

MACHINE LEARNING IN NARRATIVE INTELLIGENCE

Santhanam, Assistant Professor, Department of Information Technology, Dhanalakshmi Srinivasan College of Engineering and Technology

K. Poornima Devi, Assistant Professor, Dhanalakshmi Srinivasan College of Engineering and Technology

ABSTRACT:

In this project, we were asked to experiment with a real world dataset, and to explore how machine learning algorithms can be used to find the patterns in data. We were expected to gain experience using a common data-mining and machine learning library, Weka, and were expected to submit a report about the dataset and the algorithms used. After performing the required tasks on a dataset of my choice, herein lies my final report. Machine Learning can be thought of as the study of a list of sub-problems, viz: decision making, clustering, classification, forecasting, deep-learning, inductive logic programming, support vector machines, reinforcement learning, similarity and metric learning, genetic algorithms, sparse dictionary learning, etc. Supervised learning, or classification is the machine learning task of inferring a function from a labeled data .

Keywords: Machine Learning, Pattern Recognition, Classification, Supervised learning, Artificial Intelligence

INTRODUCTION:

Machine learning is a sub-domain of computer science which evolved from the study of pattern recognition in data, and also from the computational learning theory in artificial intelligence. It is the first-class ticket to most interesting careers in data analytics today [1]. As data sources proliferate along with the computing power to process them, going straight to the data is one of the most straightforward ways to quickly gain insights and make predictions.

Machine Learning can be thought of as the study of a list of sub-problems, viz: decision making, clustering, classification, forecasting, deep-learning, inductive logic programming, support vector machines, reinforcement learning, similarity and metric learning, genetic algorithms, sparse dictionary learning, etc. Supervised learning, or classification is the machine learning task of inferring a function from a labeled data [2].

In Supervised learning, we have a training set, and a test set. The training and test set consists of a set of examples consisting of input and output vectors, and the goal of the

supervised learning algorithm is to infer a function that maps the input vector to the output vector with minimal error. In an optimal scenario, a model trained on a set of examples will classify an unseen example in a correct fashion, which requires the model to generalize from the training set in a reasonable way. In layman's terms, supervised learning can be termed as the process of concept learning, where a brain is exposed to a set of inputs and result vectors and the brain learns the concept that relates said inputs to outputs. A wide array of supervised machine learning algorithms are available to the machine learning enthusiast, for example Neural Networks, Decision Trees, Support Vector Machines, Random Forest, Naïve Bayes Classifier, Bayes Net, Majority Classifier [4,7,8,9] etc., and each has their own merits and demerits. There is no single algorithm that works for all cases, as merited by the No free lunch theorem [3]. In this project, we try and find patterns in a dataset [2], which is a sample of males in a heart-disease high risk region of South Africa, and attempt to throw various intelligently-picked algorithms at the data, and see what sticks.

PREVIOUS WORK:

The procedure described in this contribution is based on non-adaptive meta-models. This means, the sample set is pre-computed and the meta-model is trained on the basis of this data. Adaptive meta-models integrate the sampling and training in the usage phase [Dubourg (2011)]. Due to the adaptivity, they perform better for highly non-linear problems but are more complicated in implementation and usage.

DISADVANTAGES

- Massive data sets to train on, and these should be inclusive/unbiased, and of good quality
- ML needs enough time to let the algorithms learn and develop enough to fulfill their purpose with a considerable amount of accuracy and relevancy
- Challenge is the ability . To accurately interpret results generated by the algorithms.

PROPOSED WORK:

To experiment with a real world dataset, and to explore how machine learning algorithms can be used to find the patterns in data. We were expected to gain experience using a common data-mining and machine learning library, Weka, and were expected to submit a report about the dataset and the algorithms used.

ADVANTAGES

- The ability to learn, it lets them make predictions and also improve the algorithms on their own.
- Can review large volumes of data and discover specific trends and patterns that would not be apparent to humans.
- Gain experience, they keep improving in accuracy and efficiency. This lets them make better decisions.
- Machine Learning algorithms are good at handling data that are multi-dimensional and multi-variety, and they can do this in dynamic or uncertain environments.

4. MODEL OF EXPERIMENTS

4.1 PROBLEMS AND ISSUES IN SUPERVISED LEARNING: Before we get started, we must know about how to pick a good machine learning algorithm for the given dataset. To intelligently pick an algorithm to use for a supervised learning task, we must consider the following factors [4]:

Heterogeneity of Data:

Many algorithms like neural networks and support vector machines like their feature vectors to be homogeneous numeric and normalized. The algorithms that employ distance metrics are very sensitive to this, and hence if the data is heterogeneous, these methods should be the afterthought. Decision Trees can handle heterogeneous data very easily.

Redundancy of Data:

If the data contains redundant information, i.e. contain highly correlated values, then it's useless to use distance based methods because of numerical instability. In this case, some sort of Regularization can be employed to the data to prevent this situation.

Dependent Features:

If there is some dependence between the feature vectors, then algorithms that monitor complex interactions like Neural Networks and Decision Trees fare better than other algorithms.

Bias-Variance Tradeoff:

A learning algorithm is biased for a particular input x if, when trained on each of these data sets, it is systematically incorrect when predicting the correct output for x , whereas a learning algorithm has high variance for a particular input x if it predicts different output values when

trained on different training sets. The prediction error of a learned classifier can be related to the sum of bias and variance of the learning algorithm, and neither can be high as they will make the prediction error to be high. A key feature of machine learning algorithms is that they are able to tune the balance between bias and variance automatically, or by manual tuning using bias parameters, and using such algorithms will resolve this situation.

Curse of Dimensionality:

If the problem has an input space that has a large number of dimensions, and the problem only depends on a subspace of the input space with small dimensions, the machine learning algorithm can be confused by the huge number of dimensions and hence the variance of the algorithm can be high. In practice, if the data scientist can manually remove irrelevant features from the input data, this is likely to improve the accuracy of the learned function. In addition, there are many algorithms for feature selection that seek to identify the relevant features and discard the irrelevant ones, for instance

Principle Component Analysis for unsupervised learning. This reduces the dimensionality.

Overfitting:

The programmers should know that there is a possibility that the output values may constitute of an inherent noise which is the result of human or sensor errors. In this case, the algorithm must not attempt to infer the function that exactly matches all the data. Being too careful in fitting the data can cause overfitting, after which the model will answer perfectly for all training examples but will have a very high error for unseen samples.

4.2 DATASET:

The dataset used is a sample of males in a heart-disease high-risk region of the Western Cape, South Africa. The dataset that was used for this project is a subset of a much larger dataset, as described in Rousseau et al, 1983, South African Medical Journal, and has the following feature vectors:

- Sbp - systolic blood pressure
- Tobacco - cumulative tobacco (kg)
- Ldl - low density lipoprotein in cholesterol
- Adiposity - this is the amount of fat tissue in the body
- Famhist - family history of heart disease (Present, Absent) (String)
- Typea - type-A behavior
- Obesity - State of being overweight

- Alcohol - currentalcoholconsumption
- Age - Ageatonset
- Chd - coronaryheartdisease(ClassLabeloftheDataset)

In the dataset, there are 462 example vectors. Expert Systems have been used in the field of medical science to assist the doctors in making certain diagnoses, and this can help save lives. Coronary Heart Disease is a disease where a waxy substance builds up inside the coronary arteries, and hence this may lead to heart attack, and even death. When diagnosed and treated, the treatment can go a long way in helping the patient. This classification task is important because the expert system, when correctly generalized, can tell the doctor which patient may have the disease, and the doctor can take a look at that case in more detail. Moreover, if the doctor makes a slip, i.e. misdiagnoses someone, the expert system can help rectify his mistake. It results in two doctors, one of them virtual, instead of one doctor diagnosing every case which has a greater chance of accuracy and precision.

First we perform the significance analysis of the 9 feature vectors, to see which vector has more significance in representing the classes.

We used **Principal Component Analysis**[4,7,9]

for this purpose and came up with the following results.

Attribute Evaluator(supervised, Class(nominal): 10chd):

0.6795	10.516adiposity+0.46age+0.401ity+0.334ldl+0.324sbp...
0.5465	20.543alcohol+0.459tobacco-0.392obesity-0.364ldl-0.282typ
0.4269	3-0.792typea-0.459alcohol+0.338famhist+0.1e+0.125sbp...
0.322	4-0.833famhist-0.305obesity-0.258alcohol-0.21typea-0.196st

Correlationmatrix	0.2291	50.624tobacco-
1 0.21 0.160.36 -0.09 -0.060.24		0.419alcohol+0.321typea+0.305
0.140.39		hist-0.283obesity...
	0.1446	60.781sbp-
0.211 0.160.29-0.09 -0.01		0.379alcohol+0.332typea-0.215
0.120.2 0.45		0.174obesity...
	0.0706	70.788ldl-
0.160.161 0.44-0.160.040.33-		0.333obesity+0.277alcohol+0.2
0.030.31		p-0.196adiposity...
	0.0194	80.691age-0.489tobacco-
0.360.290.441 -0.18-0.040.720.1		0.339obesity-
0.63		0.235sbp+0.187famhist...
-0.09-0.09-0.16-0.181 -0.04-0.12 -0.08 -0.24		

Rankedattributes:

Selectedattributes: 1,2,3,4,5,6,7,8:8

Here we can see that all factors are important after we do the PCA. The last feature has been deemed unworthy by the PCA implementation in WEKA, which made little sense to us as age is highly correlated to most

diseases. We further our investigation by using another attribute selector, the

Significance Attribute Evaluator [5,9], on the data to yield

Significant feature evaluator Ranked attributes:

0.3019age

0.2992tobacco

0.2936typea

Selected attributes: 9,2,6,5,3,4,1,7,8:9

Here, we see that feature 9, i.e. the age of the patient was the most significant factor for classification purposes, and factors 7 and 8, obesity and alcohol consumption were the least significant factors. Through combined results of PCA and SAE, we conclude that all the features

are relevant for our purposes. The name of the sample was removed as well. Except for the use of PCA and SAE, no other pre-processing was done on the data.

Baseline Classifier:

As the baseline classifier, we chose a Naïve Bayesian Network because it is easy to compute, and because the features in the given dataset are all aspects of a person’s physical habits or medical history, and hence can be assumed to be independent of each other, which is the primary assumption in Naïve Bayes Classifier [6,8,9]. It is a conditional probability model: given a problem instance to be classified, represented by a vector representing some features (independent variables), it assigns to this instance probabilities for each of K possible outcomes or classes. The problem with the above formulation is that if the number of features is large or if a feature can take on a large number of values, then basing

such a model on probability tables is infeasible. We therefore reformulate the model to make it more tractable. Using Bayes' theorem, the conditional probability can be decomposed as

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

So, assuming that $p(x_i | x_{i+1}, \dots, x_n, C_k) = p(x_i | C_k)$ for $i = 1, \dots, n-1$. So, under the independence assumption, we can say that

Where, the evidence is a scaling factor dependent only on C_k , that is, a constant if the values of the feature variables are known.

Until now, we have derived an independent feature model. In Naïve Bayes classifier, we combine this model with a decision rule, and one of the common rules is to pick which hypothesis is the most probable. The corresponding classifier, a Bayes classifier, is the function that assigns a class label $\hat{y} = C_k$ for some k as follows:

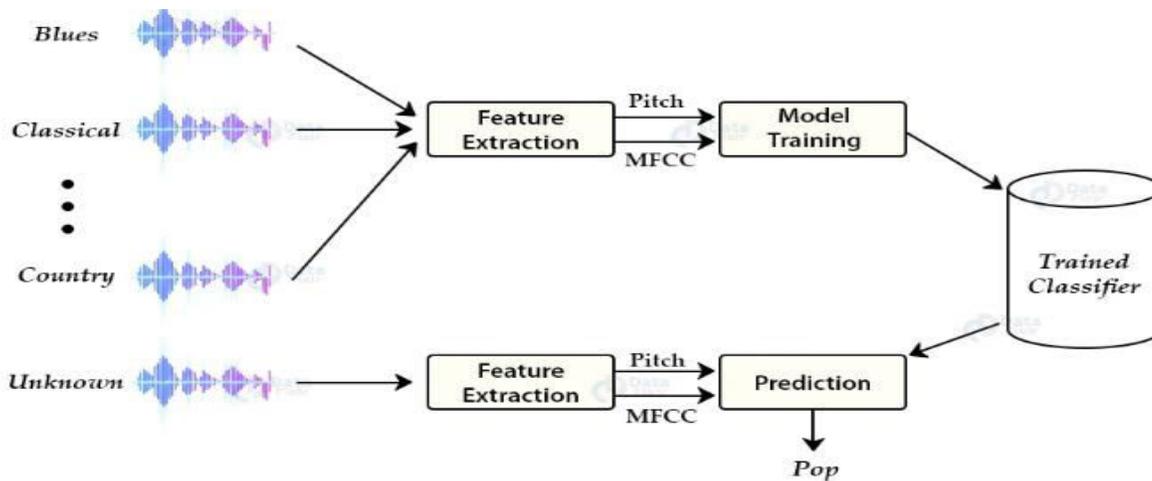
The Naïve Bayes network is already implemented in the Machine Learning library WEKA [5,9]. This was used on the aforementioned dataset, which led to the following output: As we can see that the Naïve Bayes classifier works really well with the given dataset, with the True Positive classification rate being 71.6 percent on an average, i.e. this classifier can correctly classify 71.6 percent of all the examples it sees. However, there is still a vast majority of the dataset, i.e. 28.4% which can't be correctly classified. This means that our expert medical diagnosis system still misdiagnoses one third of its cases, and one third

of the patients' symptom whomayhavethediseasearenotbeingscrutinizedbythedoctor.Wewill nowattempttoimproveresultbyusingothermoresophisticatedclassifierS.

4.3 CLASSIFIER:

Since it's a binary dataset with the class label being either the person has CHD or s/hedoesn't have CHD, and the number of samples is less than 100 times the number of features, thecorrelation matrix shows us that the correlation between various features is under . 5, we believethat support vector machines would be a viable classifier in this case.

SYSTEM ARCHITECTURE:



OUTPUT:

===DetailedAccuracy ByClass===

Class	TPRate	FPRate	Precision	Recall	F-Measure	MCC	ROCArea	PRCArea		
	0.825	0.506	0.755	0.825	0.788	0.335	0.659	0.737	0	
	0.494	0.175	0.598	0.494	0.541	0.335	0.659	0.471	1	
WeightedAvg.			0.710	0.392	0.700	0.710	0.702	0.335	0.659	0.645

===ConfusionMatrix ==

ab<-- classified as

24953|a=0

8179|b=1

Here, we can see that the said SVM performs better than the Naïve Bayes classifier for class 0, predicting 82.5% of the classes correctly, whereas it performs slightly worse than Naïve Bayes for class 1 with 49.4%. On an average, the true positive rate was achieved to be 71% as compared to 71.6% in case of Naïve Bayes. This result is surprising, as we expected SVM to perform better than the Naïve Bayes Classifier for independent non-redundant feature vectors as SVM projects low-dimensional sub-space to a higher dimensional subspace where the features are linearly separable. The RMS error for SVM was comparatively higher compared to Naïve Bayes by .10 and the kappa statistic of Naïve Bayes was lower than SVM by .05, which shows that Naïve Bayes is the better classifier. Curious about why the data was behaving the way it was, we did use other classifiers on the said dataset. We used Multilayer Perceptron, Decision Tree (J48) [8,9], Random Forest [8,9] with 100 trees, and the only classifier that got close was the J48 with true positive rate of 70.7%. performed poorly with only 63% TPR, -

learning neural net performed with 65.38% correct classifications. Curious if Lazy learning [8,9] could do any better, we tried it and found that it correctly classified 61.25% of the cases. The only thing we could now think of is that the input space was incomplete, and needed more dimensions for better predictions, and with the given feature vectors,

CONCLUSION:

We conclude that the dataset is not a complete space, and there are still other feature vectors missing from it. What we were attempting to generalize is a subspace of the actual input space, where the other dimensions are not known, and hence none of the classifiers were able to do better than 71.6% (Naïve Bayes). In the future, if similar studies are conducted to generate the dataset used in this report, more feature vectors need to be calculated so that the classifiers can form a better idea of the problem at hand.

FUTURE WORK:

The prospect of Machine Learning is not limited to the investment sector. As the Machine Learning prospect is very high, there are some of the areas where researchers are working toward revolutionizing the world for the future. Thus, the future prospect of Machine Learning will accelerate the processing power of the automation system used in various types of different collection technologies.

Further, in this blog on the future prospect of Machine Learning, we will look into the skills that are required to become an ML Engineer. In this blog on the future prospect of Machine Learning, we have looked around the need for Machine Learning. We have seen the future prospect of Machine Learning and the opportunities in the field.

REFERENCES:

- ElementsofStatisticalLearning:DataMining,Inference,andPrediction. 2ndEdition.
- Datasets:CoronaryHeartDiseaseDataset."ElementsofStatisticalLearning:DataMining, Inference, and Prediction. 2nd Edition. Accessed April 27, 2016.<http://statweb.stanford.edu/~tibs/ElemStatLearn/>.
- "NoFreeLunchTheorems."No FreeLunchTheorems.Accessed April27, 2016.<http://www.no-free-lunch.org/>.
- Hastie, Trevor, Robert Tibshirani, and J. H. Friedman. The Elements of StatisticalLearning:DataMining,Inference,andPrediction:With200Full-colorIllustrations.NewYork:Springer,2001.
- "Weka 3: Data MiningSoftware in Java."Weka 3.AccessedApril27,2016.<http://www.cs.waikato.ac.nz/ml/weka/>.
- Bozhinova,Monika,NikolaGuid,andDamjanStrnad.Naivni Bayesov Klasifikator:DiplomskoDelo.Maribor:M.Bozhinova,2015.
- Schölkopf, Bernhard, ChristopherJ. C. Burges, andAlexanderJ. Smola.Advances inKernelMethods:SupportVectorLearning.Cambridge,MA:MITPress,1999.

Norving, Peter, and
Stuart Russel. *Artificial Intelligence: A Modern Approach*. S.l.: Pearson Education Limited, 2013.

- Witten, I.H., and Eibe Frank. *Data Mining: Practical Machine Learning Tools and Techniques*. Amsterdam: Morgan Kaufman, 2005.
"Intro to Machine Learning | Udacity." *Intro to Machine Learning | Udacity*. Accessed April 27, 2016. <https://www.udacity.com/course/intro-to-machine-learning--ud120>.