



1 **Results:** To address this issue, and to improve interpretability of Random Forest predictions,  
2 we compared different methods for feature importance estimation in real and simulated  
3 datasets with non-additive interactions. As a result, we detected a discrepancy between the  
4 metrics evaluations for the real-world datasets and further established that the permutation  
5 feature importance metric provides more precise feature importance rank estimation for the  
6 simulated datasets with non-additive interactions.

7 **Conclusions:** By analyzing both real and simulated data, we established that the permutation  
8 feature importance metric provides more precise feature importance rank estimation in the  
9 presence of non-additive interactions.

## 10 **Keywords**

11 Machine learning, feature importances, random forest, epistasis, simulation, Alzheimer's  
12 disease, glaucoma.

## 13 **Background**

14 Machine learning has become a common analytical approach for modeling the relationship  
15 between measures of biological systems and clinical outcomes. Models generated from machine  
16 learning can be used for prediction, in which case the biological basis for the pattern being  
17 modeled may not be of interest. However, biological interpretation of machine learning models  
18 can be extremely important if the goal is to generate biological hypotheses that need to be  
19 validated clinically or experimentally. As such, methods for model interpretation have become  
20 an important component of machine learning research. A group of methods provide a graphical  
21 description of the model's global behavior i.e. partial dependency plots [1] and decision tree  
22 surrogate models. Several methods such as individual conditional expectation (ICE) plots [2],

1 local interpretable model-agnostic explanations or LIME [3] and shapley additive explanations  
2 or SHAP and [4] focus on explaining individual model predictions. Another class of methods  
3 focus on assigning weights to individual variables or features based on how much information  
4 they provide to the predictions being made by the machine learning model. This latter approach  
5 generates ‘feature importance scores’ that can be used to create a list of features ranked  
6 according to their importance. This allows the modeler to focus on the most important features  
7 for biological interpretation. Interpretability (along with performance) is the key quality of the  
8 machine learning model, specifically when it is applied to biomedical research goals such as  
9 biomarkers discovery and patient diagnostics.

10 It has been widely discussed that for the complex biomedical phenotypes epistatic interactions  
11 between genes could be present more frequently than we commonly think [5 – 7]. Indeed, gene-  
12 gene interactions have been detected in multiple genome-wide association studies of various  
13 disease phenotypes, including Alzheimer’s disease [8], cataracts [9], diabetes [10, 11],  
14 cardiovascular diseases [12, 13], neurological diseases [14, 15], and various cancer types [16,  
15 17]. Epistasis has been defined in several different ways [18, 19]. Here, we employ epistasis term  
16 as interactions between two or more gene loci such that the phenotype cannot be accurately  
17 predicted by simply adding the effects of individual gene loci. This is a statistical definition of  
18 epistasis as it measures the deviation from additivity using models that summarize genotypic and  
19 phenotypic variability of human population data. In contrast to that, biological epistasis is  
20 identified at the cellular level in an individual as a result of a physical interactions among  
21 molecules within the biological network. The relationship between two epistatic concepts is a  
22 complicated matter in such that statistical epistasis doesn’t literally translate into the biological  
23 [18].

1 The intrinsic complexity of the epistatic interactions creates several analytical and practical  
2 challenges for its detection via traditional statistical methods due to its inability to detect non-  
3 additive effects in the large volume of data from GWAS studies. Machine learning methods have  
4 more flexibility in its power to detect underlying complexity of genetic architecture and,  
5 therefore, has been widely used in epistasis discovery [20]. A large body of research has been  
6 accumulated on epistatic interaction detection with neural network, support vector machines,  
7 multifactor dimensionality reduction and random forest (RF) models [21, 22]. Specifically, RF  
8 algorithm is known for its ability to take into account non-additive effects through its  
9 hierarchical tree-based structure [20, 23]. A number of studies have been conducted on the  
10 integration of RF method into the epistatic interaction discovery. Among them RF with sliding  
11 window sequential forward feature selection method [24], a hybrid of RF and mutual information  
12 network methodology that reports reduction in RF bias towards the marginal effects [25],  
13 synthetic feature RF complemented with statistical epistasis network approach that incorporates  
14 genetic interactions within and among the biochemical pathways to explain its association with  
15 the disease outcome [26], mixed RF approach that accounts for both population structure and  
16 epistatic interactions [27], relative recurrency variable importance metric (r2VIM) [28] for RF  
17 that generates variable set with main and interaction effect, and permuted RF method that  
18 identifies interacting pair of SNPs by calculating the effect of disrupted interaction on the RF  
19 prediction error rate [29].

20 In addition to notorious performance on datasets with epistasis, benchmark studies have  
21 recognized the RF classification algorithm to be among the best classifiers for the majority of the  
22 real-world datasets [30, 31]. RF is conveniently interpretable with built-in feature importance  
23 scores implemented in a majority of popular programming languages and analytics platforms

1 including Python and R. These are calculated with entropy or Gini importance criterion. Despite  
2 the utility of RF feature important scores, some studies have reported a bias introduced by the RF  
3 feature importance scores when working with categorical, grouped, and varying types of features  
4 [32]. An alternative way to estimate feature importances with the RF classifier is by calculating  
5 permutation feature importance (PFI) scores. This metric employs an exhaustive permutation  
6 concept where features are permuted one at a time and the importance scores expressed via the  
7 difference in the ML algorithm's performance score. While computationally intensive for large  
8 feature sets, this approach is advantageous due to its applicability to any machine learning  
9 method of any complexity and, therefore, is independent of the characteristics of any given  
10 algorithm.

11 In this study, we aim to improve the interpretability of RF predictions for genetic data in the  
12 presence of non-additive interactions by comparing two feature importance metrics: RF's built-in  
13 feature importance coefficients (BIC) and PFI coefficients. We use two real-world datasets with  
14 previously described non-additive interactions to compare model interpretation using the two  
15 different feature importance score metrics. We also compare the metrics using simulated data  
16 with varying levels of interaction, imbalance in case/control ratio, and sample size where the  
17 ground truth is known.

18

## 19 **Results**

### 20 **Evaluation of feature importance metrics performance with HIBACHI simulated datasets.**

21 An RF classification algorithm with Gini impurity criterion was fitted on all simulated  
22 replicates with epistatic interactions. For majority of them, predictive balanced accuracy with 10-  
23 fold cross-validation was estimated as 1.0 or very close to 1.0 (Fig. S1). PFI and BIC metrics

1 were estimated for the fitted RF models and further compared to the real feature importances  
2 retrieved with the HIBACHI sensitivity analysis. The percentage of successful rank identification  
3 per 100 replicates of each experiment was reported in Table 1. For all combinations of factors  
4 that were considered in HIBACHI simulations, PFI metric consistently outperformed BIC metric  
5 in the ability to determine feature importance order. The most accurate PFI evaluation was  
6 produced for the datasets in the category with sample size 1000, two-way IG and 50% of cases,  
7 where the most important feature (F1) was identified precisely in 90% of replicates, F2 – 89%,  
8 and F3 – 77%. For the sample size 10000 category, the most accurate estimate by PFI was done  
9 for the three-way IG, with 25% of cases where F1 was identified correctly for 88% of replicates,  
10 F2 – for 79%, F3 -for 88%, F4 – for 83%. Overall, feature ranks estimated by PFI metric were at  
11 least twice more accurate than BIC estimates for nearly all the corresponding experiments.

12 Interestingly, for the majority of replicates of the experiments with the two-way IG, the PFI  
13 metric was able to identify the top three features, while for the three-way IG this came up to the  
14 top four features. We investigated further into the population of the feature importance effect  
15 sizes and discovered that experiments with the two-way interactions observe the same range of  
16 effect sizes for the top two important features, while the effect sizes of the bottom two features  
17 are 2-3 orders of magnitude smaller than the top ones or zero (Fig. 3A, B). To support this  
18 observation Table 2 report the percentage of replicates with zero effect size at each feature  
19 importance rank: for the majority of the experiments with the two-way IG, two out of five  
20 features didn't have an effect on the phenotype. A more thorough exploration of the fitness  
21 function landscape confirmed that in the majority of cases only one pair or trio of features had  
22 enough time to evolve strong epistatic interactions due to the limited simulation time. Hence, we  
23 observe that the interacting pair of features share the major informative contribution towards the

1 phenotype, while non-interacting features contribution effect is orders of magnitude smaller and  
2 insufficient to be detected accurately by both importance metrics (Fig.4A, B).

### 3 **Evaluation of the feature importances in real-world datasets with epistatic interactions.**

4 We used HIBACHI simulation framework to obtain datasets with strong epistatic interaction  
5 between the features - genetic variants and the phenotype. We were able to demonstrate that  
6 regardless of the dataset's parameters, the best method to determine features' informative  
7 contribution to the phenotype is PFI. To validate our findings for real-world data we analyzed  
8 two datasets with complex disease phenotype which previously have been identified to have  
9 epistatic interactions among SNPs. The subset of seven SNPs was selected from the Alzheimer's  
10 disease data and the subset of six SNPs was selected from the Glaucoma dataset as described  
11 above (see Methods 2.4). A ViSEN analysis and plot for the Alzheimer's dataset (Fig. 5A)  
12 revealed one SNP with large independent effect affiliated with ApoE gene (rs429358, MI 0.08) –  
13 a known risk factor for Alzheimer disease, and three strong two-way gene-gene interactions  
14 (rs4955208 - rs7782571, IG 0.05; rs12785149 - rs12209418, IG 0.05; rs2414325 - rs1931073, IG  
15 0.05). ViSEN analysis of the Glaucoma dataset revealed two strong gene-gene interactions  
16 (rs7738052 - rs1489169, IG 0.015, and rs10915315 - or rs1266924, IG 0.015). We built a  
17 predictive model for each dataset with the RF classifier: the model was tuned with grid search  
18 and had classification balanced accuracy 62.6% for the Alzheimer dataset, and 60 % for the  
19 Glaucoma dataset. We verified the significance of the cross-validated balanced accuracy scores  
20 of the optimized classifiers by doing a permutation test [33] and confirmed that both models  
21 have reliable classification performance (Figure 6). We further calculated PFI and BIC estimates  
22 for the optimized RF models. For Alzheimer dataset (Fig. 7A) the most important feature,  
23 according to both PFI and BIC metrics, is the SNP with the largest independent effect –

1 rs429358. This SNP has been discovered by ViSEN analysis and expected to have strong effect  
2 on the phenotype. However, the consecutive feature importance order diverged between two  
3 metrics. According to PFI rank, SNPs located at second and third position belong to the same  
4 interaction and valued with the highest IG out of all calculated pairwise interactions for this  
5 dataset (Fig 5A). At the same time, SNPs located at the second and third position by BIC ranking  
6 didn't create a strong pairwise interactions with each other. We observed a similar discrepancy  
7 between the metric rank evaluations for the Glaucoma datasets (Fig 7B): while SNP rs2157719  
8 has the largest main effect among all SNPs and was assigned to be the top feature by both  
9 metrics, the following order of SNPs was outlined conversely by different metrics. SNP  
10 rs10915315 has the largest number of detected pairwise interaction among all SNPs and is  
11 ranked second by PFI, while only fourth by BIC; SNP rs1489169 has three pairwise interactions  
12 detected and is ranked third by PFI and second by BIC. Interestingly, both metrics identified the  
13 less important and non-important SNPs in the same order, suggesting that the issue with the  
14 discrepancy between the metrics evaluations should be address for the top features only.

## 15 **Discussion**

16 Multiple research studies suggested that epistatic interactions are widespread by nature and  
17 that many genes work in an interactive manner [5]. Epistatic interactions are expected to be  
18 found in a variety of pathophysiological processes with one of them most likely to be  
19 Alzheimer's disease [8]. The most powerful predictor of Alzheimer's disease at this time is  
20 ApoE E4 gene variation: one or two copies of ApoE is associated with an increased risk of  
21 disease onset [34]. However, some carriers of ApoE E4 variation haven't developed an  
22 Alzheimer's disease so it is very likely that other genetic factors are involved in disease's  
23 pathophysiology. ViSEN entropy-based analysis revealed several strong pairwise genetic

1 interactions, along with the known largest independent signal from the ApoE variant (rs429358)  
2 (Fig. 5A). Furthermore, ViSEN method allocated epistatic interactions within the Glaucoma  
3 disease dataset: several strong pairwise interactions in addition to the independent main effect  
4 contribution from the SNP affiliated with retinal ganglion cells pathology (rs2157719) have been  
5 confirmed (Fig. 5B). Evaluation of an informative contribution of a single genetic factor  
6 involved in the two- or three-way interactions towards the phenotypic outcome is a challenging  
7 analytical task and it requires a non-linear solution which RF classifier is notorious for.  
8 Therefore, we aimed to identify whether RF classifier is capable to detect genetic signatures that  
9 were previously confirmed by ViSEN analysis and whether different RF's feature importance  
10 metrics will be able to agree on the rank.

11 Two feature importance metrics were considered, PFI and BIC, and each was compared after  
12 RF analysis of data derived from genome-wide association studies of Glaucoma and  
13 Alzheimer's. The resulting feature ranking confirms the lack of consensus between the studied  
14 metrics (Fig 6). Indeed, while BIC and PFI identified the SNPs with the largest independent  
15 main effects (rs429358 and rs2157719) as a top feature in both real-world datasets, genes  
16 involved in the epistatic interactions have been assigned different importance ranks by different  
17 metrics. More specifically, in the Alzheimer data, the second and the third most important  
18 feature predicted by PFI belong to the same interacting pair (rs4955208 and rs7782571), while  
19 the same rank positions predicted by BIC were occupied by SNPs that do not belong to the same  
20 genetic interaction (Fig 5A, Fig 7A). A similar discrepancy has been observed in the Glaucoma  
21 dataset: PFI and BIC indicated different interacting pairs of SNPs as the second and third most  
22 important feature (Fig. 7A), however SNP predicted to be the second most important by PFI is a

1 part of four additional SNP-SNP interactions, while SNP predicted second by BIC is a part of  
2 only three interacting SNP pairs (Fig. 5B).

3 Such uncertainty has been associated with RF predictions in the past and we attempted to  
4 reveal the true interpretation with the computational experiments driven by HIBACHI  
5 simulations. The HIBACHI framework has the ability to consider any desirable biological  
6 concept in the form of mathematical expressions that define the genotype-phenotype relationship  
7 and evolve models that can be used to simulate data consistent with that relationship. We set up a  
8 simulation goal to maximize two- or three-way interactions among features and compared RF's  
9 feature importance metrics with the sensitivity analysis results of the simulated data that  
10 provided us with the ground truth information about the feature ranks (Figure 2). In all  
11 HIBACHI experimental setups, which included such factors as the proportion of cases and  
12 controls, sample size and interaction complexity, PFI metrics produced the most precise feature  
13 ranking (Table 1, Fig. 3). Although BIC metric misplaced feature ranks for the majority of the  
14 replicates, it correctly identified features that belong to the interactive pair or trio by putting them  
15 as a top two or three features correspondingly. Therefore, we can suggest that BIC metric can  
16 still be used with some caution, however, when there is a need for an absolute precision, PFI  
17 estimation method should be used.

18 PFI has some limitations we didn't cover in our experiments – it is particularly sensitive to  
19 highly correlated features. In case such features are present in the dataset, PFI requires a special  
20 preprocessing directed onto the removal of such features. An example of this could be a  
21 hierarchical clustering based on the Spearman's rank correlation coefficients with following  
22 cluster-based filtering. In future studies, a more advanced permutation feature importance  
23 techniques may be considered which can include pairwise permutation importance, joint

1 importance by maximal subtrees, and joint variable importance as well as corrected Gini  
2 importance score [35]. In addition, a more complex epistatic schemes needs to be explored (e.g.  
3 different combinations of marginal and interaction effects as in [36], multiple strong interactions  
4 and/or combination of two- and three-way interactions). Within our simulation framework we  
5 mostly observed one or two strong interactions and this might be an overly simplistic given the  
6 complexity of common diseases such as Alzheimer’s and glaucoma.

## 7 **Conclusion**

8 In this study, we performed a comparative analysis of feature importance metrics with the aim  
9 to improve Random Forest’s interpretability in datasets with complex interactions. By analyzing  
10 both real and simulated data, we established that the permutation feature importance metric  
11 provides more precise feature importance rank estimation in the presence of non-additive  
12 interactions.

## 13 **Methods**

### 14 **Random forest and its properties**

15 RF is a popular ML algorithm because it often demonstrates good performance while  
16 remaining relatively easy to optimize and interpret. RF algorithms belong to a Bagging  
17 (Bootstrap Aggregation) type of ensemble ML methods where a group of weak learners in a  
18 form of decision trees (DT) classify the outcome using majority vote. Decision trees are sensitive  
19 to the data they are trained and often suffer from high variance problem especially when the  
20 depth of the tree is high. To address that RF algorithm trains each tree on a subset of samples  
21 drawn from the complete dataset with the bootstrap procedure. An additional source of  
22 randomness within the RF algorithm is introduced during the construction of a tree when  
23 selecting a split node from the random sample of features (in place of the greedy search through

1 all feature set like in DT). These two sources of randomness aim to decorrelate weak learners  
2 and correspondingly decrease the variance of an estimator by combining diverse trees prediction  
3 via majority vote.

4 During the construction of a DT, the decrease in the error function can be calculated for a  
5 feature at each split node. In a classification task, this is often done via estimating the Gini score  
6 or the entropy score. The function decreases can be averaged across all trees and returned as  
7 feature importances score (the greater the decrease the higher feature importance). Feature  
8 importance scores are often conveniently implemented as a RF function which makes this  
9 method more interpretable.

#### 10 **Permutation feature importances**

11 Permutation feature importance metrics were first introduced by Breiman in his Random  
12 Forest manuscript [37] and further extended by Altmann [38] to correct for the bias of the RF's  
13 Gini importance and entropy criterion for feature selection. We utilize a custom implementation  
14 of PFI which could be applied to any machine learning classification and regression algorithms  
15 (Figure 1). Here, PFI metric is calculated with following steps: 1) the dataset is shuffled and split  
16 into the training and testing datasets 2) the model is fitted on the training dataset and the  
17 balanced accuracy is estimated on the testing dataset, 3) feature 1 out of N is permuted for the  
18 testing dataset 4) the balanced accuracy is estimated on the permuted testing dataset 5) the  
19 relative decrease of the permuted and non-permuted balanced accuracies is calculated and stored  
20 as relative decrease in accuracy, 4) step 2 and 3 are repeated for the remaining N-1 features, 5)  
21 steps 1 through 4 are replicated for N-1 times with the new seeds for shuffle and split procedure,  
22 6) mean of relative decrease in accuracy per feature is calculated across the splits and is used as  
23 features' PFI value. Retrieved PFI values were normalized to sum to 1.

## 1 **Evaluation of the feature importance metrics**

### 2 **Data simulation using HIBACHI**

3 We used Heuristic Identification of Biological Architectures for simulating Complex  
4 Hierarchical Interactions (HIBACHI) software to simulate genetic datasets with non-additive  
5 epistatic interactions of different complexity. The HIBACHI method employs biological and  
6 mathematical frameworks to connect genotype and phenotype [39, 40] At the bottom of the  
7 biological framework is the concept of information transfer from the DNA sequence to a clinical  
8 phenotype through complex interactions at multiple levels: gene expression, pathway, and cell.  
9 Within this framework, a population of samples is expressed as a collection of genes (genotype)  
10 each of which has three variants 0, 1 or 2. The mathematical framework's goal is to specify a  
11 relationship between genotype and phenotype in terms of logical, arithmetical and other  
12 functions. HIBACHI merges the frameworks by evolving a mathematical expression tree which  
13 when applied to a genotype, generates a binary clinical phenotype.

14 HIBACHI employs Genetic Programming (GP) as an optimization engine. One of the key  
15 characteristics of GP is a fitness function which is represented through the mathematical  
16 expression that satisfy a specific objective of interest defined by user. This objective could be a  
17 performance of a machine learning pipeline, complexity of genetic interactions, odds ratio of  
18 genetic effect sizes, length of expression tree, etc. or a combination of those in the form of the  
19 multi-objective fitness function. Additionally, user to allowed to specify the length of the  
20 optimization (by specifying the number of generations), the population size (number of samples)  
21 and genotype size (number of genes) of the dataset.

22 At the beginning of the GP optimization process, a population of individuals with random  
23 mathematical expression trees is initialized and is further subjected to mutational and

1 recombinational processes. This process serves as a source of variation for the expression trees  
 2 and have a pre-defined rate. After that, a user-defined fitness function is calculated and the best-  
 3 fitted individuals are selected for the next round. At the last optimization round, an overall best-  
 4 fitted individual is used as an output dataset.

5 For the aims established within this study, we wanted to generate datasets with non-additive  
 6 epistatic interactions that would involve two and three genetic variants and, therefore, we have  
 7 defined a fitness function that maximizes two-way and three-way information gain (IG) term.  
 8 Entropy-based IG approach to detect epistatic interactions has been introduced by Moore et al.  
 9 [41] and extended by Fan et al. [42]. Two-way IG defined as:

$$10 \quad IG(X; Y; Z) = I(X, Y; Z) - I(X; Z) - I(Y; Z) ,$$

11 Where  $I$  is mutual information that describes the dependency between variables X, Y and Z.  
 12 It measures the reduction of uncertainty of one variable (Z) given the knowledge of others (X  
 13 and Y) It is expressed with entropy terms:

$$14 \quad I(X, Y; Z) = H(Z) - H(Z|X, Y); \quad I(X; Z) = H(Z) - H(Z|X); \quad I(Y; Z) = H(Z) - H(Z|Y),$$

15 Where entropy  $H$  defined with the probability mass function:

$$H(Z|X) = H(Z, X) - H(X); \quad H(Z) = - \sum_{z \in Z} p(z) \log p(z);$$

$$H(Z, X) = - \sum_{z \in Z} \sum_{x \in X} p(z, x) \log p(z, x)$$

16 Definitions of three-way IG can be found in Hu et al., 2013 [43]. We used the following version  
 17 of this term:

$$IG(X; Y; Z; N) = I(X, Y, Z; N) - IG(X; Y; N) - IG(X; Z; N) - IG(Y; Z; N) - I(X; N) - I(Y; N) \\ - I(Z; N)$$

1        Additionally, a second fitness objective was set up as maximization of expression tree length,  
2 to encourage multiple combinations of genetic interactions. In addition to the variability in the  
3 interaction complexity, the following factors have been considered in the HIBACHI  
4 experimental schemes: percent of cases (25% and 50%) to address an imbalanced dataset  
5 structure, and sample size (1000 and 10000). Each experimental setup was reproduced 100 times  
6 using random seed generator and the whole population of replicates was considered in the  
7 consecutive analysis. All simulated data used here is available upon request.

### 8    **Sensitivity analysis**

9        We implemented a HIBACHI-based sensitivity analysis to determine the true feature  
10 importances ranks and the effect sizes. For this the permutation-based framework was  
11 implemented with the following steps (Figure 2): 1) split HIBACHI-generated dataset into the  
12 outcome vector and the feature set 2) permute feature vector 1 out of N and re-evaluate the  
13 outcome by applying the HIBACHI-generated expression tree 3) calculate the dissimilarity of the  
14 outcome as of mismatch between the perturbed and unperturbed feature set 4) repeat the estimate  
15 2) for the remaining feature and normalize the counts by the total sum; 5) replicate steps 2-4 100  
16 times and calculate the average of the replicates per feature as a final true feature importance  
17 score. Sensitivity scores were further normalized to sum to 1.

### 18    **Real world data analysis.**

19        To examine the convergence of the RF's feature importance metrics we used two real-world  
20 datasets with evidence for non-additive interactions. The first includes preselected SNPs from a  
21 genome-wide association study of Alzheimer's Disease while the second includes preselected  
22 SNPs from a genome-wide association study of Primary Open Angle Glaucoma (POAG). The  
23 Alzheimer's dataset came from the Alzheimer's Disease Neuroimaging Initiative during which

1 functional MRI was taken every six to twelve months for patients with three health conditions  
2 (neuro-typical (identified here as controls), and mild cognitive impairment and Alzheimer's  
3 disease (identified here as cases)). A computational evolution system [44] identified a model of  
4 seven SNPs with evidence of non-additive interactions and a classification accuracy of 0.738.  
5 Among these SNPs are the SNP with the large main effect that is located in the APOE gene  
6 (rs429358) – a known risk factor for Alzheimer disease, four SNPs located within genes with  
7 known functionality/disease state (rs1931073 – an intergenic region near the PPAP2B gene that  
8 is participating in cell-cell interactions, rs7782571 – near the ISPD gene that is associated with  
9 the Walker-Walburg syndrome, rs4955208 – in the OSBPL10 genes which are expressed into  
10 intracellular lipid receptor, rs12209418 – in the PKIB gene that codes a protein kinase inhibitor)  
11 and two remaining SNPs (rs2414325, rs12785149) are located within genes with unknown  
12 functionality.

13 The glaucoma dataset came from the Glaucoma Gene Environment Initiative study and  
14 contained with POAG individuals identified as cases and healthy individuals as controls. This  
15 dataset has been previously analyzed by Moore et al. [45] with the EMERGENT algorithm that  
16 resulted in the identification of a model of six SNP's with evidence of non-additive interactions  
17 and a classification accuracy of 0.615. Two of these SNPs (rs2157719, and rs1266924) are  
18 located within the genes that were previously associated with glaucoma disease, two SNPs  
19 (rs10915315, and rs1489169) are located within the genes that are associated with glaucoma-  
20 non-related diseases and relevant pathology, and two more SNPs (rs936498, and rs7738052) are  
21 located within the genes that were not previously associated with any disease, but have a known  
22 functionality that is relevant to visual cortex and retina development.

1 We used the visualization of the statistical interaction network (ViSEN) method [46] to  
2 analyze and visualize SNP main effects, and two-way and three-way gene-gene interactions  
3 among SNPs for real-world datasets. The ViSEN method calculates the mutual information (MI)  
4 between individual SNP (genotype) and the phenotype, the pairwise interaction between every  
5 pair of SNPs and the phenotype and the three-way interaction between every combination of  
6 three SNPs and the phenotype via the IG term. A positive IG indicates synergistic (i.e. non-  
7 additive) effects of SNPs on the phenotype. The IG metrics used in the ViSEN method were  
8 designed to detect pure epistatic interactions and excluded all lower-order effects (by subtracting  
9 all main effects and pairwise synergies in cases of the three-way term).

#### 10 **List of abbreviations**

11 RF- Random Forest, PFI – permutation feature importance, HIBACHI - Heuristic  
12 Identification of Biological Architectures for simulating Complex Hierarchical Interactions, BIC  
13 – build-in coefficients

#### 14 **Declarations**

#### 15 **Ethics approval and consent to participate**

16 Human data used in preparation of this article were obtained from the Alzheimer’s Disease  
17 Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)) and The Primary Open-Angle  
18 Glaucoma Genes and Environment (GLAUGEN) Study at dbGaP database.

#### 19 **Consent for publication**

20 Not applicable

#### 21 **Availability of data materials**

1 The datasets used and/or analyzed during the current study are available from the corresponding  
2 author on request.

### 3 **Competing interests**

4 The authors declare that they have no competing interests

### 5 **Funding**

6 This work was supported by National Institutes of Health (USA) grants LM010098 and  
7 AI116794

### 8 **Authors' contributions**

9 AO conceived the study, performed data simulation and data analysis. JHM conceived and  
10 supervised the study. All authors were involved in writing the manuscript. All authors read and  
11 approved the manuscript.

### 12 **Acknowledgments**

13 Not applicable

### 14 **References**

15 1. Hastie T, Tibshirani R and Friedman J, *The Elements of Statistical Learning*, Second Edition,  
16 Section 10.13.2, Springer, 2009.

17 2. Goldstein A, et al. "Peeking inside the black box: Visualizing statistical learning with plots of  
18 individual conditional expectation." *Journal of Computational and Graphical Statistics* 24.1  
19 2015; 44-65

20 3. Ribeiro, M. T., Singh, S., and Guestrin, C. "Why should I trust you?": Explaining the predictions  
21 of any classifier. In *International Conference on Knowledge Discovery and Data Mining (KDD)*,  
22 2016.

- 1 4. Lundberg SM, Lee SI. A unified approach to interpreting model predictions. Proc. Adv. Neural  
2 Inf. Process. Syst. 2017; 4768-4777.
- 3 5. Moore, JH The ubiquitous nature of epistasis in determining susceptibility to common human  
4 diseases. Human Heredity. 2003; 56, 73–82
- 5 6. Carlborg, O., & Haley, C. S. Epistasis: Too often neglected in complex trait studies? Nature  
6 Reviews Genetics 5, 618–625 (2004) doi:10.1038/nrg1407
- 7 7. Mackay, T.F., Moore, J.H. Why epistasis is important for tackling complex human disease  
8 genetics. Genome Med. 2014; 6, 42
- 9 8. Hohman TJ, Bush WS, Jiang L, et al. Discovery of gene-gene interactions across multiple  
10 independent data sets of late onset Alzheimer disease from the Alzheimer Disease Genetics  
11 Consortium. Neurobiol Aging. 2016;38:141–150. doi:10.1016/j.neurobiolaging.2015.10.031
- 12 9. Pendergrass SA, Verma SS, Holzinger ER, Moore CB, Wallace J, Dudek SM, Huggins W,  
13 Kitchner T, Waudby C, Berg R, McCarty CA, Ritchie MD. Next-generation analysis of  
14 cataracts: determining knowledge driven gene-gene interactions using Biofilter, and gene-  
15 environment interactions using the PhenX Toolkit. Pac Symp Biocomput. 2013;147-58.  
16 Corrected and republished in: Pac Symp Biocomput. 2015;:495-505. PMID: 23424120; PMCID:  
17 PMC3615413.
- 18 10. Bell JT, Timpson NJ, Rayner NW, Zeggini E, Frayling TM, et al. Genome-wide association  
19 scan allowing for epistasis in type 2 diabetes. Ann Hum Genet. 2011;75:10–19.
- 20 11. Manduchi E, Chesi A, Hall MA, Grant SFA, Moore JH (2018) Leveraging putative enhancer-  
21 promoter interactions to investigate two-way epistasis in Type 2 Diabetes GWAS. Pac Symp  
22 Biocomput 2018:548–558

1 12. Lippert C, Listgarten J, Davidson RI, et al. An exhaustive epistatic SNP association analysis on  
2 expanded Wellcome Trust data [published correction appears in *Sci Rep.* 2013 Feb 18;3:1321.  
3 Poong H. *Sci Rep.* 2013;3:1099. doi:10.1038/srep01099

4 13. Meng Y, Groth S, Quinn JR, Bisognano J, Wu TT. 2017 An exploration of gene-gene  
5 interactions and their effects on hypertension. *Int. J. Genomics.* 2017, 7208318.

6 14. Sha Q, Zhang Z, Schymick JC, Traynor BJ, Zhang S. Genome-Wide Association Reveals Three  
7 SNPs Associated With Sporadic Amyotrophic Lateral Sclerosis Through a Two-Locus  
8 Analysis. *BMC Med Genet.* 2009;10:86. 86.

9 15. Steffens M, Becker T, Sander T, Fimmers R, Herold C, Holler D, A, Leu C, Herms S, Cichon S,  
10 Bohn B, Gerstner T, Griebel M, Nöthen M, M, Wienker T, F, Baur M, P: Feasible and  
11 Successful: Genome-Wide Interaction Analysis Involving All  $1.9 \times 10^{11}$  Pair-Wise Interaction  
12 Tests. *Hum Hered* 2010;69:268-284. doi: 10.1159/000295896

13 16. Chu M, Zhang R, Zhao Y, Wu C, Guo H, Zhou B, Lu J, Shi Y, Dai J, Jin G, Ma H, Dong J, Wei  
14 Y, Wang C, Gong J, Sun C, Zhu M, Qiu Y, Wu T, Hu Z, Lin D, Shen H, Chen F. A genome-  
15 wide gene-gene interaction analysis identifies an epistatic gene pair for lung cancer susceptibility  
16 in Han Chinese. *Carcinogenesis.* 2014 Mar;35(3):572-7. doi: 10.1093/carcin/bgt400. Epub 2013  
17 Dec 9. PMID: 24325914; PMCID: PMC3941747.

18 17. Shen Z, Li Z, Song J, Chen Y Shi. Genome-wide two-locus interaction analysis identifies  
19 multiple epistatic SNP pairs that confer risk of prostate cancer: a cross-population study *Int. J.*  
20 *Cancer*, 140 (9) (2017), pp. 2075-2084, 10.1002/ijc.30622

21 18. Moore JH and Williams SM (2005), Traversing the conceptual divide between biological and  
22 statistical epistasis: systems biology and a more modern synthesis. *Bioessays*, 27: 637-646.  
23 doi:10.1002/bies.20236

- 1 19. Phillips P. Epistasis — the essential role of gene interactions in the structure and evolution of  
2 genetic systems. *Nat Rev Genet* 9, 855–867 (2008). <https://doi.org/10.1038/nrg2452>
- 3 20. McKinney BA, Reif DM, Ritchie MD, Moore JH. Machine learning for detecting gene-gene  
4 interactions: a review. *Appl Bioinformatics*. 2006;5(2):77–88. doi:10.2165/00822942-  
5 200605020-00002
- 6 21. Moore JH, Asselbergs FW, Williams SM. Bioinformatics challenges for genome-wide  
7 association studies. *Bioinformatics*. 2010;26(4):445–455. doi:10.1093/bioinformatics/btp713
- 8 22. Koo CL, Liew MJ, Mohamad MS, and Salleh AH. (2013). A Review for detecting gene-gene  
9 interactions using machine learning methods in genetic epidemiology. *Biomed. Res. Int.*  
10 2013:432375. doi: 10.1155/2013/432375
- 11 23. Goldstein BA, Polley EC, Briggs FBS. Random forests for genetic association studies. *Stat Appl*  
12 *Genet Mol Biol*. 2011;10(1):32.
- 13 24. Jiang R, Tang W, Wu X, Fu W. A random forest approach to the detection of epistatic  
14 interactions in case-control studies. *BMC Bioinformatics*. 2009
- 15 25. Pan QX, T. Hu, J. D. Malley, A. S. Andrew, M. R. Karagas, and J. H. Moore, “Supervising  
16 random forest using attribute interaction networks,” in *Evolutionary Computation, Machine*  
17 *Learning and Data Mining in Bioinformatics*, L. Vanneschi, W. S. Bush, and M. Giacobini, Eds.,  
18 vol. 7833 of *Lecture Notes in Computer Science*. 2013;104–116.
- 19 26. Pan Q, Hu T, Malley JD, Andrew, A.S., Karagas, M.R. and Moore, J.H. A System-Level  
20 Pathway-Phenotype Association Analysis Using Synthetic Feature Random Forest. *Genet.*  
21 *Epidemiol*. 2014; 38: 209-219. doi:10.1002/gepi.21794

- 1 27. Stephan J, Stegle O and Beyer A. A random forest approach to capture genetic effects in the  
2 presence of population structure. *Nat Communications*. **6**, 7432 (2015). [https://doi-](https://doi.org.proxy.library.upenn.edu/10.1038/ncomms8432)  
3 [org.proxy.library.upenn.edu/10.1038/ncomms8432](https://doi.org.proxy.library.upenn.edu/10.1038/ncomms8432)
- 4 28. Holzinger E.M., Szymczk S., Dasgupta A., Malley J., Li Q., Bailey-Wilson J.E. (Jan 4–8, 2015)  
5 Variable Selection Method for the Identification of Epistatic Models. Paper presented at the  
6 Pacific Symposium on Biocomputing (PSB), Maui, HI.
- 7 29. Li J, Malley JD, Andrew AS, Kargas MR, Moore JH. Detecting gene-gene interactions using a  
8 permutation-based random forest method. *BioDataMining*. 2016; **9**, 14  
9 <https://doi.org/10.1186/s13040-016-0093-5>
- 10 30. Fernandez-Delgado, M., Cernadas, E., Barro, S. & Amorim, D. Do we need hundreds of  
11 classifiers to solve real world classification problems? *J. Mach. Learn. Res.* 2014;15, 3133–3181
- 12 31. Olson RS, La Cava W, Orzechowski P. et al. PMLB: a large benchmark suite for machine  
13 learning evaluation and comparison. *BioData Mining*. 2017; **10**, 36.  
14 <https://doi.org/10.1186/s13040-017-0154-4>
- 15 32. Strobl C, Boulesteix AL, Zeileis A, Hothorn T. Bias in random forest variable importance  
16 measures: illustrations, sources and a solution. *BMC Bioinformatics*. 2007;8:25. Published 2007  
17 Jan 25. doi:10.1186/1471-2105-8-25
- 18 33. Ojala and Garriga. Permutation Tests for Studying Classifier Performance. *The Journal of*  
19 *Machine Learning Research* (2010) vol. 1
- 20 34. Corder, E. H., et al. Gene dose of apolipoprotein E type 4 allele and the risk of Alzheimer's  
21 disease in late onset families. *Science* 261, 921–923 (1993)
- 22 35. Nembrini S, König IR, Wright MN, The revival of the Gini importance?, *Bioinformatics*, Volume  
23 34, Issue 21, 01 November 2018, Pages 3711–3718, <https://doi.org/10.1093/bioinformatics/bty373>

- 1 36. Wright MN, Ziegler A, König IR. Do little interactions get lost in dark random forests? BMC  
2 Bioinformatics. 2016;17:145. doi:10.1186/s12859-016-0995-8
- 3 37. Breiman L, Random Forests. Machine Learning, 45, 5-32 (2001)
- 4 38. Altmann, A., Toloşi, L., Sander, O., Lengauer, T.: Permutation importance: a corrected feature  
5 importance measure. Bioinformatics. 2010; 26(10):1340-7
- 6 39. Moore JH, Amos R, Kiralis J, Andrews PC. Heuristic identification of biological architectures  
7 for simulating complex hierarchical genetic interactions. *Genet Epidemiol.* 2015;39(1):25–34.  
8 doi:10.1002/gepi.21865
- 9 40. Moore JH, Shestov M, Schmitt P, Olson RS. A heuristic method for simulating open-data of  
10 arbitrary complexity that can be used to compare and evaluate machine learning methods. Pac  
11 Symp Biocomput. 2018;23:259–267.
- 12 41. Moore JH, Gilbert JC, Tsai CT, Chiang FT, Holden T, Barney N, White BC. 2006 A flexible  
13 computational framework for detecting, characterizing, and interpreting statistical patterns of  
14 epistasis in genetic studies of human disease susceptibility. *J Theor Biol* 241: 252– 261.
- 15 42. Fan R, Zhong M, Wang S et al. Entropy-based information gain approaches to detect and to  
16 characterize gene–gene and gene–environment interactions/correlations of complex diseases.  
17 *Genet Epidemiol.* 2011; 35:706–721
- 18 43. Hu T, Chen Y, Kiralis JW, Collins RL, Wejse C, Sirugo G, Williams SM, Moore JH: An  
19 information-gain approach to detecting three-way epistatic interactions in genetic association  
20 studies. *J Am Med Inform Assoc.* 2013
- 21 44. Moore JH, Douglas P, Saykin, Andrew, Shen Li. 2013. Exploring Interestingness in a  
22 Computational Evolution System for the Genome-Wide Genetic Analysis of Alzheimer's  
23 Disease. Third Indonesian-American Kavli Frontiers of Science Symposium. Bali, Indonesia.

- 1 45. Moore JH, Greene CS, Hill DP. Identification of Novel Genetic Models of Glaucoma Using the  
2 “EMERGENT” Genetic Programming-Based Artificial Intelligence System. In: Riolo R, Worzel  
3 WP, Kotanchek M, editors. *Genet. Program. Theory Pract.* XII [Internet]. Springer International  
4 Publishing; 2015 [cited 2016 Dec 12]. p. 17–35.
- 5 46. Hu T, Chen Y, Kiralis JW et al (2013) ViSEN: methodology and software for visualization of  
6 statistical epistasis networks. *Genet Epidemiol* 37:283–285

7

8

9