees (*J48*), *SVM*, *AdaBoost*, *Bagging*, *Logistic Regression* (LR) and Naive Bayes (*NB*).

SRL approaches output probabilities instead of labels. Hence, to measure accuracy, we computed the mean probability of the true class i.e., $\frac{1}{n} \sum_i P(y_i | \mathbf{x}_i)$ where y_i is the true class label for subject i . Accuracies in this data set do not really reflect the true performance of the algorithms and the accuracies of all the algorithms are very close to each other and there is no significance in the results (every accuracy was around 0.667). A key property of most relational data sets is the number of negative examples. This is also seen in our data set since most CAC-levels are zero and hence the number

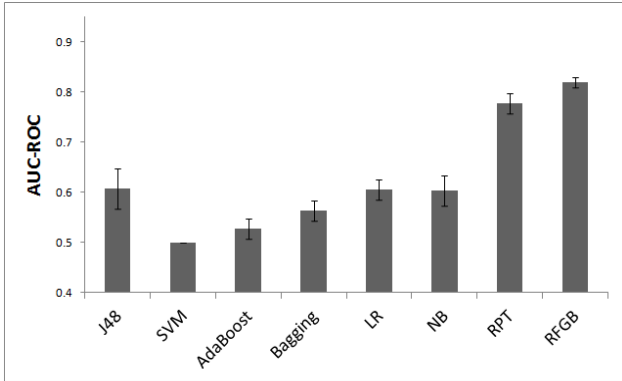


Figure 2. AUC ROC values for the different algorithms.

of negatives can be order of magnitude more than the number of positives. In these cases, simply measuring accuracy or conditional loglikelihood (CLL) over the entire data set can be misleading. It can be shown easily that predicting all the examples as the majority class (when the number of examples in one class are far greater than the other) can have a very good CLL value, but a very low AUC-ROC or AUC-PR value (nearly 0).

The AUC-ROC results presented in Figure 2 clearly show that the SRL approaches dominate the propositional ones. Most of the standard algorithms classify nearly all the examples as negative and hence presenting accuracies can be misleading. We chose to present AUC-ROC instead. SVM and AdaBoost classify all examples as negative while Bagging, LR, Naive Bayes and J48 classify a very small number of examples (nearly 5% of positive examples correctly). In contrast, the SRL approaches have a smoother classification performance and hence have a higher AUC-ROC with RFGB having the best ROC.

Detailed Analysis of SRL Algorithms: We present the Precision Recall curves for the SRL algorithms in Figure 3.a. We did not include the other algorithms since their PR values were very low. The boosting approach has a better performance particularly in the medically-relevant high recall region. Evaluating precision at high recalls assesses an algorithm’s ability to predict while disallowing many false negatives, which is the critical component to a good screening tool. In the case of predicting CAC levels, a false negative means categorizing a patient as “low-risk” who goes on to have a heart attack, a costly outcome we wish to avoid. It is clear that RFGB has a better precision in high recall regions.

The effect of the number of trees on the performance

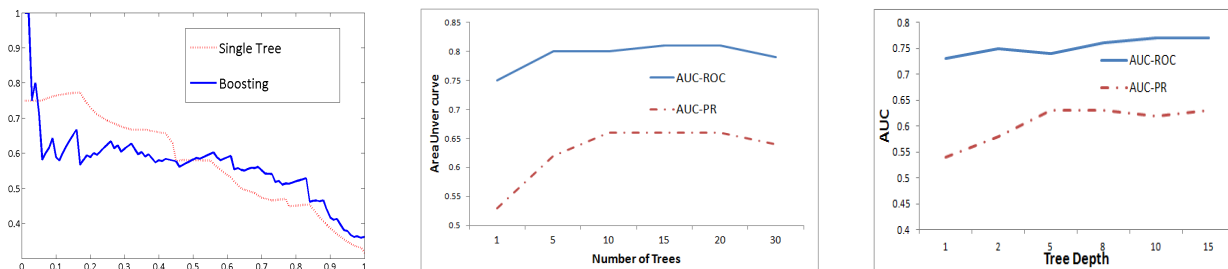


Figure 3. (a.) Precision-Recall curves comparing SRL algorithms. (b.) Effect of number of trees in performance of RFGB. (c.) Effect of the depth of the tree in the performance of a single RRT

of RFGB is presented in Figure 3.b. We have presented both the AUC-ROC and AUC-PR values as a function of the number of trees. As the number of trees increase, there is an increase in the performance of RFGB. Also, it can be noted that beyond a certain number of trees (in our case 10), there is not a significant increase in the performance. When we look at the trees closer, it appears that with larger number of trees (say 30), the last few trees are picking up random correlations in the data (though the regression values in the leaves are quite low). Figure 3.c presents the effect of the depth of the tree when learning a single tree (i.e., RPT). It appears that the performance of the algorithm stabilizes around a depth of 5. Increasing beyond 5 does not have a statistically significant impact on the performance showing that interactions between 5 risk factors is sufficient to predict the CAC-levels.

Analysis of Learned Relational Models: Figure 4 illustrates a part of one tree learned. The first argument a of every predicate is the subject’s ID and the last argument of every predicate (except sex) indicates the year. The left branch out of every node is the *true* branch, the right branch the *false* branch. The leaves indicate the probability of cac-level (say p) being greater than 0. We use $_{bw}$ in predicates to indicate that the value of a certain variable is between two values. For instance, $ldl_{bw}(a, b, 0, 100, 10)$ indicates that the ldl level of the person a is b and is between 0 and 100 in year 10.

We are not presenting the entire tree and indicate the missing branches by dots. As one can see, the first test checks the sex of the subject. So the first right probability (box with value 0.05) indicates that if the person is a female and she does not smoke in year 5, $smoke(a, No, 5)$ where No is the value of smoke, then $p = 0.05$. However, if she smokes in year 5, but has a low ldl cholesterol level in year 7 and does not have high blood pressure in year 7, then $p = 0.2$. This indicates that even if she is a smoker, if her cholesterol level and blood pressure are under control, she has

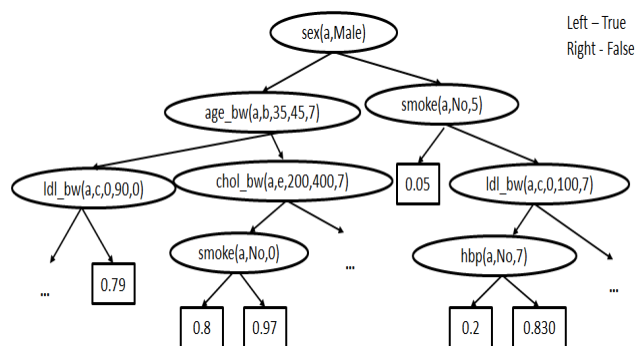


Figure 4. Learned Tree for predicting CAC-level greater than 0. The leaves indicates $P(cac(a)) > 0$.

low risk. However, if the subject is a male, things are significantly different. If he is in middle age in year 7 and has a high ldl level, $p = 0.79$. If he is young and has a high cholesterol level in year 7 (under the age of 35) and he smokes in year 0, then $p = 0.97$. When boosting was employed, the shorter trees had very similar structures. The first tree in fact splits the population based on sex and then on smoking and cholesterol respectively.

Prediction Based on Early Adulthood Data only: We repeated the experiment with one major change. Instead of considering all the risk factors at all years, we considered the measurements only till year 10 i.e., only the risk factors from young adulthood. The goal is again to predict the CAC-level in year 20. The average AUC-ROC values are 0.779 ± 0.01 and are not significantly different from the ones learned using the entire data set. This confirms our hypothesis that the risk factors in young age are responsible for the cardio-vascular risks in older age thus validating our claim that changing the lifestyle in younger years can lead to a healthier life in older years. This validates the observations made by Loria et al. (Loria et al., 2007) where individual correlations between risk factors at different years and CAC-level at year 15 are measured to show that year 0 risk factors are as informative as

later years.

4.1. Assessment of the Results

The results were verified by our *radiologist*, and are very interesting from a medical perspective for several reasons: First, as our last set of experiments show, the risk of CAC levels in later years is mostly indicated by risk factors in early years (ages 25 through 45). This is very significant from the point of view of the CARDIA study since the goal is to identify risk factors in early adult stage so as to prevent cardio-vascular issues in late adulthood. Second, the tree conforms to some known or hypothesized facts. For instance, it is believed that females are less prone to cardio-vascular issues than males. The tree identifies sex as the top splitting criterion. Similarly, in men, it is believed that the ldl and hdl levels are very predictive and the tree confirms this. It can be seen that all the cases presented in the earlier paragraph make sense w.r.t clinical knowledge. Third, the tree also identifies complex interaction between risk factors at different years. For instance – (i) smoking in year 5 interacts with cholesterol level in later years in the case of females, and (ii) the triglyceride level in year 5 interacts with the cholesterol level in year 7 for males. Finally, the structure of the tree enables the physician to identify treatable risk factors and plan preventive treatments leading to a healthier lifestyle.

5. Conclusion

Coronary heart disease (CHD) kills millions of people each year. The broadening availability of longitudinal studies and electronic medical records presents both opportunities and challenges to apply AI techniques to improve CHD treatment. We discussed them for the important problem of identifying risk factors in young adults that can lead to cardio-vascular issues in their late adulthood. We addressed the specific problem of uncovering interactions among risk factors and of using them for predicting CAC levels in adults given the risk factor measurements of their youth. Our experimental results indicate that statistical relational models are superior to non-relational ones. More importantly, our learned models were verified by the domain expert and the results conform to the current clinical knowledge.

Motivated by the initial success of our work, we plan to pursue research in several different directions. First, we plan to include all the collected features for training the models. This will allow one to identify complex relationships between different types of features such as demographics and psychosocial etc. Second, while the boosted set of trees have high predictive accuracy,

they may not necessarily be easy to interpret by physicians. Hence our goal is to convert the set of trees into a single tree. Finally, we are focussing on understanding the development of risks over time (i.e., explicit temporal modeling).

References

- Blockeel, H. and Raedt, L. De. Top-down induction of first-order logical decision trees. *Artificial Intelligence*, 101: 285–297, 1998.
- De Raedt, L. *Logical and Relational Learning*. Springer, 2008.
- Freund, Y. and Schapire, R. Experiments with a new boosting algorithm. In *ICML*, 1996.
- Getoor, L. and Taskar, B. *Introduction to Statistical Relational Learning*. MIT Press, 2007.
- Gutmann, B. and Kersting, K. TildeCRF: Conditional Random Fields for Logical sequences. In *ECML*, 2006.
- Hung, J., Chaitman, B.R., Lam, J., Lesperance, J., Dupras, G., Fines, P., Cherakaoui, O., Robert, P., and Bourassa, M.G. A logistic regression analysis of multiple noninvasive tests for the prediction of the presence and extent of coronary artery disease in men. *American Heart Journal*, 110:460–469, 1985.
- Kersting, K. and Driessens, K. Non-parametric policy gradients: A unified treatment of propositional and relational domains. In *ICML*, 2008.
- Kok, S. and Domingos, P. Learning Markov Logic networks using structural motifs. In *ICML*, 2010.
- Lewis, T., Everson-Rose, S., Powell, L., Matthews, K., Brown, C., Karavolos, K., Sutton-Tyrell, K., Jacobs, E., and Wesley, D. Chronic exposure to everyday discrimination and coronary artery calcification in african-american women: The swan heart study. *Psychosomatic Medicine*, 68:362–368, 2006.
- Loria, C.M., Liu, K., and et al, C. E. Lewis. Early adult risk factor levels and subsequent coronary artery calcification - the cardia study. *Journal of the American College of Cardiology*, 49, 2007.
- Natarajan, S., Khot, T., Kersting, K., Guttmann, B., and Shavlik, J. Gradient-based boosting for statistical relational learning: The relational dependency network case. *Machine Learning*, 2011.
- Neville, J., Jensen, D., Friedland, L., and Hay, M. Learning Relational Probability trees. In *KDD*, 2003.
- Sun, Y., Beilak, L., Peyser, P., Turner, T., Sheedy, P., Boerwinkle, E., and Kardia, S. Application of machine learning algorithms to predict coronary artery calcification with a sibship-based design. *Epidemiology*, 32: 350–360, 2008.