

Computational intelligence and asthma

Pablo Vicente-Munuera

Bioinformatics MSc

Email: pablovm1990@gmail.com

Abstract—The aim of this work is apply some algorithms typical used in statistics and machine learning to a given dataset. With this, we'll try to find if some of these algorithms can classify our dataset and try to figure something out. The dataset selected for this purpose is an asthma dataset.

I. INTRODUCTION

Asthma is a chronic (long-lasting) inflammatory disease of the airways. In those susceptible to asthma, this inflammation causes the airways to spasm and swell periodically so that the airways narrow. The individual then must wheeze or gasp for air. Obstruction to air flow either resolves spontaneously or responds to a wide range of treatments, but continuing inflammation makes the airways hyper-responsive to stimuli such as cold air, exercise, dust mites, pollutants in the air, and even stress and anxiety.

A few keys from AAAAI (American Academy of Allergy Asthma and Immunology)[3] who collects all of this information:

- 1) From 2001 through 2009 asthma rates rose the most among black children, almost a 50% increase, in U.S.A. [4]
- 2) More than half (53%) of people with asthma had an asthma attack in 2008. More children (57%) than adults (51%) had an attack. 185 children and 3,262 adults died from asthma in 2007, in U.S.A. [4]
- 3) An estimated 300 million people worldwide suffer from asthma, with 250,000 annual deaths attributed to the disease.[5]
- 4) About 70% of asthmatics also have allergies.[5]
- 5) It is estimated that the number of people with asthma will grow by more than 100 million by 2025. [5]

But it's not just a problem of the (illogically) called "first world". It also affects to other undeveloped continents such as the African one [6]:

- In 2010, 49.7 million (13.9%; 95% CI 9.6-18.3) among children ;15 years, 102.9 million (13.8%; 95% CI 6.2-21.4) among people aged ;45 years, and 119.3 million (12.8%; 95% CI 8.2-17.1) in the total population.

A. Dataset

The dataset chosen is the one used by Voraphani N. (2014)[1], in which the main purpose was to identify differentially expressed genes. These genes belongs to different subjects with different classes:

- Control: Subjects with no asthma.
- MMA: Subjects with mild-moderate asthma.

- SA: severe asthmatic patients.

All of these classes represent an expression array with bronchial epithelial cells of Homo Sapiens.

Because we get an expression file, we had to convert it to an arff file (used by weka [9]). In order to get this file, we have to do some operations before we could some research about it. The first thing, and the most important thing, you must do when you get an expression file, is the normalization of the data. With R and bioconductor [?] (citation here), this purpose is solved. We do the normalization with the background of the array with R. We could, also, have done differential expression in order to obtain just a few genes (differentially expressed over the others).

After all of this, we obtain 43377 (without array controls) genes and 108 subjects. With this dataset we'll do all of the classification's problems.

II. METHODS

As said before, the main purpose is classify well and see which genes could contribute more than others We'll analyze these genes so as to see if any of those important genes to our classification problem, have the same relevancy in biology, and, being more specific, in asthma.

Why machine learning & statistics and not other methods? Well, machine learning and statistics are so powerful. That could be a double-edged sword, but in this case it is not. Although, machine learning, by now, is widely used in the field of bioinformatics and biology.

A. Supervised classification

In a classification problem, we have a set of elements divided into classes [7]. Given an element (or instance) of the set, a class is assigned according to some of the element's features and a set of classification rules. In our case, we have three classes (Control, MMA, SA) and 108 instances (subjects with asthma or not). Our subjects are labelled with their own class. So we proceed to divide the dataset in two subsets: the training dataset and the test dataset. The training one will be the input (labelled as well) of the classifier. The classifier learn to classify this training dataset, and the output will be a model. With this model, we run again the classifier, but, at this time, we input the test dataset without labels (i.e. no classes). After all of this procedure, we'll obtain a percentage of how good is our classifier.

In order to reduce the bias with the division of the dataset (in training and test), we'll execute the k-fold-cross-validation [8]. In this case, the dataset is partitioned into k folds. Each

fold is left out of the design process and used as a testing set. The estimate of the error is the overall proportion of the errors committed on all folds.

We've used several paradigms (classifiers) to see which could be fitted more to our data. In the next subsections, they are going to be explained.

1) *Naive Bayes*:

2) *KNN, K-nearest neighbours (IBk)*: The nearest-neighbour rule [75] to classify x is to assign it to the label associated with the prototype nearest to the test point (Figure 6). An obvious extension of the nearest-neighbour rule is the k -nearest-neighbour rule. This rule classifies x by assigning it to the label most frequently represented among the k nearest samples. In other words, a decision is made by examining the labels on the k -nearest-neighbours and voting. A practical problem with this simple method is that it tends to be slow for large training sets because the entire set must be searched for each test instance. A strategy to avoid the computational complexity of the nearest neighbour algorithm is to classify each example with respect to the examples already seen and to save only those that are misclassified. This strategy is known as condensing.

3) *Decision tree: C4.5 (j48)*:

4) *Logistic regression*:

$$P(C = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \sum_{i=1}^n \beta_i x_i)}} \quad (1)$$

Where x represents an instance to be classified, and $\beta_0, \beta_1, \dots, \beta_n$ are the parameters of the model. These parameters should be estimated from the data in order to obtain a concrete model.

5) *Bayesian Networks, TAN*:

B. Complementary

book wittenyfrank.pdf

1) *Random forest*: RandomForest constructs random forests by bagging ensembles of random trees

2) *Adaboost*: There are many variants on the idea of boosting. We describe a widely used method called AdaBoost.M1 that is designed specifically for classification. Like bagging, it can be applied to any classification learning algorithm. To simplify matters we assume that the learning algorithm can handle weighted instances, where the weight of an instance is a positive number. (We revisit this assumption later.) The presence of instance weights changes the way in which a classifier's error is calculated: it is the sum of the weights of the misclassified instances divided by the total weight of all instances, instead of the fraction of instances that are misclassified.

3) *Multilayer perceptron*: Section 4.6 explained that a perceptron represents a hyperplane in instance space. We mentioned there that it is sometimes described as an artificial neuron. Of course, human and animal brains successfully undertake very complex classification tasks for example, image recognition. The functionality of each individual neuron in a brain is certainly not sufficient to perform these feats. How can

they be solved by brain-like structures? The answer lies in the fact that the neurons in the brain are massively interconnected, allowing a problem to be decomposed into subproblems that can be solved at the neuron level. This observation inspired the development of networks of artificial neurons neural nets. Consider the simple datasets in Figure 6.10. Figure 6.10(a) shows a two-dimensional instance space with four instances that have classes 0 and 1, represented by white and black dots, respectively. No matter how you draw a straight line through this space, you will not be able to find one that separates all the black points from all the white ones. In other words, the problem is not linearly separable, and the simple perceptron algorithm will fail to generate a separating hyperplane (in this two-dimensional instance space a hyperplane is just a straight line). The situation is different in Figure 6.10(b) and Figure 6.10(c): both these problems are linearly separable. The same holds for Figure 6.10(d), which shows two points in a one-dimensional instance space (in the case of one dimension the separating hyperplane degenerates to a separating point). If you are familiar with propositional logic, you may have noticed that the four situations in Figure 6.10 correspond to four types of logical connectives. Figure 6.10(a) represents a logical XOR, where the class is 1 if and only if exactly one of the attributes has value 1. Figure 6.10(b) represents logical AND, where the class is 1 if and only if both attributes have value 1. Figure 6.10(c) represents OR, where the class is 0 only if both attributes have value 0. Figure 6.10(d) represents NOT, where the class is 0 if and only if the attribute has value 1. Because the last three are linearly separable, a perceptron can represent AND, OR, and NOT. Indeed, perceptrons for the corresponding datasets are shown in Figure 6.10(f) through (h) respectively. However, a simple perceptron cannot represent XOR, because that is not linearly separable. To build a classifier for this type of problem a single perceptron is not sufficient: we need several of them.

4) *RBF Network*: Other kernel functions can be used instead to implement different nonlinear mappings. Two that are often suggested are the radial basis function (RBF) kernel and the sigmoid kernel. Which one produces the best results depends on the application, although the differences are rarely large in practice. It is interesting to note that a support vector machine with the RBF kernel is simply a type of neural network called an RBF network (which we describe later)

C. Filter selection subset

Bioinformatics-2007-Saeys.pdf

1) *Univariate*:

2) *Multivariate*:

3) *Wrapper*:

III. RESULTS

In this section we'll show all the results we get from the expression set used. The software used to get all of those results is weka [9]. And all of the results are compared by the ROC curve [10], which is a trusty measure.

TABLE I

*: THE X VALUES MEANS NO RESULTS COULD BE OBTAINED DUE TO COMPUTATIONAL PROBLEMS. **THE K USED HERE IS EQUAL TO ONE.

Method	No filter	FSS		
		text	Univariate	Wrapper
Naive Bayes	0.436		0.626	0.623
KNN**	0.594		0.644	0.67
Logistic	X*		0.731	0.7
Decision Tree	0.486		0.639	0.698
Bayesian Net	X*		0.657	0.663

We had these problems mentioned before because the initial dataset had, approximately, 43300 variables. One of the first approaches to solve this, was a first univariate filter in which we selected the information gain ratio with respect to their own class. After this filter, all of these classify problems became more time-relaxed, due to the reductio of the number of variables (43300 to 32).

TABLE II

*: K EQUALS TO 4; ** K EQUALS TO 1.

Method	Multivariate		
	CFS Best first	CFS Genetic	Wrapper
Naive Bayes	0.623	0.599	0.74
KNN	0.773*	0.674*	0.687**
Logistic	0.699	0.675	0.824
Decision Tree	0.668	0.76	0.641
Bayesian Net	0.667	0.674	0.678

TABLE III

RESULTS DEPENDING THE K VALUE ON THE KNN CLASSIFIER

K	Univariate	Multivariate
1	0.644	0.67
2	0.735	0.72
3	0.73	0.745
4	0.744	0.773
5	0.762	0.807
8	0.77	0.771

TABLE IV

BEST RESULTS ON THE COMPLEMENTARY METHODS.

Method	No filter	FSS		Filtered by gain ratio		
		Univariate	CFS	CFS	CFS (genetic)	Wrapper
Adaboost (random forest)	0.573	0.801	0.795	0.798	0.829	0.835
Adaboost (J48)	0.544	0.758	0.736	0.736	0.8	X
Multilayer perceptron	X*	0.788	0.792	0.778	0.697	0.83
RBF Network	0.393	0.675	0.726	0.726	0.693	X
Random forest	0.569	0.811	0.783	0.803	0.809	0.843

Another important thing which showed up with the output of the results was the classification between one class and the other ones.

IV. DISCUSSION

REFERENCES

- [1] Voraphani N, Gladwin MT, Contreras AU, Kaminski N et al. An airway epithelial iNOS-DUOX2-thyroid peroxidase metabolome drives

TABLE V

SOME EXAMPLES OF THE DIFFERENCE BETWEEN CLASSES. MEASURED BY THE ROC CURVE.

Method - Filter	Classes			Overall
	Control	MMA	SA	
KNN (k=1) - Gain ratio	0.642	0.589	0.717	0.644
Bayesian Net - Gain ratio & CFS	0.802	0.591	0.697	0.667
Logistic - Wrapper	0.82	0.755	0.829	0.793
Decision Tree - No filter	0.55	0.473	0.471	0.486
Naive Bayes - CFS Best first	0.766	0.534	0.664	0.623

- Th1/Th2 nitrate stress in human severe asthma. *Mucosal Immunol* 2014 Sep;7(5):1175-85.
- [2] "Asthma." MedlinePlus. January 16, 2009 [cited January 20, 2009]. <http://www.nlm.nih.gov/medlineplus/asthma.html> .
- [3] American Academy of Allergy, Asthma, and Immunology (AAAAI) 555 East Wells Street, Suite 1100, Milwaukee, WI 53202-3823. Telephone: (414) 272-6071. <http://www.aaaai.org>.
- [4] Centers for Disease Control and Prevention, Vital Signs, May 2011.
- [5] World Health Organization. Global surveillance, prevention and control of chronic respiratory diseases: a comprehensive approach, 2007.
- [6] Adeyoye, Davies et al. An Estimate of Asthma Prevalence in Africa: A Systematic Analysis. *Croatian Medical Journal* 54.6 (2013): 519531. PMC. Web. 21 Apr. 2015.
- [7] Larrañaga, P., Calvo, B., Santana, R., Bielza, C., Galdiano, J., Inza, I., Lozano, J.A., Armañanzas, R., Santafé, G., Pérez, A. et al. (2006). Machine Learning in Bioinformatics. *Briefings in Bioinformatics*, 17(1), 86-112.
- [8] Stone M. Cross-validated choice and assessment of statistical predictions. *Journal of the Royal Statistical Society Series B* 1974;36:11147.
- [9] Ian H. Witten, Eibe Frank, Len Trigg, Mark Hall Geoffrey Holmes, and Sally Jo Cunningham. *Weka: Practical machine learning tools and techniques with java implementations*. Department of Computer Science. University of Waikato. New Zealand.
- [10] Green DM, Swets JA. *Signal Detection Theory and Psychophysics*. Wiley, 1974.