

# Improving eQTL Analysis using a Machine Learning Approach for Data Integration: A Logistic Model Tree Solution

**Stefano Beretta**<sup>1;2</sup> stefano.beretta@disco.unimib.it (+39)02 6448 7832

**Mauro Castelli**<sup>3</sup> mcastelli@novaims.unl.pt (+351)213 828 610

**Ivo Gonçalves**<sup>3;4</sup> igoncalves@novaims.unl.pt (+351)213 828 206

**Ivan Kel**<sup>2</sup> ivan.kel@itb.cnr.it (+39)02 2642 2606

**Valentina Giansanti**<sup>2</sup> valentina.giansanti@itb.cnr.it (+39)02 2642 2606

**Ivan Merelli**<sup>2;\*</sup> ivan.merelli@itb.cnr.it (+39)02 2642 2606

Addresses:

1 Dipartimento di Informatica Sistemistica e Comunicazione, Università degli Studi di Milano-Bicocca, Viale Sarca 336, 20126 Milan, Italy.

2 Istituto di Tecnologie Biomediche, Consiglio Nazionale Ricerche, Via F.lli Cervi 93, 20090 Segrate, Italy.

3 NOVA IMS, Universidade Nova de Lisboa, Campus de Campolide, 1070-312, Lisboa, Portugal.

4 CISUC, Department of Informatics Engineering, University of Coimbra, Coimbra, Portugal.

\* Corresponding author

**This is the author accepted manuscript version of the article published by EMERALD as:** Beretta, S., Castelli, M., Gonçalves, I., Kel, I., Giansanti, V., & Merelli, I. (2018). Improving eQTL Analysis Using a Machine Learning Approach for Data Integration: A Logistic Model Tree Solution. *Journal of Computational Biology*, 25(10), 1091-1105. <https://doi.org/10.1089/cmb.2017.0167>



This work is licensed under a [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/).



# Improving eQTL Analysis using a Machine Learning Approach for Data Integration: A Logistic Model Tree Solution

## Authors:

	E-Mail	Telephone
Stefano Beretta <sup>1,2</sup>	stefano.beretta@disco.unimib.it	(+39)02 6448 7832
Mauro Castelli <sup>3</sup>	mcastelli@novaims.unl.pt	(+351)213 828 610
Ivo Gonçalves <sup>3,4</sup>	igoncalves@novaims.unl.pt	(+351)213 828 206
Ivan Kel <sup>2</sup>	ivan.kel@itb.cnr.it	(+39)02 2642 2606
Valentina Giansanti <sup>2</sup>	valentina.giansanti@itb.cnr.it	(+39)02 2642 2606
Ivan Merelli <sup>2,*</sup>	ivan.merelli@itb.cnr.it	(+39)02 2642 2606

## Addresses:

- <sup>1</sup> Dipartimento di Informatica Sistemistica e Comunicazione,  
Università degli Studi di Milano-Bicocca,  
Viale Sarca 336, 20126 Milan, Italy.
- <sup>2</sup> Istituto di Tecnologie Biomediche,  
Consiglio Nazionale Ricerche,  
Via F.lli Cervi 93, 20090 Segrate, Italy.
- <sup>3</sup> NOVA IMS,  
Universidade Nova de Lisboa,  
Campus de Campolide, 1070-312, Lisboa, Portugal.
- <sup>4</sup> CISUC, Department of Informatics Engineering,  
University of Coimbra,  
Coimbra, Portugal.

★ Corresponding author

## Abstract

eQTL analysis is an emerging method for establishing the impact of genetic variations (such as single nucleotide polymorphisms) on the expression levels of genes. Although different methods for evaluating the impact of these variations are proposed in the literature, the results obtained are mostly in disagreement, entailing a considerable number of false positive predictions. For this reason, we propose an approach based on Logistic Model Trees that integrates the predictions of different eQTL mapping tools in order to produce more reliable results. More precisely, we employ a machine learning based method using logistic functions to perform a linear regression able to classify the predictions of three eQTL analysis tools (namely, R/qtl, MatrixEQTL, and mRMR). Given the lack of a reference dataset and that computational predictions are not so easy to test experimentally, the performance of our approach is assessed using data from the DREAM5 challenge. The results show the quality of the aggregated prediction is better than that obtained by each single tool in terms of both precision and recall. We also performed a test on real data, employing genotypes and microRNA expression profiles from *C. elegans*, which proved that we were able to correctly classify all the experimentally validated eQTLs. This good results come both from the integration of the different predictions, and from the ability of this machine learning algorithm to find the best cut-off thresholds for each tool. This combination makes our integration approach suitable for improving eQTL predictions for testing in a laboratory, reducing the number of false positive results.

1  
2  
3  
4  
5  
6  
7  
8       **Keywords:** Machine Learning, Evolutionary Algorithm, Genetic  
9 Programming, eQTL Analysis, Data Integration  
10  
11

## 14   **1 Introduction**

16  
17       Due to the importance of their role in a variety of regulatory processes  
18 and in several diseases, single nucleotide polymorphisms (SNPs) are  
19 widely studied, with particular attention given to their interactions  
20 with genes and pathologies [Merelli (2013)]. For this purpose, one  
21 of the methods most used to link the expression of genes/proteins  
22 to different genotypes is the expression quantitative trait loci (eQTL)  
23 mapping, which studies the impact of SNPs on differential measurable  
24 gene transcript levels.  
25  
26

27       More precisely, eQTL analyses seek to identify genomic variations  
28 whose genotypes affect the expression of specific genes. In the last few  
29 years, much effort has been made to define methods for performing  
30 eQTL analysis, with many different approaches being developed.  
31

32       In particular, most eQTL studies separately test for each SNP-  
33 gene pair. The association between expression and genotype is mainly  
34 tested using the linear regression and ANOVA models, as well as non-  
35 linear techniques including generalized linear and mixed models. For  
36 example, the technique implemented in MatrixEQTL [Shabalin (2012)]  
37 tests the associations between each possible pair of SNP and transcript  
38 using two models: linear and ANOVA. In the former model, the ef-  
39 fect of genotype is additive and its significance is evaluated using a  
40 t-statistic, while in the latter the genotype is modeled as a categorical  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 variable and the significance is evaluated using the ANOVA approach.  
9

10 Bayesian regression can also be used for eQTL analysis [Servin and Stephens (2007)]  
11 as well as models accounting for pedigree [Abecasis (2002)] and la-  
12 tent variables [Leek and Storey (2007)]. Moreover, several methods  
13 have been developed to find groups of SNPs associated with the ex-  
14 pression of a single gene [Hoggart (2008), Kao (1999), Lee (2008),  
15 Zeng (1993)].  
16

17 Another approach, employed by one the most popular tools for  
18 eQTL analysis, R/qtl [Broman (2003)], is to use Hidden Markov Mod-  
19 els (HMM) to deal with missing genotype data. In this way, the  
20 method can deal with the presence of genotyping errors, backcrosses,  
21 intercrosses, and phase-known four-way crosses when performing eQTL  
22 analysis.  
23

24 Recently, a tool called mRMR was proposed in the literature [Huang and Cai (2013)].  
25 This machine learning based method tests all the possible types of de-  
26 pendencies by taking advantage of Mutual Information (MI) to assess  
27 the association between genotypes and gene expressions. More pre-  
28 cisely, the key point is to consider not only the interaction between an  
29 SNP and a gene, but also the redundancy among genes, which is used  
30 to detect indirect associations.  
31

32 Although several computational methods have been developed for  
33 predicting eQTLs [Wright (2012)] and the results obtained in several  
34 studies using state-of-art methods reveal interesting aspects, there are  
35 substantial differences when different outputs are compared, since  
36 all the algorithms provide some false positive results. Some tools  
37 have been developed to compare results obtained with different tech-  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

niques [Kel (2016)]. This approach allows the possibility of performing a genome wide analysis in parallel fashion, but an exhaustive integration of the results is lacking. As observed in some analysis, there are significant differences in the results, also considering the list of top eQTL predictions (see [Huang and Cai (2013)] and the Supplementary Material of [Kel (2016)]). Although such comparisons are inconclusive since for a full evaluation different scenarios should be considered (such as the quantity of data, the number of variables, and the missing information), it is clear the tools are often in disagreement and robust results are still missing. The same issue has been successfully addressed in the prediction of miRNA-target interactions [Beretta (2017)] and also in the reconstruction of cancer networks by combining Bayesian approaches and evolutionary techniques [Beretta (2016)]. This lack of consensus in the eQTL predictions opens up the possibility of exploiting machine learning techniques to reduce the number of false positives. The biggest challenge in performing this task relates to the lack of standard reference datasets [Michaelson (2010)].

In response to this issue, the DREAM (Dialogue on Reverse Engineering Assessments and Methods) community, whose aim is to assess the performance of different techniques to solve problems coming from the biological field, proposed a challenge for eQTL mapping called DREAM5 (<https://www.synapse.org/#!Synapse:syn2787209>). The goal of one of these challenges (DREAM5 SYSGEN A) was to reverse-engineer the interaction networks, starting from synthetic genetic variations and gene expression data. This problem simulates the mapping of both *cis*- and *trans*-eQTL, on *in silico* networks and, to this end,

1  
2  
3  
4  
5  
6  
7  
8 it can be used to evaluate the accuracy of the existing methods and  
9  
10 tools. Finally, after the DREAM5 challenge was concluded, the ref-  
11  
12 erence networks were released, which we used as *ground truth* for the  
13  
14 predictions.

15  
16 Machine learning techniques have been applied on the same dataset  
17  
18 as described in [Ackermann (2012)]. More precisely, in this work *Ran-*  
19  
20 *dom Forests* and *LASSO* methods have been applied to the DREAM5  
21  
22 datasets to directly perform the eQTL mapping.

23  
24 On the other hand, in this work we present an integrative machine  
25  
26 learning approach to combine the results obtained by the tools pre-  
27  
28 viously introduced for eQTL analysis based on a supervised learning  
29  
30 process. The aim is to build a system to combine the predictions of  
31  
32 existing tools with the specific goal of establishing whether an inter-  
33  
34 action between an SNP and a gene is likely to be real or not. To  
35  
36 accomplish this task, the technique described in this study considers  
37  
38 two different phases: initially the problem related with data unbal-  
39  
40 ancing is addressed using an under-sampling technique; in the second  
41  
42 phase, the transformed data are used to build a predictive model (i.e.,  
43  
44 a classifier) by means of a Model Tree [Landwehr (2005)], in which the  
45  
46 linear regression functions take the form of logistic functions. These  
47  
48 two steps of the proposed method, as well as the design choices, are  
49  
50 described in the subsequent sections of this paper. By taking ad-  
51  
52 vantage of the DREAM5 datasets, for which we have the standard  
53  
54 reference (i.e. we know, of all the possible pairs of SNP-gene inter-  
55  
56 actions, those that are true), we trained the proposed classifier by  
57  
58 taking into account the predictions obtained with R/qtl, MatrixE-  
59  
60



1  
2  
3  
4  
5  
6  
7  
8 QTL, and mRMR, namely the tools most commonly used for eQTL  
9 analysis [Wright (2012)].

11 Using the trained model, we also performed a test on real data,  
12 using genotypes and microRNA expression profiles from *C. elegans*,  
13 showing that we were able to correctly classify all the experimentally  
14 validated eQTLs. The advantage of our method is that it is able  
15 not only to take advantage of different results that are combined, but  
16 also to find the best threshold values for each tool in each dataset to  
17 identify the correct SNP-gene pairs.  
18  
19  
20  
21  
22  
23

24 We wish to point out that since the main focus of this work is on  
25 integrating predictions with a machine learning technique based on  
26 Logistic Model Trees, we did not consider other sources of information.  
27 In any case, we highlight the fact it would be possible to integrate  
28 additional sources of data such as chromatin accessibility, and histone  
29 modification and methylation to confirm, for example, trans-eQTL  
30 predictions [Merelli (2015), Duggal (2014)].  
31  
32  
33  
34  
35  
36  
37

38 The code and the datasets used for the experimental evaluation  
39 of our method available at <https://github.com/beretta/eQTL-LMT>.  
40 Additionally, a script to fully replicate the experimental analysis on  
41 the available DREAM5 datasets is present.  
42  
43  
44

45 The paper is structured as follows: Section 2 describes the two  
46 steps of the proposed machine learning approach based on Logistic  
47 Model Trees, while Section 3 presents the experimental analysis and  
48 discusses the results. Section 4 summarizes the main contributions of  
49 this work and suggests future research directions.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 2 Method

This section outlines the two main components of the proposed system. Section 2.1 introduces the unbalanced data problem and describes the sampling technique considered in this study. Section 2.2 describes the logistic model tree algorithm used to analyze the resulting balanced datasets.

### 2.1 Data Sampling

One of the core issues to be addressed before building a classifier relates to the distribution of the data among the classes. A dataset is unbalanced when at least one class is represented by only a small number of training examples (called the minority class) while the other class(es) make up the majority. In this work, we consider a binary classification problem, that is, a problem having two classes where one contains most of the examples (majority class) while the other one only has a few examples (minority class). In the eQTL mapping problem, the minority class is composed of the true SNP-gene interactions, while the majority class consists of all the other possible combinations of SNPs and genes. In this scenario, classifiers usually have good accuracy on the majority class, but very low accuracy with respect to the minority class due to the influence the larger majority class has on traditional training criteria [Ganganwar (2012)]. In fact, most classification algorithms seek to minimize the error rate, corresponding to the percentage of the incorrect predictions of class labels, but ignore the difference between types of misclassification errors, implicitly assuming that all misclassification errors come at equal

1  
2  
3  
4  
5  
6  
7  
8 cost.

9  
10 In several real-world applications this assumption may be not true  
11 and classifiers built on these unbalanced data may thus produce unsat-  
12 isfactory results. When an application is characterized by unbalanced  
13 data, two main options can be considered: (1) ignoring the problem;  
14 and (2) balancing the dataset. The first option typically results in a  
15 final classifier that is biased, with all the instances classified in the  
16 class corresponding to the majority class. The motivation is related  
17 to the fact the model looks at the data and cleverly decides the best  
18 thing to do is to always predict “Majority class” to achieve a high  
19 level of accuracy. The second option distinguishes two different ap-  
20 proaches [Ganganwar (2012)]: data-level techniques and algorithmic-  
21 level techniques. At the data level, several solutions proposing many  
22 different forms of resampling have been developed, while at the algo-  
23 rithmic level the presented solutions are related to the particular ma-  
24 chine learning technique under investigation [Nguyen (2009)]. At the  
25 data level, two main approaches may be identified: (1) oversampling  
26 techniques; and (2) undersampling techniques. The former family  
27 contains methods designed to increase the number of instances in the  
28 minority class by replicating them in order to include a higher repre-  
29 sentation of the minority class in the sample. This process can be per-  
30 formed by randomly selecting the instances of the minority class that  
31 must be replicated, or by using more complex criteria [Guo (2008)].  
32 While oversampling methods lead to no information loss, they increase  
33 the likelihood of overfitting since they replicate the minority class  
34 events [Japkowicz and Stephen (2002)]. In addition, the criterion used  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 to select the observations to be replicated has an important impact  
9 on the final model's performance [Japkowicz and Stephen (2002)]. On  
10 the other hand, undersampling methods aim to balance class distri-  
11 bution by selecting majority class examples until the majority and  
12 minority class instances are balanced out. Also in this case it is pos-  
13 sible to define a random selection strategy, or more complex crite-  
14 ria [Guo (2008)]. The biggest disadvantage of undersampling relates  
15 to the fact the sample chosen might be an inaccurate representation  
16 of the class distribution. On the other hand, undersampling does not  
17 require synthetic data to be generated and improves run time by re-  
18 ducing the number of training data samples.

19  
20  
21  
22  
23  
24  
25  
26  
27  
28 In this study, considering the large amount of biological data avail-  
29 able nowadays, we focus on a specific data level method, which tries  
30 to balance the data by undersampling the majority class. The method  
31 selects a number of training instances from the majority class equal  
32 to the number of training instances in the minority class.

33  
34  
35  
36  
37  
38 To select the undersampling method, the “unbalance” R package  
39 was used and the following undersampling techniques were evaluated:

- 40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
- Random. This method aims to balance class distribution through the random elimination of instances belonging to the majority class. The data considered in this work are composed of two classes, in which the majority class covers 99.8% of the available data. Hence, random sampling of the observations of the majority class will probably result in a poor selection of the instances. For this reason, we did not take this undersampling strategy into account.

- 1  
2  
3  
4  
5  
6  
7  
8 • Neighborhood Cleaning Rule (NCL) [Laurikkala (2001)]. NCL  
9 works by first removing negatives examples which are misclassi-  
10 fied by their 3-nearest neighbors. Second, the neighbors of each  
11 positive examples are found and those belonging to the majority  
12 class are removed. While this approach is effective, it can create  
13 a computational bottleneck for very large datasets, with a large  
14 majority class. Hence, it is an unsuitable option given the nature  
15 of the data considered in this paper.  
16  
17  
18  
19  
20  
21  
22
- 23 • Tomek Links [Tomek (1976)]. Tomek links consist of points that  
24 are each other's closest neighbors, but do not share the same class  
25 label. More formally, let us assume a dataset  $\{E_1, \dots, E_n\} \in$   
26  $[\mathbb{R}]^k$ , in which each  $E_i$  has exactly one of the two labels “+” or  
27 “-”. A pair  $(E_i, E_j)$  is called a Tomek link if  $E_i$  and  $E_j$  have  
28 different labels, and there is not an  $E_k$  such that  $d(E_i, E_k) <$   
29  $d(E_i, E_j)$  or  $d(E_j, E_k) < d(E_i, E_j)$ , where  $d(x, y)$  is the distance  
30 between  $x$  and  $y$ . Tomek links can be used as an undersampling  
31 method, where only instances belonging to the majority class are  
32 eliminated.  
33  
34  
35  
36  
37  
38  
39  
40  
41
- 42 • Condensed Nearest Neighbor Rule (CNN) [Hart (1968)]. CNN  
43 is used to find a consistent subset of instances, not necessarily  
44 the smallest one. A subset  $S' \subset S$  is consistent with  $S$  if using a  
45 1-Nearest Neighbor (1-NN),  $S'$  correctly classifies the instances  
46 in  $S$ . The rationale is to eliminate the instances from the major-  
47 ity class that are distant from the decision border, because they  
48 can be considered less relevant for the learning task.  
49  
50  
51  
52  
53  
54
- 55 • One-Side Selection (OSS) [Kubat and Matwin (1997)]. This method  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 is an undersampling method resulting from the application of  
9 Tomek links [Tomek (1976)] followed by the CNN rule. Basi-  
10 cally, it combines the advantages of the two previously presented  
11 methods.  
12  
13  
14

15  
16 Considering the advantages and limitations of the proposed tech-  
17 niques, and taking account of the nature of the available data, in our  
18 work we decided to employ the One-Side Selection (OSS) method. As  
19 explained in [Kubat and Matwin (1997)], this undersampling method  
20 aims to create a dataset composed only of “safe” instances. In order  
21 to do this, OSS removes (from the majority class) instances that are  
22 noisy, redundant, or near the decision border.  
23  
24  
25  
26  
27  
28

29 To detect less reliable instances of this class, one can classify the  
30 instances in four different categories (see Figure 1): (i) instances that  
31 suffer from class-label noise (the points in the bottom left corner of  
32 Figure 1); (ii) borderline examples that are close to the boundary  
33 between minority and majority regions (these points are unreliable  
34 because even a small amount of noise can push the instance into the  
35 wrong region); (iii) redundant instances that can be safely removed  
36 (points in the upper right corner of Figure 1); and (iv) safe examples  
37 that are worth maintaining for future classification tasks.  
38  
39  
40  
41  
42  
43  
44  
45

46 Considering these four cases, an intelligent learner may try to re-  
47 move the borderline instances as well as those suffering from class-  
48 noise. Finally, the learner will remove the redundant instances. In  
49 particular, Tomek links allow the removal of the noisy and borderline  
50 examples, while the CNN technique removes examples from the ma-  
51 jority class that are far away from the decision border. Regarding the  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 second component of the OSS technique, CNN works by selecting the  
9  
10 set  $S$  of reference instances (data points that are needed for an ac-  
11 curate classification) from the dataset obtained after the Tomek link  
12 phase, such that a 1-NN (1-Nearest Neighbors) with  $S$  can classify the  
13 examples almost as accurately as the same 1-NN does with the whole  
14 dataset.  
15  
16  
17  
18

19 Here is the pseudocode of the OOS algorithm:

- 21 1. let  $S$  be the original training set;
- 22  
23 2. initially,  $C$  contains all the minority class instances of  $S$  and one  
24 randomly selected instance of the majority class;
- 25  
26 3. classify  $S$  with the 1-NN rule using the instances in  $C$  and com-  
27 pare the assigned labels with the original ones;
- 28  
29 4. move all the misclassified instances into  $C$ ;
- 30  
31 5. remove from set  $C$  all the majority class instances participating  
32 in Tomek links (to remove borderline and noisy instances); and  
33  
34 6. return the newly obtained set.  
35  
36  
37  
38  
39

40 All related procedural details are outlined in [Kubat and Matwin (1997)].  
41  
42  
43

## 44 2.2 Logistic Model Trees

45  
46  
47 Once the dataset is created by exploiting the above mentioned tech-  
48 nique, a supervised classifier must be trained on it in order to ob-  
49 tain the final results. Among the existing methods in the literature  
50 to solve supervised learning problems [Hastie (2009), Witten (2016),  
51 Caruana and Niculescu-Mizil (2006)], linear models and tree induc-  
52 tion methods [Loh (2011)] have gained popularity in the data mining  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 community, both for the prediction of nominal classes and numeric  
9 values. As reported in [Landwehr (2005)], the former approaches fit a  
10 simple (linear) model to the data, where the process of model fitting  
11 is quite stable, resulting in low variance, but potentially high bias.  
12 The latter ones exhibit low bias but often high variance: they search  
13 a less restricted space of models, allowing them to capture nonlinear  
14 patterns in the data, but making them less stable and prone to over-  
15 fitting. It is therefore not surprising that neither of these two methods  
16 is generally superior [Bishop (2006)].

17  
18 In recent years, the possibility to combine these two schemes into  
19 model trees, i.e. trees that contain linear regression functions on  
20 the leaves, has been investigated [Sumner (2005), Hosmer Jr (2013),  
21 Harrell Jr (2015)]. More precisely, as reported in [Landwehr (2005)],  
22 the main idea behind model trees is to combine the advantages of tree  
23 induction methods and linear models. Hence, model trees rely on sim-  
24 ple regression models if only little and/or noisy data are available and  
25 add a more complex tree structure if there is enough data to warrant  
26 such a structure. Using model trees with linear regression functions  
27 might not be the best choice when a classification task must be ad-  
28 dressed because this approach produces several trees (one per class)  
29 and thus makes the final model hard to interpret.

30  
31 A better solution to tackle classification tasks is to use a combi-  
32 nation of a tree structure and logistic regression models resulting in a  
33 single tree. In detail, logistic regression is a regression model originally  
34 proposed for predicting the value of a binomially distributed response  
35 variable  $Y$ . Given a training set  $TS = \{(x, y) \in X \times Y \mid y = g(x)\}$   
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6  
7  
8 where  $X$  represents the search space spanned by  $m$  independent con-  
9 tinuous variables  $X_i$ , a logistic regression model  $M$  is induced by  
10 generalizing observations in  $TS$  in order to estimate the posterior  
11 class probability  $P(C_i | x)$  that any unlabeled example  $x \in X$  be-  
12 longs to  $C_i$  (one of the possible class labels). Hence, as outlined  
13 in [Appice (2008)], differently from the classical regression settings  
14 where the value of a response variable is directly predicted, in logistic  
15 regression the response to be predicted is the probability that an ex-  
16 ample belongs to a given class. This probability can then be used for  
17 classification purposes.  
18  
19  
20  
21  
22  
23  
24  
25

26 Model trees in which the linear regression functions take the form  
27 of logistic functions are called Logistic Model Trees (LMT). As ex-  
28 plained in [Landwehr (2005)], a logistic model tree consists of a stan-  
29 dard decision tree structure with logistic regression functions on the  
30 leaves, much like a model tree is a regression tree with regression  
31 functions on the leaves. As in standard decision trees, a test on one  
32 of the attributes is associated with every inner node. For a nominal  
33 attribute with  $k$  values, the node has  $k$  child nodes, and instances are  
34 sorted down one of the  $k$  branches depending on their value of the  
35 attribute. For numeric attributes, the node has two child nodes and  
36 the test consists of comparing the attribute value to a threshold: an  
37 instance is sorted down the left branch if its value for that attribute  
38 is smaller than the threshold and sorted down the right branch other-  
39 wise. In this study, we take the LMT learning algorithm introduced  
40 in [Landwehr (2005)] into account.  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

54 The LMT learning algorithm employs the LogitBoost algorithm  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 proposed in [Friedman (2000)] for building the logistic regression func-  
9 tions at the nodes of a tree. LogitBoost performs an additive logistic  
10 regression. In particular, as explained in [Friedman (2000)], at each  
11 iteration it fits a simple regression function by going through all the  
12 attributes, finding the simple regression function with the smallest er-  
13 ror, and adding it into the additive model. By running the LogitBoost  
14 algorithm until convergence, the result is a maximum likelihood mul-  
15 tiple logistic regression model [Witten (2016)]. However, to achieve  
16 a good generalization ability (i.e., optimum performance on unseen  
17 data) it is usually detrimental to wait for convergence. Hence, an  
18 appropriate number of iterations for the LogitBoost algorithm is de-  
19 termined by estimating the expected performance for a given number  
20 of iterations using cross-validation and stopping the process when no  
21 performance improvement is noticed.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

34 Logistic model trees are built by considering a simple extension  
35 of LogitBoost. As described in [Friedman (2000)], the boosting al-  
36 gorithm ends when no additional pattern in the data can be mod-  
37 eled by means of a linear logistic regression function. However, there  
38 may still be a pattern that linear models can fit by restricting the  
39 attention to a subset of the data. This subset can be obtained, for  
40 instance, by a standard decision tree criterion such as information  
41 gain [Witten (2016)]. Then, once no further improvement can be ob-  
42 tained by adding simpler linear models, the data are split and boosting  
43 is resumed separately in each subset. At each split, the logistic regres-  
44 sions of the parent node are passed to the child nodes. Hence, the  
45 logistic models generated so far are refined separately for the data in  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 each subset. Again, cross-validation is run in each subset to deter-  
9 mine an appropriate number of iterations to perform in that subset.  
10 The process is applied recursively until the subsets become too small.  
11 The final model in the leaf nodes accumulates all parent models and  
12 creates probability estimates for each class.  
13  
14  
15  
16

17 Finally, a pruning algorithm [Hastie (2009)] is applied to reduce  
18 the tree size and increase model's generalization. After the pruning  
19 operation, the algorithm produces small but very accurate trees with  
20 linear logistic models at the leaves.  
21  
22  
23  
24

25 The steps of the algorithm for building a LMT (also summarized  
26 in Figure 2) are the following:  
27  
28

- 29 1. Create a logistic regression model at root node. In this phase all  
30 the training observations are used for building an initial logistic  
31 regression model. The number of iterations (and simple regres-  
32 sion functions  $f_{mj}$  to add to  $F_j$ ) is determined using a five fold  
33 cross-validation. In particular, the number of iterations showing  
34 the lowest sum of errors is used to train the LogitBoost algo-  
35 rithm on all the data. This gives the logistic regression model at  
36 the root of the tree.  
37  
38  
39  
40  
41  
42  
43
- 44 2. Splitting step. When deciding which attribute to conduct a split  
45 on, the algorithm considers the C4.5 splitting criterion [Hastie (2009)].  
46 In particular, C4.5 chooses the attribute that maximizes the nor-  
47 malized information gain.  
48  
49  
50  
51
- 52 3. Tree growing. Tree growing continues in the following way: for  
53 each node resulting from the split just created, the logistic regres-  
54 sion function is refined based only on the subset of observations  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 that reached that node.

- 9  
10 4. Iterative step. The splitting process continues as described in the  
11 previous point, until more than 10 instances are at a node and  
12 a useful split can be found by using the C4.5 splitting criterion.  
13  
14 5. Model pruning. To increase the generalization ability of the  
15 model and to avoid overfitting, a pruning procedure is applied.  
16 The resulting model consists of small and accurate trees with  
17 linear logistic models on the leaves.  
18  
19  
20  
21  
22  
23

24 The reader is referred to [Landwehr (2005)] for a complete analysis  
25 of the LMT learning algorithm.  
26  
27  
28  
29

### 30 **3 Experimental Results**

31  
32  
33 To evaluate the performance obtained by the presented approach, we  
34 used the same data as in [Ackermann (2012)]. More specifically, we  
35 considered the 15 datasets from the DREAM5 systems genetics *in*  
36 *silico* network challenge A, in which the goal was to reconstruct gene-  
37 regulatory networks starting from (synthetic) genetic variations and  
38 gene expression data. Since the aim of this work is to show the im-  
39 provements achieved by the proposed machine learning technique, we  
40 measured the quality of its integrated results with respect to those  
41 obtained with individual prediction tools.  
42  
43  
44  
45  
46  
47  
48  
49

50 Each of these gene-regulatory networks was obtained by consider-  
51 ing 1000 markers, with each corresponding to a mutation of one of the  
52 1000 considered genes (corresponding to the nodes of the network),  
53 possibly having a different number of edges. Moreover, three different  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 sample sizes were considered: 100, 300, and 999; they correspond to  
9  
10 the three different sub-challenges, namely SysGenA100, SysGenA300,  
11  
12 and SysGenA999, respectively, and for each of these 5 different net-  
13  
14 works were generated. In the simulation process, the variations were  
15  
16 evenly distributed on 20 chromosomes by also considering the possibil-  
17  
18 ity of a local linkage between adjacent positions on the chromosomes.  
19  
20 More precisely, starting from homozygous recombinant inbred lines  
21  
22 (RILs), a randomized population was obtained by introducing ran-  
23  
24 dom mutations to simulate both *cis*- and *trans*-effects.

25  
26 As anticipated, the final goal of the DREAM5 SYSGEN A chal-  
27  
28 lenge was to perform an eQTL analysis, that is, to establish the rela-  
29  
30 tions in each regulatory network by exploiting the information of the  
31  
32 simulated gene expression levels of the 1000 genes and the simulated  
33  
34 genotype data. (eQTL mapping). To obtain the initial datasets, we  
35  
36 ran each of the aforementioned tools, namely, mRMR, MatrixEQTL in  
37  
38 both its *Linear* and *ANOVA* variants, and R/qtl in both its *maximum*  
39  
40 *likelihood* (EM) and *Haley-Knott* (HK) models, on the DREAM5 net-  
41  
42 works, with default parameters. Moreover, since the input datasets  
43  
44 involve both *cis*- and *trans*-eQTL, we did not impose any condition  
45  
46 (which could be done by setting specific tool parameters), such as the  
47  
48 maximum distance of the genomic positions, when running the afore-  
49  
50 mentioned tools to find the mappings. Anyway, it must be noticed that  
51  
52 there is no indication about the coordinates of the eQTLs, neither for  
53  
54 genes nor for SNPs. The predictions of these 5 tools/variants on the  
55  
56 15 networks constitute the initial datasets on which we performed the  
57  
58 experimental evaluation of our approach involving two steps: under-  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 sampling to balance the two classes, and combinations of the results  
9  
10 to improve the quality of the predictions. For each of the 15 datasets  
11 taken into account, the OSS sampling technique was used to produce  
12 a balanced dataset so that each resulting dataset (called a *reduced*  
13 *dataset*) has an equal number of positive and negative instances. In  
14 fact, in the datasets, for each SNP-gene pair the prediction values of  
15 the 5 tools/variants and a binary value indicating if the interaction is  
16 real or not (i.e. the *truth*) are reported.  
17  
18  
19  
20  
21  
22

23 We recall that the aim of this work is exactly to compare, inte-  
24 grate, and improve the results achieved by the different tools and we  
25 thus considered the DREAM5 *in silico* datasets for which the real  
26 predictions are known (ground truth).  
27  
28  
29

30 Before applying the LMT algorithm, we split the reduced dataset  
31 in such a way that 70% of the instances were used for training pur-  
32 poses, while the remaining 30% were used to assess the performance of  
33 the classifier on unseen instances. For each reduced dataset, consisting  
34 of the predictions of the 5 tools/variants, 30 different independent par-  
35 titions between training and test instances were performed. These 30  
36 partitions of the reduced dataset were created in order to take account  
37 of any bias introduced by the random sampling. Here, we stress that  
38 the training sets were used to train the logistic model which combines  
39 the outputs of the tools that were taken into account, and the test  
40 sets to assess its performance on unseen data. Hence, the considered  
41 tools were used only to produce values that are fed into the logistic  
42 model created, and no training was needed for these tools.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

54 After creating the training and test instances, the LMT algo-  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 rithm was conducted, considering the implementation provided by the  
9  
10 WEKA machine learning tool [Hall (2009)]. In particular, we used the  
11  
12 default parameters provided by WEKA, with the exception of the “fas-  
13  
14 tRegression” option that was used in the experiments we performed.  
15  
16 Use of this option in the LMT algorithm includes an heuristic that  
17  
18 avoids cross-validating the number of LogitBoost iterations at every  
19  
20 node. This allows the computational time needed to build the model  
21  
22 to be reduced, without significantly affecting the classifier’s final per-  
23  
24 formance.

25  
26 Statistical results, in terms of precision and recall, for both the  
27  
28 training and the test instances are obtained considering, for each of  
29  
30 the 15 datasets, the median over the 30 independent runs that were  
31  
32 performed. We preferred the median over the average for its higher  
33  
34 robustness to outliers. Table 1 summarizes all the results we obtained.

35  
36 As one can notice, the proposed system can achieve a good classifi-  
37  
38 cation performance on both training and test instances, hence showing  
39  
40 that it produces robust classifiers. In more in detail, all the median  
41  
42 values of precision and recall obtained on the 15 networks are higher  
43  
44 than 0.89 (except for the first network which has values around 0.76  
45  
46 for the recall). The overall median values on all networks computed in  
47  
48 the training phase are 1 for the precision and 0.931 for the recall. On  
49  
50 unseen instances (i.e. test phase) the values are 0.997 for the precision  
51  
52 and 0.916 for the recall. A graphic presentation of the results achieved  
53  
54 on all 30 runs on the 15 datasets is shown as a scatter plot in Figure 5,  
55  
56 in which the area from 0.5 to 1 for both axes is shown. Red dots show  
57  
58 the results obtained on the training sets, while light blue dots show  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 those obtained on the test sets.

9  
10 To better understand the improvements stemming from combining  
11 the predictions of the 5 tools/variants, we computed the precision and  
12 recall values for each of them on the 30 independent runs on the 15  
13 datasets. To identify the SNP-gene pairs predicted as positive by each  
14 tool/variant, we adopted as a threshold to split the two classes the  
15 default values of each tool as suggested in the corresponding paper:  
16 mRMR scores higher than 0,  $p$ -values of MatrixEQTL (both linear  
17 and ANOVA) lower than 0.05, and R/qtl scores (both EM and HK)  
18 higher than 2. By reporting the precision and recall values on each  
19 independent run of each dataset, one may observe in Figure 3 that  
20 the results are quite different and, in some cases, not so good. These  
21 scatter plots show for the 3 DREAM5 sample sizes (on the columns)  
22 and for each tested tool/variant (on the rows) the precision and recall  
23 values achieved on the 30 independent runs for each of the 5 networks  
24 on both the training and test sets. Note that while MatrixEQTL  
25 and R/qtl reveal a similar trend on the networks with more samples  
26 (SysGenA300 and SysGenA999), mRMR has a very low recall and  
27 also a precision often under 0.5. Moreover, as expected, the results  
28 achieved on the training instances are better than those obtained for  
29 the unseen ones (test instances).

30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47 To quantify the improvement in the proposed system's perfor-  
48 mance, we decided to apply our LMT system to each tool/variant  
49 separately. The aim of this analysis was to determine the best thresh-  
50 old for each tool/variant on every run of the tested datasets. The  
51 values of the computed thresholds reported in Table 2 differ in each  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6  
7  
8 dataset, highlighting the fact that it is hard to find default values with-  
9  
10 out considering the specific problem. The results shown in Figure 4  
11  
12 show a significant improvement in both precision and recall in almost  
13  
14 all tested tools/variants, emphasizing the importance of the choice of  
15  
16 the thresholds. In fact, although for each tool the authors suggest the  
17  
18 ideal threshold value(s), each dataset varies from the others and so it  
19  
20 is important to select the threshold values accordingly.

21  
22 From this point of view, our method not only exploits the infor-  
23  
24 mation provided by the different tools to achieve a better prediction,  
25  
26 but is unconstrained by the particular threshold values. Moreover,  
27  
28 these results show the the proposed method's suitability for address-  
29  
30 ing the problem under scrutiny: in particular, in the vast majority of  
31  
32 datasets, the training and test performance are comparable. Hence,  
33  
34 the method can not only extract a model of the data that produces a  
35  
36 good classification of the training instances, but it also shows a good  
37  
38 generalization ability.

39  
40 This is also supported by the experimental analysis we performed  
41  
42 on a real case study, consisting of a dataset of miRNA expression of  
43  
44 30 *C. elegans* recombinant in-bred lines [Kel (2016)]. These lines were  
45  
46 obtained from the crosses between two different *C. elegans* strains: the  
47  
48 N2 wild-type ancestral worms and the wild CB4856 nematodes [Hodgkin and Doniach (1997)],  
49  
50 isolated from Hawaii (HW). All individuals were genotyped across  
51  
52 1455 SNP markers [Rockman (2010)], while the miRNA expression  
53  
54 profiles were estimated using the small RNA sequencing technology of  
55  
56 Illumina [Kel (2016)].

57  
58 Table 3 gives an overview of the experimentally validated eQTL  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 results, which were all found by the proposed machine learning ap-  
9  
10 proach. These predictions, concerning novel putative eQTLs, originate  
11 from the comparison between the HW and N2 strains of *C. elegans*  
12 and were experimentally validated. To this end, the effective miRNA  
13 differential expression was quantified on the 16 RILs that were most  
14 representative of the original set of 30, by using a Rotor Gene Real-  
15 Time PCR System (Qiagen) with Sybr Green miRNA assays (Qiagen)  
16 in accordance with the manufacturer's instructions.  
17  
18  
19  
20  
21  
22

23 The predictions reported in Table 3 are experimentally confirmed,  
24 since the identified fold change supports our results considering the  
25 widely adopted 1.5 threshold. Moreover, as discussed in [Kel (2016)],  
26 none of the tools considered in our integration approach was able alone  
27 to predict all these eQTLs, while using our integration approach all  
28 the validated eQTLs are correctly predicted. This proves the superior  
29 performance of our machine learning approach, and that this method  
30 is generally applicable also to real datasets. We would like to point out  
31 that although our machine learning approach bases its classification  
32 process on a function, its value can not be used to score the predictions  
33 since it does not correspond to a quality measure. For this reason  
34 we only classified the instances and verified that the experimentally  
35 validated ones were correctly predicted.  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49

## 50 4 Conclusions

51  
52 In this work, we proposed a supervised learning approach to integrate  
53 eQTL predictions provided by different tools. In particular, we com-  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 bined two approaches, linear models and tree induction methods, that  
9  
10 are very popular in the data mining community for solving a super-  
11 vised learning problem.  
12

13       Linear model approaches fit a simple (linear) model to the data,  
14 with this process being quite stable, resulting in low variance, but  
15 potentially high bias. Tree induction methods exhibit low bias but  
16 often high variance: they search a less restricted space of models,  
17 allowing it to capture nonlinear patterns in the data, but making it  
18 less stable and prone to overfitting. Not surprisingly, neither method  
19 is generally superior. Much work has been done to combine these two  
20 schemes into model trees, namely, trees containing linear regression  
21 functions on the leaves. Hence, if only little and/or noisy data are  
22 available model trees rely on simple regression models, while if there  
23 are enough data more complex tree structures are used.  
24  
25

26       Moreover, since real SNP-gene interactions correspond to a small  
27 fraction of all the possible combinations ( $\sim 0.2\%$ ), another critical  
28 problem addressed in this work is that the set of elements of a spe-  
29 cific class (positive results) is very small relative to the overall dataset  
30 (whole predictions). The predictions obtained with most of the clas-  
31 sifiers are influenced by the fact that the classes are not equally bal-  
32 anced. These approaches assume that all misclassification errors come  
33 at equal cost, while in several applications this assumption may not  
34 be true. To avoid this problem, we used a method called One Side Se-  
35 lection, which undersamples the class containing the highest number  
36 of elements.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

54       The results obtained using our Logistic Model Trees on the DREAM5  
55  
56  
57

1  
2  
3  
4  
5  
6  
7  
8 challenge datasets are encouraging, showing we are always able to  
9  
10 achieve better results than the original tools. This is achieved by  
11  
12 combining the predictions obtained by each tool, but also by com-  
13  
14 puting the “best” threshold values for each of them on the analyzed  
15  
16 dataset. For this reason, we plan to extend our integrative approach  
17  
18 to other software to further reduce the number of false positive results,  
19  
20 as well as to apply the presented machine learning technique to other  
21  
22 datasets to assess its power on real case studies.  
23  
24

## 25 **Acknowledgments**

26  
27  
28 This work has been supported by the Italian Ministry of Education  
29  
30 and Research (MIUR) through the Flagship (PB05) InterOmics and  
31  
32 HIRMA (RBAP11YS7K), and by the FP7-HEALTH through the Eu-  
33  
34 ropean MIMOMICS projects (no. 305280).  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## References

- [Abecasis (2002)] Abecasis, G.R., Cherny, S.S., Cookson, W.O., et al. 2002. Merlinrapid analysis of dense genetic maps using sparse gene flow trees. *Nature Genetics*, 30(1), 97–101.
- [Ackermann (2012)] Ackermann, M., Clment-Ziza, M., Michaelson, J.J., et al. 2012. Teamwork: improved eQTL mapping using combinations of machine learning methods. *PloS One*, 7(7), e40916.
- [Appice (2008)] Appice, A., Ceci, M., Malerba, D., et al. 2008. Stepwise induction of logistic model trees. *Foundations of Intelligent Systems*, 68–77.
- [Beretta (2016)] Beretta, S., Castelli, M., Goncalves, I., et al. 2016. Combining Bayesian Approaches and Evolutionary Techniques for the Inference of Breast Cancer Networks, 217–224. In *Proceedings of the 8th International Joint Conference on Computational Intelligence (IJCCI)*. SciTePress.
- [Beretta (2017)] Beretta, S., Maj, C., and Merelli, I. 2017. Rank miRNA: a web tool for identifying polymorphisms altering miRNA target sites, 1125–1134. In *Proceedings of International Conference on Computational Science (ICCS)*. Procedia Computer Science.
- [Bishop (2006)] Bishop, C.M. 2006. *Pattern Recognition and Machine Learning*. Springer-Verlag, New York.
- [Broman (2003)] Broman, K.W., Wu, H., Sen, Ś., et al. 2003. R/qtl: QTL mapping in experimental crosses. *Bioinformatics*, 19(7), 889–890.

- 1  
2  
3  
4  
5  
6  
7  
8 [Caruana and Niculescu-Mizil (2006)] Caruana, R., and Niculescu-  
9 Mizil, A. 2006. An empirical comparison of supervised learn-  
10 ing algorithms, 161–168. In *Proceedings of the 23rd international*  
11 *conference on Machine learning*. ACM, New York.  
12  
13  
14  
15  
16 [Duggal (2014)] Duggal, G., Wang, H., and Kingsford, C. 2014.  
17 Higher-order chromatin domains link eQTLs with the expression  
18 of far-away genes. *Nucleic Acids Research*, 42(1), 87–96.  
19  
20  
21  
22 [Friedman (2000)] Friedman, J., Hastie, T., Tibshirani, R. et al. 2000.  
23 Additive logistic regression: a statistical view of boosting (with  
24 discussion and a rejoinder by the authors). *The annals of statis-*  
25 *tics*, 28(2), 337–407.  
26  
27  
28  
29  
30 [Ganganwar (2012)] Ganganwar, V. 2012. An overview of classifica-  
31 tion algorithms for imbalanced datasets. *International Journal*  
32 *of Emerging Technology and Advanced Engineering*, 2(4), 42–47.  
33  
34  
35  
36 [Guo (2008)] Guo, X., Yin, Y., Dong, C., Yang, et al. 2008. On the  
37 class imbalance problem, 192–201. In Guo, M., Zhao, L., and  
38 Wang L. eds. *Natural Computation, ICNC’08. Fourth Interna-*  
39 *tional Conference on*. IEEE.  
40  
41  
42  
43  
44 [Hall (2009)] Hall, M., Frank, E. Holmes, G. et al. 2009. The WEKA  
45 data mining software: an update. *ACM SIGKDD Explorations*  
46 *Newsletter*, 11(1), 10–18.  
47  
48  
49  
50 [Harrell Jr (2015)] Harrell Jr, F.E. 2015. *Regression Modeling Strate-*  
51 *gies: With Applications to Linear Models, Logistic and Ordinal*  
52 *Regression, and Survival Analysis*. Springer, New-York.  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8 [Hart (1968)] Hart, P. 1968. The condensed nearest neighbor rule  
9 (Corresp.). *IEEE Transactions on Information Theory*, 14(3),  
10 515–51.  
11  
12  
13  
14 [Hastie (2009)] Hastie, T., Tibshirani, R., and Friedman, J. 2009.  
15 *Overview of Supervised Learning*. Springer-Verlag, New York.  
16  
17  
18 [Hodgkin and Doniach (1997)] Hodgkin, J., and Doniach, T. 1997.  
19 Natural variation and copulatory plug formation in *Caenorhab-*  
20 *ditis elegans*. *Genetics*, 146(1), 149–164.  
21  
22  
23  
24 [Hoggart (2008)] Hoggart, C.J., Whittaker, J.C., De Iorio, M., et al.  
25 2008. Simultaneous analysis of all SNPs in genome-wide and re-  
26 sequencing association studies. *PLoS Genetics*, 4(7), e1000130.  
27  
28  
29  
30 [Hosmer Jr (2013)] Hosmer Jr, D.W., Lemeshow, S., and Sturdivant,  
31 R.X. 2013. *Applied Logistic Regression*. John Wiley & Sons.  
32  
33  
34 [Huang and Cai (2013)] Huang, T., and Cai, Y.D. 2013. An  
35 information-theoretic machine learning approach to expression  
36 QTL analysis. *PloS One*, 8(6), e67899.  
37  
38  
39  
40 [Japkowicz and Stephen (2002)] Japkowicz, N., and Stephen, S. 2002.  
41 The class imbalance problem: A systematic study. *Intelligent*  
42 *Data Analysis*, 5(6), 429–449.  
43  
44  
45  
46 [Kao (1999)] Kao, C.H., Zeng, Z.B., Teasdale, R.D. 1999. Multiple  
47 interval mapping for quantitative trait loci. *Genetics*, 152(3),  
48 1203–1216.  
49  
50  
51  
52 [Kel (2016)] Kel, I., Chang, Z., Galluccio, N., et al. 2016. SPIRE, a  
53 modular pipeline for eQTL analysis of RNA-Seq data, reveals a  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8 regulatory hotspot controlling miRNA expression in *C. elegans*.  
9 *Molecular BioSystems*, 12(11), 3447–3458.

10  
11  
12 [Kubat and Matwin (1997)] Kubat, M., and Matwin, S. 1997. Ad-  
13 dressing the curse of imbalanced training sets: one-sided selec-  
14 tion, 179–186. In Fisher, D.M. eds. *International Conference on*  
15 *Machine Learning*, Morgan Kaufmann, Burlington.

16  
17  
18 [Landwehr (2005)] Landwehr, N., Hall, M., Frank, E. 2005. Logistic  
19 model trees. *Machine Learning*, 59(1-2), 161–205.

20  
21  
22 [Laurikkala (2001)] Laurikkala, J. 2001. Improving identification of  
23 difficult small classes by balancing class distribution. *Artificial*  
24 *Intelligence in Medicine*, 63–66.

25  
26  
27 [Lee (2008)] Lee, S.H., van der Werf, J.H., Hayes, B.J., et al. 2008.  
28 Predicting unobserved phenotypes for complex traits from whole-  
29 genome SNP data. *PLoS Genetics*, 4(10), e1000231.

30  
31  
32 [Leek and Storey (2007)] Leek, J.T., and Storey, J.D. 2007. Captur-  
33 ing heterogeneity in gene expression studies by surrogate variable  
34 analysis. *PLoS Genetics*, 3(9), e161.

35  
36  
37 [Loh (2011)] Loh, W.Y. 2011. Classification and regression trees. *Wi-*  
38 *ley Interdisciplinary Reviews: Data Mining and Knowledge Dis-*  
39 *covery*, 1(1), 14–23.

40  
41  
42 [Merelli (2013)] Merelli, I., Calabria, A., Cozzi, P., et al. 2013.  
43 SNPranker 2.0: a gene-centric data mining tool for diseases asso-  
44 ciated SNP prioritization in GWAS. *BMC Bioinformatics*, 14(1),  
45 S9.



- 1  
2  
3  
4  
5  
6  
7  
8 [Merelli (2015)] Merelli, I., Tordini, F., Drocco, M., et al. 2015. In-  
9 tegrating multi-omic features exploiting Chromosome Conforma-  
10 tion Capture data. *Frontiers in Genetics*, 40, 1664–8021.  
11  
12 [Michaelson (2010)] Michaelson, J.J., Alberts, R., Schughart, K., et  
13 al. 2010. Data-driven assessment of eQTL mapping methods.  
14 *BMC Genomics*, 11(1), 502.  
15  
16 [Nguyen (2009)] Nguyen, G.H., Bouzerdoum, A., Phung, S. 2009.  
17 Learning pattern classification tasks with imbalanced data sets,  
18 193–208. In Yin, P. eds. *Pattern recognition*, Vukovar, Croatia:  
19 In-Teh.  
20  
21 [Rockman (2010)] Rockman, M.V., Skrovanek, S.S., and Kruglyak,  
22 L. 2010. Selection at linked sites shapes heritable phenotypic  
23 variation in *C. elegans*. *Science*, 330(6002), 372–376.  
24  
25 [Shabalín (2012)] Shabalín, A.A. 2012. Matrix eQTL: ultra fast  
26 eQTL analysis via large matrix operations. *Bioinformatics*,  
27 28(10):1353–1358.  
28  
29 [Servin and Stephens (2007)] Servin, B., and Stephens, M. 2007.  
30 Imputation-based analysis of association studies: candidate re-  
31 gions and quantitative traits. *PLoS Genetics*, 3(7), e114.  
32  
33 [Sumner (2005)] Sumner, M., Frank, E., and Hall, M. 2005. Speeding  
34 up logistic model tree induction, 675–683. In A. Jorge et al.  
35 eds. *Knowledge Discovery in Databases: PKDD 2005*. Springer-  
36 Verlag, Berlin.  
37  
38 [Tomek (1976)] Tomek, I. 1976. Two modifications of CNN. *IEEE*  
39 *Trans. Systems, Man and Cybernetics*, 6, 769–772.  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8 [Witten (2016)] Witten, I.H. , Eibe, F., Hall, M.A. et al. 2016. *Data*  
9 *Mining: Practical Machine Learning Tools and Techniques*. Mor-  
10 gan Kaufmann, Burlington.  
11  
12  
13  
14 [Wright (2012)] Wright, F.A., Shabalin, A.A., and Rusyn, I. 2012.  
15 Computational tools for discovery and interpretation of expres-  
16 sion quantitative trait loci. *Pharmacogenomics*, 13(3), 343–352.  
17  
18  
19  
20 [Zeng (1993)] Zeng, Z.B. 1993. Theoretical basis for separation of  
21 multiple linked gene effects in mapping quantitative trait loci.  
22 *Proceedings of the National Academy of Sciences*, 90(23), 10972–  
23 10976.  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Table 1: Classification performance on the training and test instances for datasets considered. Averages values of Precision and Recall over the 30 independent runs are reported for each of the 15 datasets.

Dataset	Train		Test		
	Precision	Recall	Precision	Recall	
DREAM5 SysGenA100	Network1	0.971	0.762	0.952	0.765
	Network2	1.000	0.952	0.993	0.949
	Network3	1.000	0.971	0.999	0.968
	Network4	1.000	0.971	1.000	0.971
	Network5	0.999	0.963	0.997	0.959
DREAM5 SysGenA300	Network1	0.999	0.918	0.997	0.916
	Network2	1.000	0.904	1.000	0.910
	Network3	0.999	0.916	0.998	0.911
	Network4	1.000	0.897	1.000	0.897
	Network5	1.000	0.902	0.998	0.901
DREAM5 SysGenA999	Network1	1.000	0.952	0.997	0.949
	Network2	1.000	0.953	0.996	0.950
	Network3	0.998	0.948	0.991	0.944
	Network4	0.959	0.914	0.904	0.895
	Network5	0.999	0.931	0.997	0.931

Table 2: Threshold values obtained with the machine learning classifier on each single tool. For each dataset we reported the threshold computed with the LMT classifier by considering the predictions of each tool/variant independently. Default values for each prediction tool/variant are shown in parentheses.

		MatrixEQTL		mRMR	R/qt1	
Dataset		ANOVA (0.05)	Linear (0.05)	(0)	EM (2)	HK (2)
DREAM5 SysGenA100	Network1	0.000179	0.000195	0.111	1.92866	2.36070
	Network2	0.000185	0.000240	0.108	4.87492	5.47845
	Network3	0.000169	0.000135	0.125	0.12500	0.10672
	Network4	0.000184	0.000223	0.118	0.00011	0.08902
	Network5	0.000223	0.009998	0.152	1.67575	1.47080
DREAM5 SysGenA300	Network1	0.000144	0.000136	0.118	6.23438	7.07576
	Network2	0.000150	0.000279	0.040	5.21553	6.20687
	Network3	0.000178	0.000175	0	5.27126	5.81573
	Network4	0.000247	0.000215	0.033	4.50471	5.16540
	Network5	0.000246	0.000332	0.029	4.25885	0.02900
DREAM5 SysGenA999	Network1	0.000143	0.000163	0.159	8.66660	10.83790
	Network2	0.000219	0.000137	0.067	4.66071	5.11208
	Network3	0.000123	0.000253	0.038	7.13816	7.81666
	Network4	0.000220	0.000281	0.015	8.32116	3.88823
	Network5	0.000344	0.000310	0.153	6.44104	7.26180

Table 3: Experimentally validated eQTL predictions for the *C. elegans* dataset. For each prediction miRNA and QTL coordinates are reported, and also the support interval of the miRNA. Predictions are sorted according to the expression fold change calculated using ddCT values on the results of qPCR validation.

Id	miRNA		QTL		Support Interval	Exp. Fold Change
	Chr	Pos	Chr	Pos		
mir-799	X	8600630	X	8665271	[124, 132]	12.441
mir-8201	4	6294591	4	6544139	[969, 978]	3.88
lin-4	2	5902266	5	16623881	[1390, 1390]	3.06
mir-4932	1	9512352	5	6259748	[1238, 1239]	2.115
mir-787	X	11294680	1	14267212	[457, 462]	1.905
mir-4936	3	3249194	1	1563141	[271, 280]	1.857
mir-242	4	4274287	4	4265784	[944, 945]	1.848
mir-793	X	13857930	1	11722283	[419, 424]	1.823
mir-357	5	8580573	4	13532205	[1078, 1078]	1.61

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

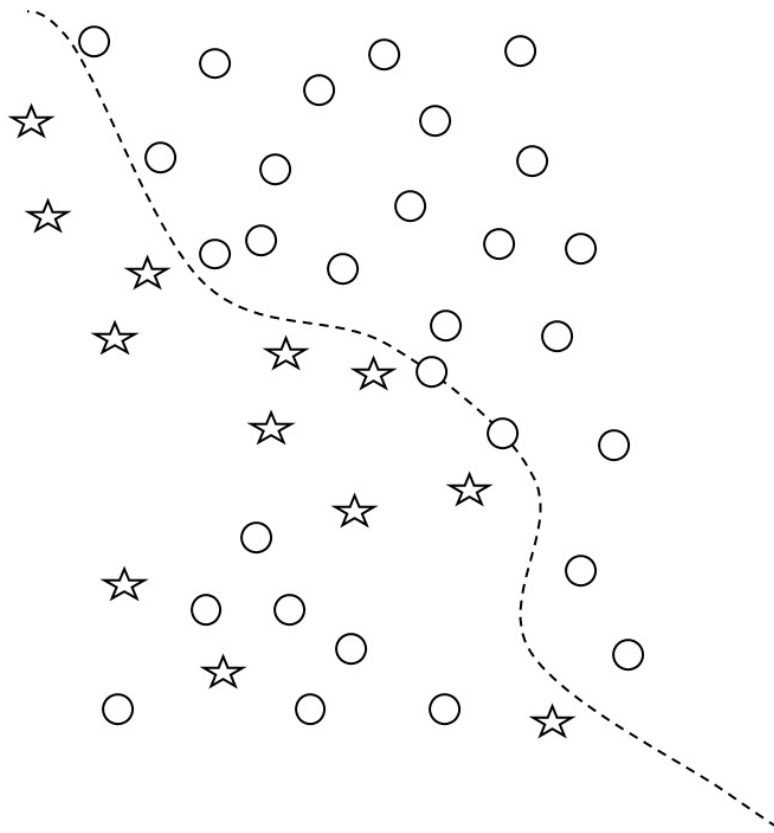


Figure 1: Example of the distribution of instances in a two-class dataset: stars denote instances in the minority class, circles denote instances in the majority class.

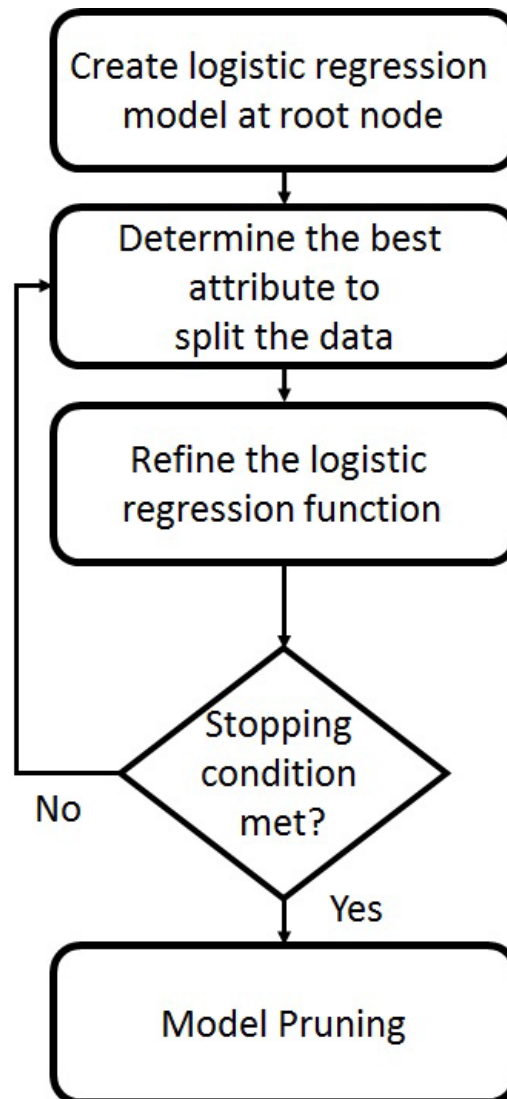


Figure 2: Graphical presentation of the LMT learning algorithm.

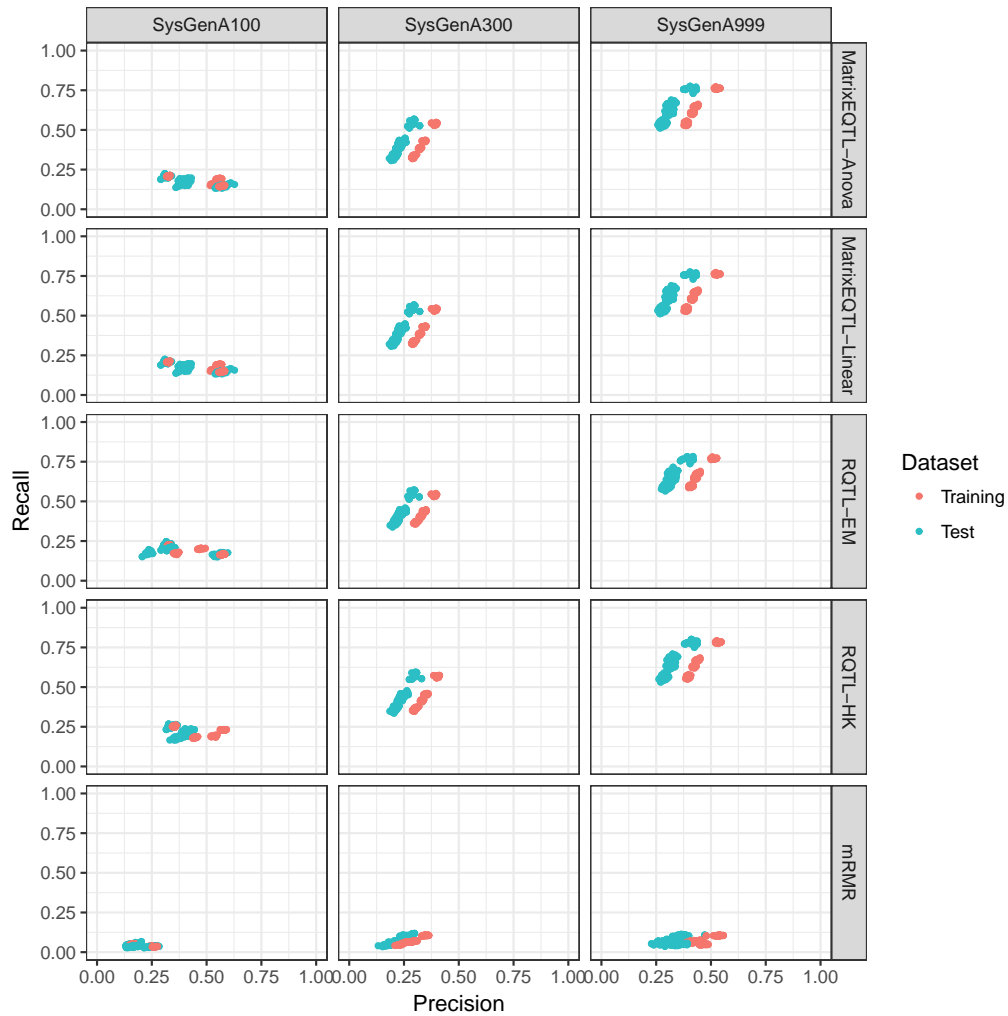


Figure 3: Scatter plots showing the precision and recall achieved by each considered tool/variant with *suggested* threshold values (rows) on the 30 independent runs of the 15 DREAM5 datasets, split by sample size (columns). Colors highlight the two sets: Training (red) and Test (light blue).



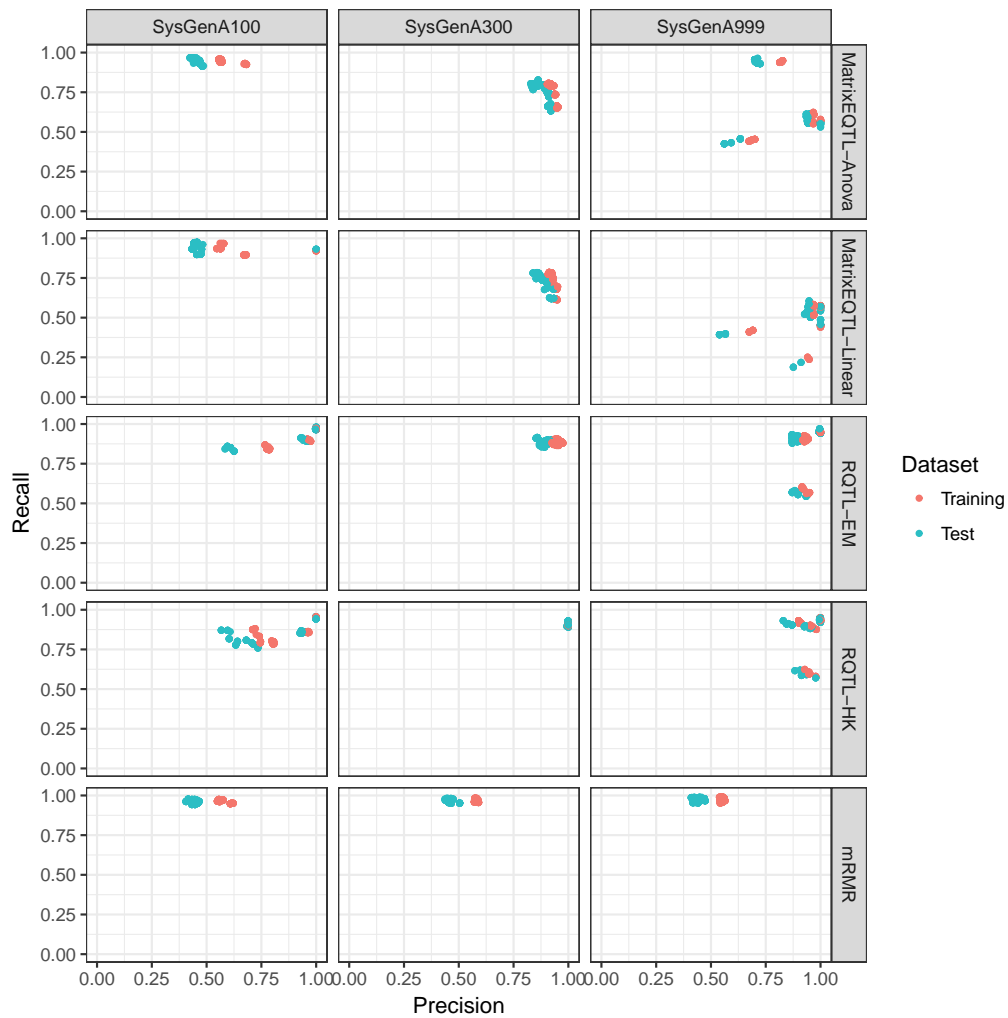


Figure 4: Scatter plots showing the precision and recall achieved by each considered tool/variant with *optimal* threshold values (rows) on the 30 independent runs of the 15 DREAM5 datasets, split by sample size (columns). Colors highlight the two sets: Training (red) and Test (light blue).

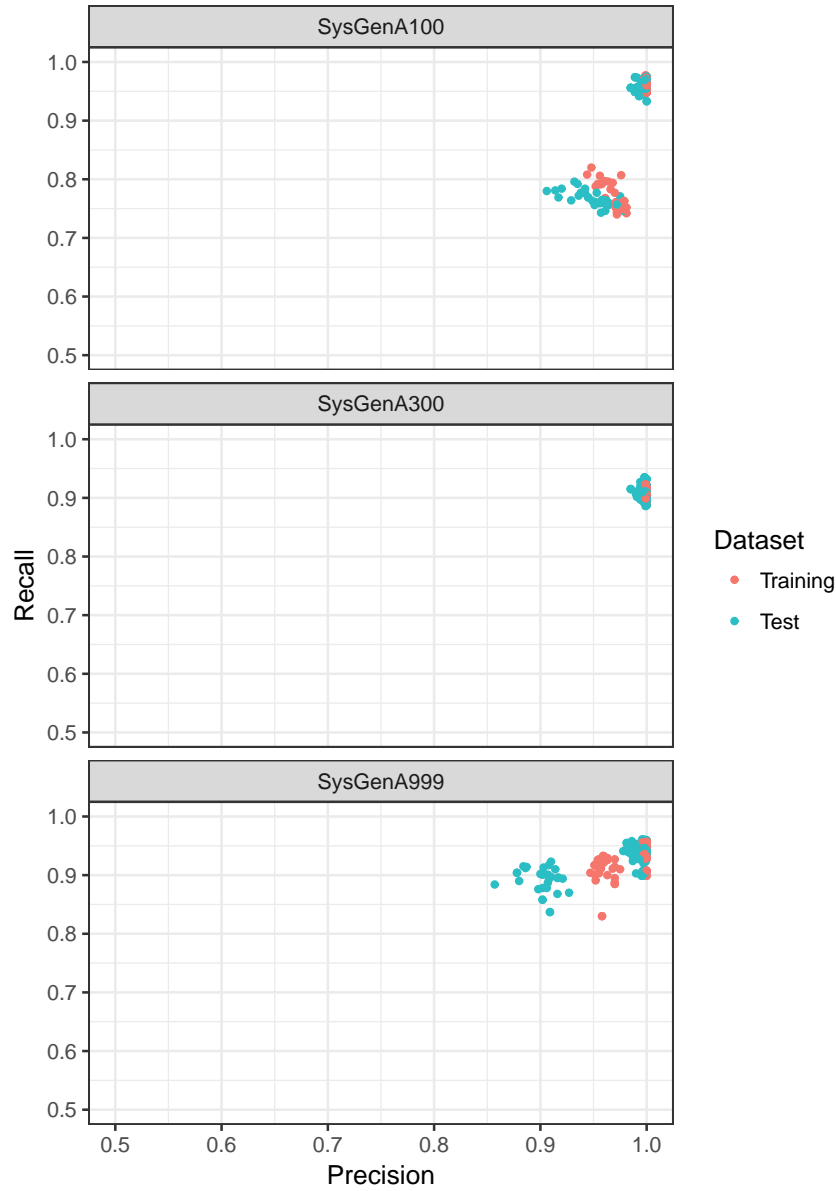


Figure 5: Scatter plots showing the Precision and Recall achieved by our approach on the 30 independent runs of the 15 DREAM5 datasets, split by sample size (rows), in both Training (red) and Test (light blue) sets. Here, only the area from 0.5 to 1 on both axes is shown.