



Published in final edited form as:

Appl Soft Comput. 2007 January ; 7(1): 471–479. doi:10.1016/j.asoc.2006.01.013.

Genetic Programming Neural Networks: A Powerful Bioinformatics Tool for Human Genetics

Marylyn D Ritchie¹, Alison A. Motsinger¹, William S Bush¹, Christopher S Coffey², and Jason H Moore³

Marylyn D Ritchie: ritchie@chgr.mc.vanderbilt.edu; Alison A. Motsinger: motsinger@chgr.mc.vanderbilt.edu; William S Bush: wbush@chgr.mc.vanderbilt.edu; Christopher S Coffey: CCoffey@ms.soph.uab.edu; Jason H Moore: Jason.H.Moore@dartmouth.edu

¹Center for Human Genetics Research, Department of Molecular Physiology and Biophysics, Vanderbilt University, 519 Light Hall, Nashville, TN 37232

²Department of Biostatistics, University of Alabama at Birmingham, Ryals Public Health Bldg., Rm. 327M, Birmingham, AL, 35294

³Computational Genetics Laboratory, Department of Genetics 706 Ruben Bldg., HB7937, Dartmouth-Hitchcock Medical Center, One Medical Center Dr. Lebanon, NH 03756

Abstract

The identification of genes that influence the risk of common, complex disease primarily through interactions with other genes and environmental factors remains a statistical and computational challenge in genetic epidemiology. This challenge is partly due to the limitations of parametric statistical methods for detecting genetic effects that are dependent solely or partially on interactions. We have previously introduced a genetic programming neural network (GPNN) as a method for optimizing the architecture of a neural network to improve the identification of genetic and gene-environment combinations associated with disease risk. Previous empirical studies suggest GPNN has excellent power for identifying gene-gene and gene-environment interactions. The goal of this study was to compare the power of GPNN to stepwise logistic regression (SLR) and classification and regression trees (CART) for identifying gene-gene and gene-environment interactions. SLR and CART are standard methods of analysis for genetic association studies. Using simulated data, we show that GPNN has higher power to identify gene-gene and gene-environment interactions than SLR and CART. These results indicate that GPNN may be a useful pattern recognition approach for detecting gene-gene and gene-environment interactions in studies of human disease.

1 Introduction

One goal of genetic epidemiology is to identify genes associated with common, complex multifactorial diseases. Success in achieving this goal will depend on a research strategy that recognizes and addresses the importance of interactions among multiple genetic and environmental factors in the etiology of diseases such as essential hypertension [8,15]. One traditional approach to modeling the relationship between discrete predictors such as genotypes and discrete clinical outcomes is logistic regression [7]. Logistic regression is a parametric statistical approach for relating one or more independent or explanatory variables (e.g. genotypes) to a dependent or outcome variable (e.g. disease status) that follows a binomial distribution. However, as reviewed by Moore and Williams [15], the number of possible interaction terms grows exponentially as each additional main effect is included in the logistic regression model. Thus, logistic regression is limited in its ability to deal with interactions involving many factors. Having too many independent variables in relation to the number of

observed outcome events is a well-recognized problem [3,17] and is an example of the curse of dimensionality [2].

In response to this limitation, Ritchie et al. [19] developed a genetic programming optimized neural network (GPNN). Neural networks (NN) have been utilized in genetic epidemiology, however, with little success. A potential weakness in the previous NN applications is the poor specification of NN architecture. GPNN was developed in an attempt to improve upon the trial-and-error process of choosing an optimal architecture for a pure feed-forward back propagation neural network. The GPNN optimizes the inputs from a larger pool of variables, the weights, and the connectivity of the network including the number of hidden layers and the number of nodes in the hidden layer. Thus, the algorithm attempts to generate optimal neural network architecture for a given data set. This is an advantage over the traditional back propagation NN in which the inputs and architecture are pre-specified and only the weights are optimized.

Although previous empirical studies suggest GPNN has excellent power for identifying gene-gene interactions, a comparison of GPNN with traditional statistical methods has not yet been performed. The goal of the present study was to compare the power of GPNN to that of stepwise logistic regression (SLR) and classification and regression trees (CART) for identifying gene-gene and gene-environment interactions using data simulated from a variety of interaction models. This study is motivated by the number of studies in human genetics where SLR and CART have been applied. We wanted to determine if GPNN is more powerful than the status quo in the field. We find that GPNN has higher power to detect gene-gene and gene-environment interactions than stepwise logistic regression and classification and regression trees. These results demonstrate that GPNN may be an important pattern recognition tool for future studies in genetic epidemiology.

2 Methods

2.1 A Genetic Programming Neural Network Approach

GPNN was developed to improve upon the trial-and-error process of choosing an optimal architecture for a pure feed-forward back propagation neural network (NN) [19]. Optimization of NN architecture using genetic programming (GP) was first proposed by Koza and Rice [9]. The goal of this approach is to use the evolutionary features of genetic programming to evolve the architecture of a NN. The use of binary expression trees allow for the flexibility of the GP to evolve a tree-like structure that adheres to the components of a NN. Figure 1 shows an example of a binary expression tree representation of a NN generated by GPNN. The GP is constrained such that it uses standard GP operators but retains the typical structure of a feed-forward NN. While GP could be implemented without constraints, the goal was to evolve NN since they were being explored as a tool for genetic epidemiology. Thus, we wanted to make an improvement to a method already being used. A set of rules is defined prior to network evolution to ensure that the GP tree maintains a structure that represents a NN. The rules used for this GPNN implementation are consistent with those described by Koza and Rice [9]. The flexibility of the GPNN allows optimal network architectures to be generated that consist of the appropriate inputs, connections, and weights for a given data set.

The GPNN method has been described in detail [19]. The steps of the GPNN method are shown in Figure 2 and described in brief as follows. First, GPNN has a set of parameters that must be initialized before beginning the evolution of NN models. These include an independent variable input set, a list of mathematical functions, a fitness function, and finally the operating parameters of the GP. These operating parameters include number of demes (or populations), population size, number of generations, reproduction rate, crossover rate, mutation rate, and migration [19]. Second, the data are divided into 10 equal parts for 10-fold cross-validation.

Here, we will train the GPNN on 9/10 of the data to develop a NN model. Later, we will test this model on the 1/10 of the data left out to evaluate the predictive ability of the model.

Third, training of the GPNN begins by generating an initial population of random solutions. Each solution is a binary expression tree representation of a NN, similar to that shown in Figure 1. Fourth, each GPNN is evaluated on the training set and its fitness recorded. Fifth, the best solutions are selected for crossover and reproduction using a fitness-proportionate selection technique, called roulette wheel selection, based on the classification error of the training data [10]. Classification error is defined as the proportion of individuals where the disease status was incorrectly specified. A predefined proportion of the best solutions will be directly copied (reproduced) into the new generation. Another proportion of the solutions will be used for crossover with other best solutions. The new generation, which is equal in size to the original population, begins the cycle again. This continues until some criterion is met at which point the GPNN stops. This criterion is either a classification error of zero or the maximum number of generations having been reached. A “best-so-far” solution is chosen after each generation. At the end of the GPNN evolution, the one “best-so-far” solution is selected as the optimal NN. Sixth, this best GPNN model is tested on the 1/10 of the data left out to estimate the prediction error of the model. Prediction error is a measure of the ability to predict disease status in the 1/10 of the data. Steps two through six are performed ten times with the same parameters settings, each time using a different 9/10 of the data for training and 1/10 of the data for testing.

The results of a GPNN analysis include 10 GPNN models, one for each split of the data. In addition, a classification error and prediction error is recorded for each of the models. A cross-validation consistency can be measured to determine those variables which have a strong signal in the gene-gene interaction model [12,14,18,19]. Cross-validation consistency is the number of times a particular combination of variables are present in the GPNN model out of the ten cross-validation data splits. Thus a high cross-validation consistency, ~10, would indicate a strong signal, whereas a low cross-validation consistency, ~1, would indicate a weak signal and a potentially false positive result. We select the best model as the combination of loci with the highest cross-validation consistency across the ten cross-validation intervals.

2.2 Stepwise Logistic Regression

A common traditional method for developing genetic models from association data is logistic regression. Logistic regression is a derivative of linear regression that fits a function to continuous or discrete independent variables based on a dichotomous dependent variable [7]. Logistic regression uses a transformation of the logistic distribution to develop a function based on the independent variables. The logistic distribution provides a flexible platform with a clinically meaningful interpretation. The distribution is transformed by the logit function, allowing traditional regression techniques to be applied. Using this formulation, the predicted values from regression are dichotomous with binomially distributed errors. Once a regression function has been derived using iterative fitting techniques, least squares is used to minimize errors, based on the binomial distribution and is given as a logarithmic transformation of the maximum likelihood, called a log-likelihood. Like linear regression, partial derivatives of the likelihood function are evaluated to minimize error.

Variable selection occurs by first developing a full model using all variables. Using this model's log-likelihood as a basis for comparison, variables are sequentially removed and significance is tested to determine if that variable contributed to the model. At each level, additional variables are re-entered into the model and evaluated for significant contributions. Using this stepwise selection, the most parsimonious model producing the best fit is selected as the functional model and the variables used are the functional loci in the final genetic model.

Logistic regression is a *de facto* standard for traditional association studies. Using independent variables to predict a dichotomous dependent variable, logistic regression by definition lacks the ability to characterize interactive effects. Only variables that contain an independent main effect will be included in the final model. To properly evaluate non-linear interactive effects, combinations of variables must be encoded as a single variable for inclusion in the analysis. Such an encoding scheme can be computationally expensive, depending on the number of variables used.

2.3 CART

Classification and Regression Trees (CART) is a technique employing binary decision trees to partition training data into smaller and smaller subsets [5]. Classification trees have a categorical outcome and regression trees have a continuous outcome, hence for genetic association studies, a classification tree is used. Because data is split into subsets based on independent variables, tree generation is a form of variable selection [22]. Each endpoint or leaf node represents a labeled partition of the original set of observations whose lineage can be traced by identifying the binary split conditions used to arrive at this state. Each node is assessed for its impurity – the number of individuals in each descendant node that will be misclassified by the proposed split. While several metrics for impurity have been developed, they are all variants of a one-step look ahead, which is a greedy method, prone to dwell in local optima rather than searching for a global optimum.

Cross-validation is an integral part of CART, and once a tree has been generated, the testing set is used to prune the tree. Pruning is the process of analyzing splits and removing pairs of leaf nodes (sharing a common antecedent) where the splitting did not significantly decrease impurity. Cost-complexity pruning is generally more efficient, removing entire sub-trees whose leaf nodes do not sufficiently optimize the tree. Cost-complexity pruning weighs both the size of the tree and the relative decrease in impurity over the sub-tree to optimize both parsimony and average impurity. Once the optimal tree is developed, the variables used in data splitting are the functional loci in the final genetic model.

CART has been hailed by some as an ideal method to identify gene-gene interactions in association studies. While CART analyzes the compounded effects of multiple variables, it fails to account for completely non-linear interactions. In the case of the XOR model, CART would not accurately partition the data according to the effects of this interaction. CART would first split on a single variable, partitioning the data into two subsets, and because it conditions on the effect of this single variable, it cannot detect the non-linear interactions of two or more independent variables that themselves lack main effects.

2.4 Data Simulation

The goal of the simulation was to generate data sets that exhibit gene-gene interactions for the purpose of evaluating the power of GPNN in comparison to the power of SLR and CART. We simulated a collection of models varying several conditions including number of interacting genes, allele frequency, and heritability. Heritability is defined in the broad sense as the proportion of phenotypic variation that is attributed to genetic factors. Loosely, this means the strength of the genetic effect. Thus a higher heritability will be a larger effect and easier to detect. Heritability is calculated using equations described in [4]. Additionally, we used a constant sample size for all simulations. We selected the sample size of 200 cases (individuals with disease) and 200 controls (individuals without disease) because this is a typical sample that is used in many genetic epidemiology studies.

As discussed by Templeton [21], epistasis, or gene-gene interaction, occurs when the combined effect of two or more genes on a phenotype could not have been predicted from their

independent effects. It is anticipated that epistasis is likely to be a ubiquitous component of the genetic architecture of common human diseases [11]. Current statistical approaches in human genetics focus primarily on detecting the main effects and rarely consider the possibility of interactions [21]. In contrast, we are interested in simulating data using different epistasis models that exhibit minimal independent main effects, but produce an association with disease primarily through interactions. In this study, we use penetrance functions as genetic models. Penetrance functions model the relationship between genetic variations and disease risk. Penetrance is defined as the probability of disease given a particular combination of genotypes.

To evaluate the power of GPNN, SLR, and CART for detecting gene-gene and gene-environment interactions, we simulated case-control data using a variety of epistasis models in which the functional genes are single-nucleotide polymorphisms (SNPs). We selected models that exhibit interaction effects in the absence of any main effects. Interactions without main effects are desirable because they provide a high degree of complexity to challenge the ability of a method to identify gene-gene interactions. If main effects were present, it could be difficult to evaluate whether particular genes were detected due to the main effects or the interactions or both. In addition, it is likely that a method that can detect interacting genes in the absence of main effects will be able to detect main effect genes as well. All models simulated assumed that the contributing factors were single nucleotide polymorphisms with three genotype levels. These models could also assume categorical environmental risk factors with three levels. Thus, all simulations done here can be extended to include categorical environmental factors.

To generate a variety of epistasis models for this study, we selected three criteria for variation. First, we selected epistasis models with a varying number of interacting genes: either two or three. Previous studies had only investigated the power of GPNN using two-gene models [19]. We speculate that common diseases will be comprised of complex interactions among many genes. The number of interacting genes simulated here may still be too few to be biologically relevant. However, few, if any complex gene-gene interaction models are known at this time. Next, we selected two different allele frequencies. An allele frequency of 0.2/0.8 was selected so that we could evaluate the ability of GPNN in situations where there is a relatively rare allele. In addition, the frequency of 0.4/0.6 was selected to allow for the situation where both alleles are relatively common. Finally, we selected a range of heritability values including 3%, 2%, 1.5%, 1%, and 0.5%. These heritability values fall into the realm of very small genetic effects. In comparison, the heritability of many common diseases is much higher. For example, Alzheimer's disease is estimated to have heritability exceeding 60% [1] while breast, colorectal, and prostate cancers are 27%, 35%, and 42% respectively [6]. We chose to simulate data using epistasis models with such small heritability values to test the lower limits of GPNN. Based on previous studies, GPNN has over 80% power when the heritability is between 2%–5% [19]. For this particular study, we wanted to explore even smaller genetic effects to identify the point at which GPNN loses power.

We generated models using software described by Moore et al. [13]. We selected models from all possible combinations of number of interacting genes, allele frequency, and heritability, resulting in 20 total models. The penetrance tables for combinations of two SNPs are shown in Tables 1–10. The penetrance tables for the three SNP models are available from the authors by request. All 20 models were selected because they exhibit interaction effects in the absence of any main effects when genotypes are generated using the Hardy-Weinberg equation. Although the biological plausibility of these models is unknown, they represent the worst-case scenario for a disease-detection method because they have minimal main effects. If a method works well with minimal main effects, presumably the method will continue to work well in the presence of main effects.

Table 1 is an example of a penetrance function for a two-gene epistasis model with no main effects. Each gene is a single SNP with two alleles and three genotypes. In this example, the alleles each have a biological population frequency of $p = 0.2$ $q = 0.8$ with genotype frequencies of p^2 for AA and BB , $2pq$ for Aa and Bb , and q^2 for aa and bb , consistent with Hardy-Weinberg equilibrium. Thus, assuming the frequency of the AA genotype is 0.16, the frequency of Aa is 0.32, and the frequency of aa is 0.64, then the marginal penetrance of BB (i.e. the effect of just the BB genotype on disease risk) can be calculated as $(0.04 * 0.0998) + (0.32 * 0.0984) + (0.64 * 0.0022) = 0.03$. This means that the probability of disease given the BB genotype is 0.03, regardless of the genotype at the other genetic variation. Similarly, the marginal penetrance of Bb can be calculated as $(0.04 * 0.0933) + (0.32 * 0.0996) + (0.64 * 0.0002) = 0.03$. Note that for this model, all of the marginal penetrance values (i.e. the probability of disease given a single genotype, independent of the others) are equal, which indicates the absence of main effects (i.e. the genetic variations do not independently affect disease risk). This is true despite the table penetrance values not being equal. Here, risk of disease is greatly increased by inheriting one of the following high-risk genotype combinations: $AABB$, $AABb$, $AaBB$, $AaBb$, and slightly increased by inheriting genotype combination $aaBb$.

Each data set consisted of 400 cases and 400 controls. We simulated 100 data sets of each model consisting of the functional SNPs and either seven or eight non-functional SNPs for a total of ten SNPs. This resulted in 2000 total datasets. We used a dummy variable encoding for the genotypes where $n-1$ dummy variables are used for n levels (or genotypes) [16]. Based on the dummy coding, these data would have 20 input variables.

2.5 Data Analysis

Next, we used GPNN, SLR, and CART to analyze 100 data sets for each of the epistasis models. The GP parameter settings for GPNN included 10 demes, population size of 200 per deme, 50 generations, reproduction rate of 0.10, crossover rate of 0.90, mutation rate of 0.0, and migration every 25 generations. GPNN is not required to use all the variables as inputs. Here, GPNN performed random variable selection in the initial population of solutions. Through evolution, GPNN selects those variables that are most relevant. We calculated a cross-validation consistency for each SNP in each data set. This measure is defined as the number of times each SNP is in the GPNN model across the ten cross validation intervals. Thus, one would expect a strong signal to be consistent across all ten or most of the data splits, where a false positive signal may be present in only one or a few of the cross validation intervals. We estimated the power of GPNN as the number of times the combination of the correct functional SNPs had a cross-validation consistency that was higher than all other combinations of SNPs in the dataset, divided by the total number of datasets for each epistasis model. Either one or both of the dummy variables could be selected to consider a gene present in the model.

SLR is based on a statistical algorithm that determines the importance of variables and either includes them or excludes them from the model. The importance is determined by the statistical significance of the variable based on a chi-squared test [7]. Here, we used a p-Value of 0.20 to enter the model, and a p-Value of 0.10 to remain in the model. This type of model building procedure can also be referred to as hierarchical model building because to consider interactions among the variables, each variable must remain in the model due to its statistical significance on its own. Thus, using this approach, one can only detect interactions in the presence of main effects of each of the interacting variables. We performed this SLR procedure on each data set. We estimated power of SLR as the number of times the interaction term for the correct functional SNPs was statistically significant in the final SLR model.

CART is a tree-based algorithm that determines the importance of each independent variable and selects the most important variables to split (partition) the data. The importance is determined by the ability of the variable to classify the data. Each split is followed by an

evaluation of all remaining variables to select the subsequent splits. This decision tree process continues until no information is gained from further splits. Once the tree is built, the testing set is used to prune the tree. Using this approach, one can only detect interactions in the presence of main effects of each of the interacting variables. This is due to the fact that the first split is made based on the strongest effect. Thus, if there are no significant effects, no split will be made and the algorithm will terminate. We performed this CART procedure on each data set. We estimated power of CART as the number of times the correct genes were present in the final model selected after tree building and pruning.

3 Results

The results of this study are shown in Tables 11 and 12, and Figures 3 and 4. Here, we list the 20 epistasis models sorted by number of genes, allele frequency, and heritability along the vertical axis. Table 11 and Figure 3 report the power results of the three methodologies. Here, power refers to the method correctly identifying the functional genes. SLR has no power to detect the functional genes in any of the models studied. These results led to some skepticism that logistic regression (LR) may not be able to model the interactions that we had simulated. To be certain that LR was able to model these non-linear interactions, we performed a forward selection LR analysis using explicitly the two or three functional SNPs and their corresponding interaction term (Table 11 column 3 - eLR). We estimated the power of eLR based on the number of data sets where the interaction term was statistically significant. In this study, eLR had between 5–100% and 0–25% power for the two and three gene models respectively. Thus, LR was theoretically able to model these interactions. CART had minimal power to detect the functional variables in the models studied. This is likely due to the fact that the purely epistatic models may have slight main effects in certain splits of the data. Since CART uses a training set and a testing set, it could pick up a significant effect. Once it selects one of the functional genes, it is highly likely to detect the second gene. GPNN, on the other hand, has higher power than SLR and CART for most of the epistasis models. The power of GPNN is higher for the models with two functional genes, and similarly for the models with higher heritability values.

Table 12 and Figure 4 show the false positive results using each method. Here false positives include all models identified as the best model that included genes other than the functional genes. This is different from false positives in random data generated under the null hypothesis. The goal here is to see how often the methods detect genes other than the functional genes. SLR was not evaluated for false positive results because using the stepwise procedure none of the variables are included in any models because there were no significant main effects. Thus, both the power and false positive rate is zero. For the two-gene models, GPNN almost always had fewer false positive genes identified than CART. For the three-gene models, GPNN had higher numbers of false positive genes identified in eight out of ten models. This is due to the fact that CART was not able to make the first split in many datasets since there were no main effects. In contrast to the power results, the false positive results for the three-gene models are quite high.

4 Discussion

Identifying disease susceptibility genes associated with common complex, multifactorial diseases is a major challenge for genetic epidemiology. One of the dominating factors in this challenge is the difficulty in detecting gene-gene and gene-environment interactions with currently available statistical approaches. To deal with this issue, new statistical approaches have been developed such as the GPNN. GPNN has been shown to have higher power than a back propagation NN using simulated data generated under five two-gene epistasis models [19]. The goal of the current study was to compare the power of GPNN to SLR and CART for detecting gene-gene and gene-environment interactions using data simulated from a variety of

epistasis models. Computationally, GPNN is more burdensome than SLR and CART. However, in human genetics the goal is to identify disease susceptibility genes. If one method is more powerful, even if it is more computationally expensive, it may be time well spent. Based on the results shown in Table 11, SLR had no power to detect a statistically significant interaction term, and CART had minimal power. In comparison, GPNN had high power for most of the models examined. Unfortunately, GPNN also had a high false positive rate in the three locus models. This is likely due to the fact that for GPNN, false positives are the reciprocal of power. If the correct genes were not identified, then some set of incorrect genes must have been selected. To decrease the false positive rate, future studies exploring larger population sizes and longer generation time are warranted. The low power and high false positive rate is evidence that GPNN did not converge for the three gene models.

While these results demonstrate the lower limits of GPNN's power to detect gene-gene interactions, there are still many more questions to be addressed. First, it will be important to extend the simulation studies to include more interacting genes, larger sample sizes and a larger range of higher heritability values. In addition, a larger set of epistasis models including those with a small degree of main effect would provide further evidence of the power of GPNN. Finally, it would be interesting to use a different model validation procedure, such as the three-way data split [20], instead of ten-fold cross validation.

The results of this study show that GPNN has higher power than SLR and CART to detect gene-gene interactions in models with very small heritability values. Since most common diseases have overall heritability estimates greater than 20%, and GPNN was shown to have 100% power for heritability of 5% due to the genes examined [19], GPNN should have high power for detecting interactions in most common diseases. The GPNN approach can be applied to studies exploring complex nonlinear interactions associated with binary endpoints in any number of disciplines. It is certainly applicable to any situation where LR or CART had previously been applied. Thus, the GPNN software could be useful for many quantitative studies. And in particular, GPNN is likely to be a powerful pattern recognition approach for the detection of gene-gene and gene-environment interactions in future studies of common human disease.

Acknowledgments

This work was supported by National Institutes of Health grants HL65234, HL65962, GM31304, AG19085, AG20135, AI59694, HD047447, and LM007450.

References

1. Ashford JW, Mortimer JA. Non-familial Alzheimer's disease is mainly due to genetic factors. *J Alzheimers Dis* 2002;4:169–177. [PubMed: 12226536]
2. Bellman, R. *Adaptive Control Processes*. Princeton: Princeton University Press; 1961.
3. Concato J, Feinstein AR, Holford TR. The risk of determining risk with multivariable models. *Ann. Int. Med* 1996;118:201–210. [PubMed: 8417638]
4. Culverhouse R, Suarez BK, Lin J, Reich T. A Perspective on Epistasis: Limits of Models Displaying No Main Effect. *Am J Hum Genet* 2002;70:461–471. [PubMed: 11791213]
5. Duda, RO.; Hart, PE.; Stork, DG. *Pattern Classification*. 2nd Edition. Wiley Interscience; 2000.
6. Hemminki K, Mutanen P. Genetic epidemiology of multistage carcinogenesis. *Mutat. Res* 2001;473:11–21. [PubMed: 11166023]
7. Hosmer, DW.; Lemeshow, S. *Applied Logistic Regression*. New York: John Wiley & Sons Inc.; 2000.
8. Kardia SLR. Context-dependent genetic effects in hypertension. *Curr. Hypertens. Reports* 2000;2:32–38.

9. Koza JR, Rice JP. Genetic generation of both the weights and architecture for a neural network. IEEE Press 1991;Vol II:397–404.
10. Mitchell, M. An Introduction to Genetic Algorithms. Cambridge: MIT Press; 1996.
11. Moore JH. The ubiquitous nature of epistasis in determining susceptibility to common human diseases. Hum Hered 2003;56:73–82. [PubMed: 14614241]
12. Moore, JH. Cross validation consistency for the assessment of genetic programming results in microarray studies. In: Raidl, G., et al., editors. Lecture Notes in Computer Science. Vol. Vol 2611. Berlin: Springer-Verlag; 2003. p. 99-106.
13. Moore JH, Hahn LW, Ritchie MD, Thornton TA, White B. Routine Discovery of High-Order Epistasis Models for Computational Studies in Human Genetics. Applied Soft Computing 2004;4:79–86. [PubMed: 20948983]
14. Moore JH, Parker JS, Olsen NJ, Aune TS. Symbolic discriminant analysis of microarray data in autoimmune disease. Genet Epidemiol 2002;23:57–69. [PubMed: 12112248]
15. Moore JH, Williams SM. New strategies for identifying gene-gene interactions in hypertension. Ann. Med 2002;34:88–95. [PubMed: 12108579]
16. Ott J. Neural networks and disease association. Am. J. Med. Genet 2001;105:60–61. [PubMed: 11425001]
17. Peduzzi P, Concato J, Kemper E, Holford TR, Feinstein AR. A simulation study of the number of events per variable in logistic regression analysis. J. Clin. Epidemiol 1996;49:1373–1379. [PubMed: 8970487]
18. Ritchie MD, Hahn LW, Roodi N, Bailey LR, Dupont WD, Parl FF, Moore JH. Multifactor dimensionality reduction reveals high-order interactions among estrogen metabolism genes in sporadic breast cancer. Am. J. Hum. Genet 2001;69:138–147. [PubMed: 11404819]
19. Ritchie MD, White BC, Parker JS, Hahn LW, Moore JH. Optimization of neural network architecture using genetic programming improves detection of gene-gene interactions in studies of human diseases. BMC Bioinformatics 2003;4:28. [PubMed: 12846935]
20. Rowland, JJ. Generalisation and model selection in supervised learning with evolutionary computation. In: Raidl, G., et al., editors. Lecture Notes in Computer Science. Vol. Vol 2611. Berlin: Springer-Verlag; 2003. p. 119-130.
21. Templeton, AR. Epistasis and complex traits. In: Wolf, J.; Brodie, B., III; Wade, M., editors. Epistasis and Evolutionary Process. Oxford: Oxford University Press; 2000.
22. Venables, WN.; Ripley, BD. Modern Applied Statistics with S. 4th Edition. Springer Publishing; 2002.

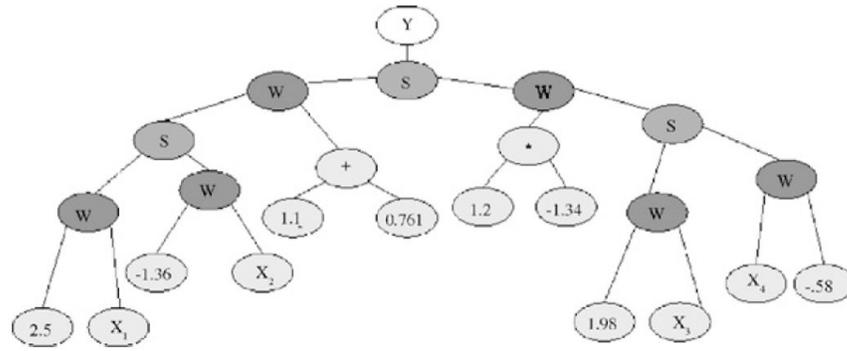


Figure 1. An example of a NN evolved by GPNN. The Y is the output node, S indicates the activation function, W indicates a weight, and X_1 - X_4 are the NN inputs.

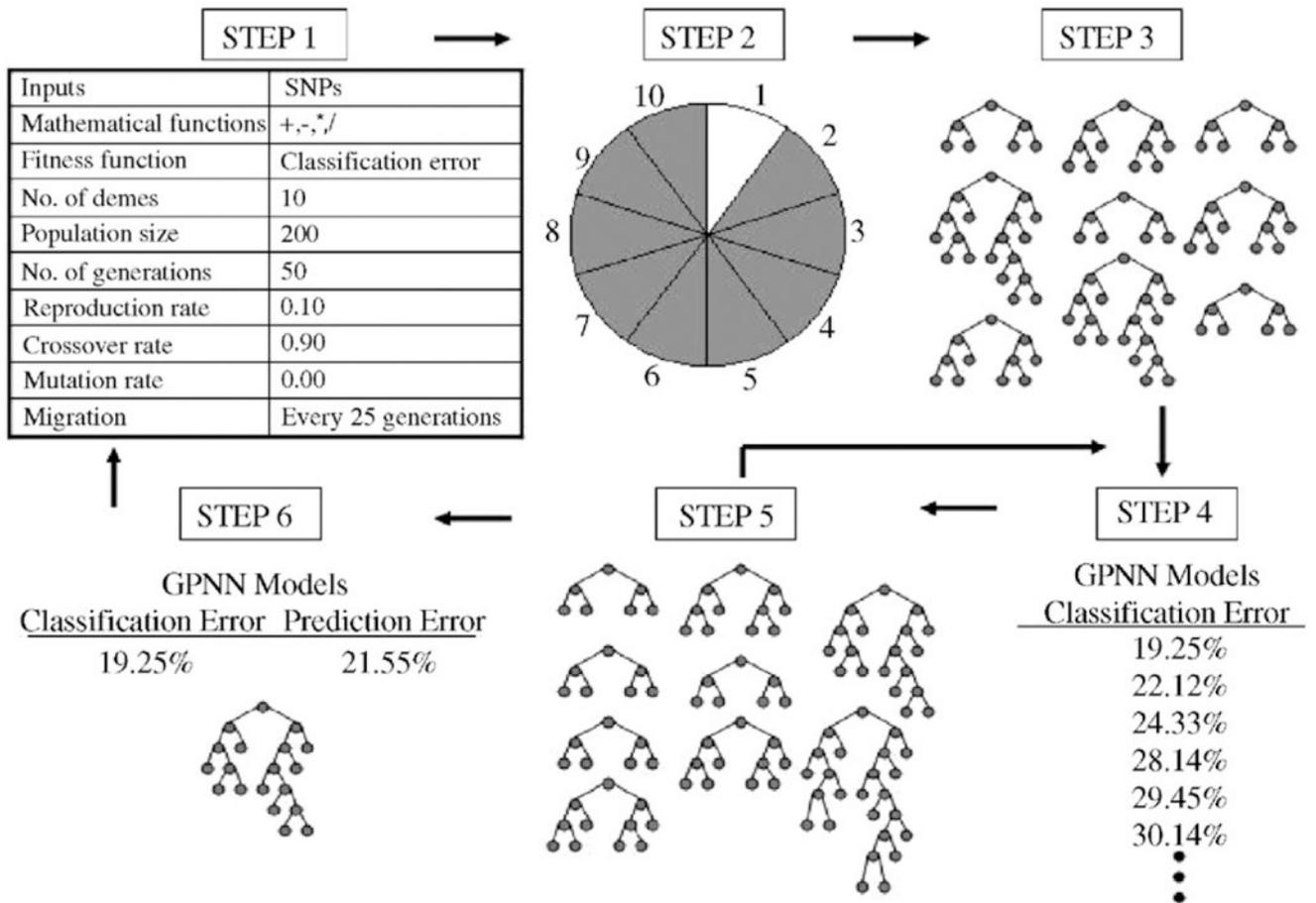


Figure 2.
The Steps of the GPNN algorithm

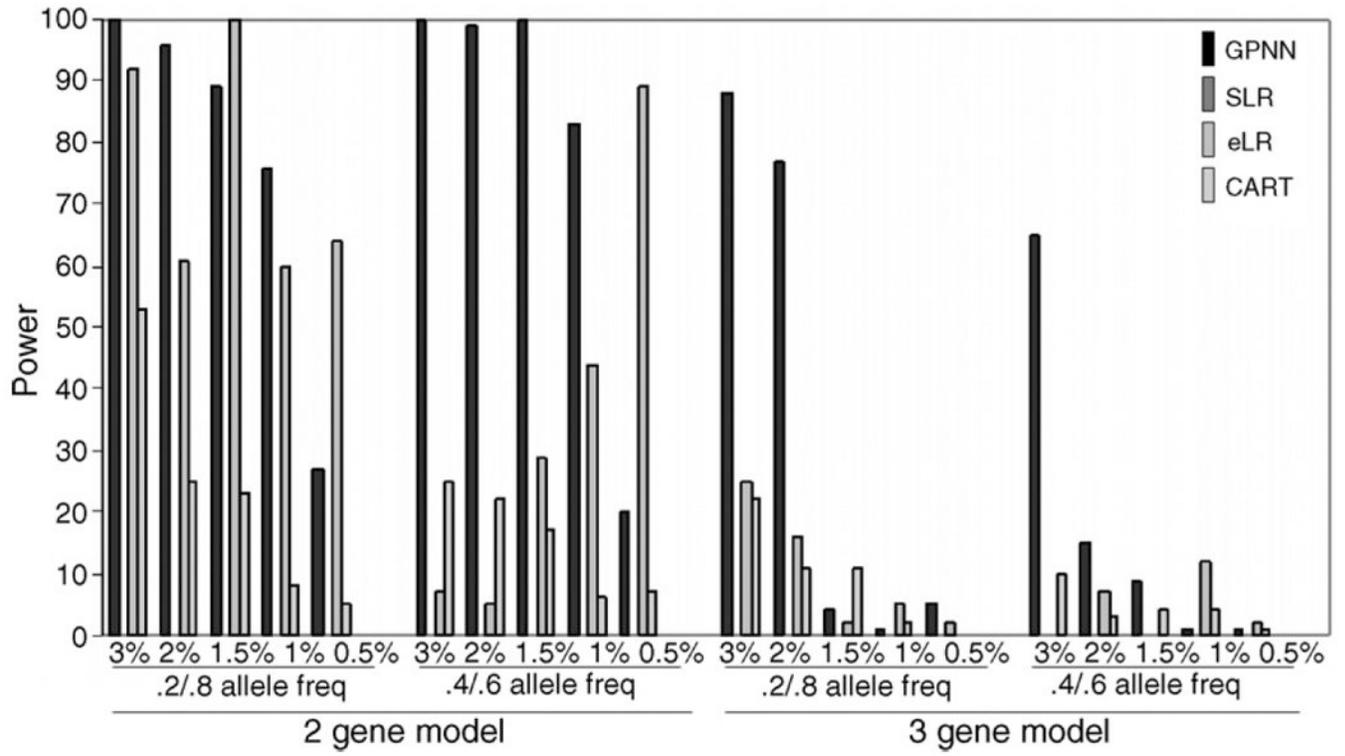


Figure 3.
Power Results of GPNN, SLR, eLR, and CART

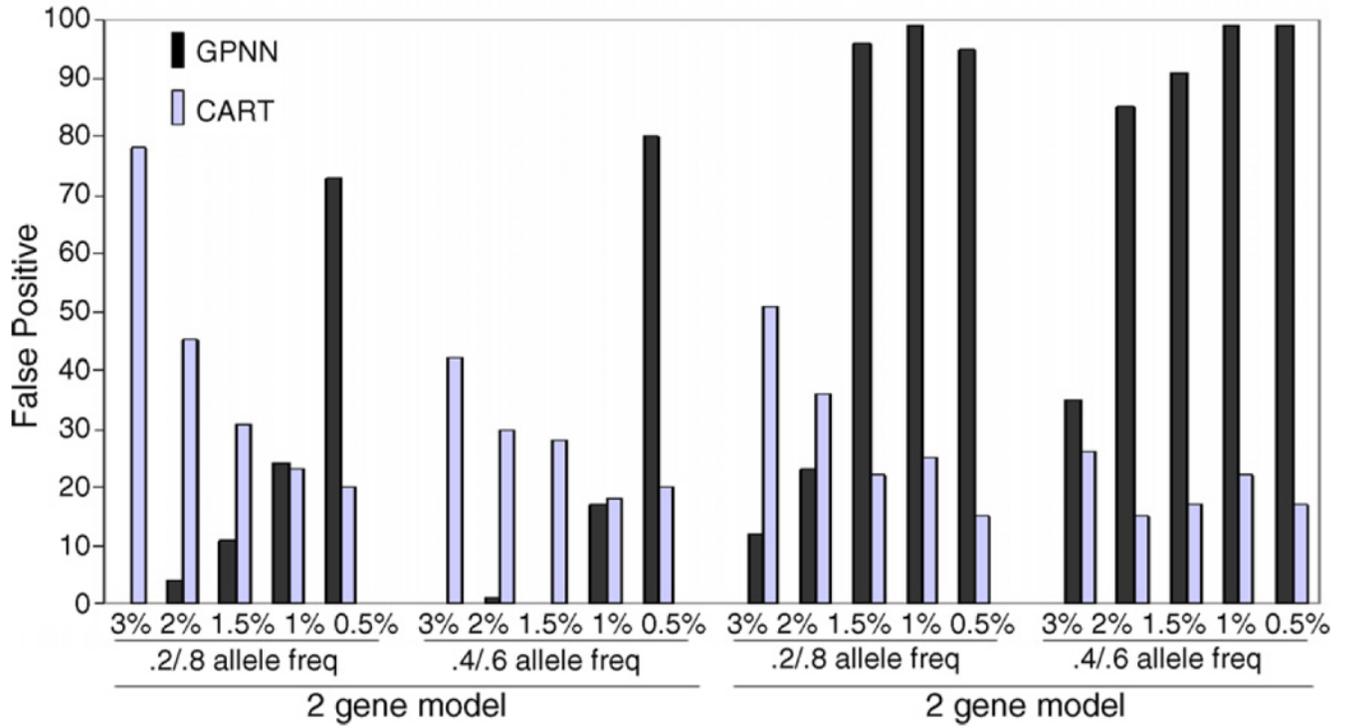


Figure 4.
False Positive Results of GPNN and CART

Table 1Model 1 – Two SNPs, allele frequency 0.2/0.8, $h^2 = 0.030$

	AA	Aa	aa
<i>BB</i>	0.0998	0.0984	0.0022
<i>Bb</i>	0.0933	0.0996	0.0002
<i>bb</i>	0.0028	0.0000	0.0574

Table 2Model 2 - Two SNPs, allele frequency 0.2/0.8, $h^2 = 0.020$

	AA	Aa	aa
<i>BB</i>	0.0786	0.0003	0.0967
<i>Bb</i>	0.0010	0.0013	0.1001
<i>bb</i>	0.0948	0.0998	0.0428

Table 3Model 3 - Two SNPs, allele frequency 0.2/0.8, $h^2 = 0.015$

	AA	Aa	aa
<i>BB</i>	0.0276	0.0942	0.0287
<i>Bb</i>	0.0941	0.0996	0.0226
<i>bb</i>	0.0277	0.0198	0.0657

Table 4Model 4 - Two SNPs, allele frequency 0.2/0.8, $h^2 = 0.010$

	AA	Aa	aa
<i>BB</i>	0.0884	0.0894	0.0307
<i>Bb</i>	0.0710	0.0036	0.0737
<i>bb</i>	0.0368	0.0711	0.0404

Table 5Model 5 - Two SNPs, allele frequency 0.2/0.8, $h^2 = 0.005$

	AA	Aa	aa
<i>BB</i>	0.0539	0.0732	0.0416
<i>Bb</i>	0.007	0.0207	0.0685
<i>bb</i>	0.0732	0.066	0.044

Table 6Model 6 - Two SNPs, allele frequency 0.4/0.6, $h^2 = 0.030$

	AA	Aa	aa
<i>BB</i>	0.0848	0.0754	0.0053
<i>Bb</i>	0.0705	0.0135	0.0967
<i>bb</i>	0.0118	0.0937	0.0131

Table 7Model 7 - Two SNPs, allele frequency 0.4/0.6, $h^2 = 0.020$

	AA	Aa	aa
<i>BB</i>	0.0093	0.0281	0.0902
<i>Bb</i>	0.0491	0.0763	0.0063
<i>bb</i>	0.0625	0.0161	0.0824

Table 8Model 8 - Two SNPs, allele frequency 0.4/0.6, $h^2 = 0.015$

	AA	Aa	aa
<i>BB</i>	0.0381	0.0151	0.073
<i>Bb</i>	0.0485	0.0618	0.0067
<i>bb</i>	0.0288	0.0209	0.0693

Table 9Model 9 - Two SNPs, allele frequency 0.4/0.6, $h^2 = 0.010$

	AA	Aa	aa
<i>BB</i>	0.0465	0.0368	0.0706
<i>Bb</i>	0.0666	0.0691	0.02
<i>bb</i>	0.0314	0.0329	0.0818

Table 10Model 10 - Two SNPs, allele frequency 0.4/0.6, $h^2 = 0.005$

	AA	Aa	aa
<i>BB</i>	0.0161	0.0514	0.0573
<i>Bb</i>	0.0287	0.0442	0.0614
<i>bb</i>	0.0867	0.0511	0.0253

Table 11

Power comparison of GPNN, Stepwise Logistic Regression (SLR), Explicit Logistic Regression (eLR), and Classification and Regression Trees (CART)

# Genes	Model		Power (%)			
	Allele frequency	h^2	GPNN	SLR	eLR	CART
2	0.2/0.8	0.030	100	0	92	53
2	0.2/0.8	0.020	96	0	61	25
2	0.2/0.8	0.015	89	0	100	23
2	0.2/0.8	0.010	76	0	60	8
2	0.2/0.8	0.005	27	0	64	5
2	0.4/0.6	0.030	100	0	7	25
2	0.4/0.6	0.020	99	0	5	22
2	0.4/0.6	0.015	100	0	29	17
2	0.4/0.6	0.010	83	0	44	6
2	0.4/0.6	0.005	20	0	89	7
3	0.2/0.8	0.030	88	0	25	22
3	0.2/0.8	0.020	77	0	16	11
3	0.2/0.8	0.015	4	0	2	11
3	0.2/0.8	0.010	1	0	5	2
3	0.2/0.8	0.005	5	0	2	0
3	0.4/0.6	0.030	65	0	0	10
3	0.4/0.6	0.020	15	0	7	3
3	0.4/0.6	0.015	9	0	0	4
3	0.4/0.6	0.010	1	0	12	4
3	0.4/0.6	0.005	1	0	2	1

Table 12

False Positive Results of GPNN and CART

# Genes	Model		False Positives (%)	
	Allele frequency	h^2	GPNN	CART
2	0.2/0.8	0.030	0	78
2	0.2/0.8	0.020	4	45
2	0.2/0.8	0.015	11	31
2	0.2/0.8	0.010	24	23
2	0.2/0.8	0.005	73	20
2	0.4/0.6	0.030	0	42
2	0.4/0.6	0.020	1	30
2	0.4/0.6	0.015	0	28
2	0.4/0.6	0.010	17	18
2	0.4/0.6	0.005	80	20
3	0.2/0.8	0.030	12	51
3	0.2/0.8	0.020	23	36
3	0.2/0.8	0.015	96	22
3	0.2/0.8	0.010	99	25
3	0.2/0.8	0.005	95	15
3	0.4/0.6	0.030	35	26
3	0.4/0.6	0.020	85	15
3	0.4/0.6	0.015	91	17
3	0.4/0.6	0.010	99	22
3	0.4/0.6	0.005	99	17