

STATISTICAL MODEL IN MACHINE LEARNING TO PREDICT SEVERE RETINOPATHY OF PREMATURITY

G. Pavani

B. Jyothi

Dr. M. Sandeep

Assistant professor Assistant professor Associate professor
Dept of EEE

Sree Dattha Institute of Engineering and Science

ABSTRACT

Using statistical analysis and logistic regression as a type of Generalized Additive Model (GAM) with pair-wise interaction terms (GA2M), we investigated the risk factors for severe retinopathy of prematurity. We talk about the trade-off between these machine learning techniques' accuracy and interpretability on clinical data in this process. We also confirm expert neonatologists' instincts regarding a few risk factors, such as gender, that were previously thought to be clinically insignificant in RoP prediction. In many middle-income countries, retinopathy of Prematurity (ROP), a vasoproliferative disorder of the immature retina in premature infants, is a significant cause of blindness. In high-income nations, risk factors like oxygen administration and blood oxygen saturation are closely monitored, resulting in a lower incidence of ROP. In developed nations, severe ROP is typically found in infants born at a very early Gestational Age (GA). Because there is insufficient awareness of the disease process's risk factors, a lack of skilled professionals, and/or a lack of essential equipment for infant care, heavier, more mature babies can also develop ROP in developing nations.

KEYWORDS: interpretability of machine learning in healthcare, generalized additive model, logistic regression Retinopathy of Prematurity (RoP), neonatology

I INTRODUCTION

A specific type of artificial intelligence known as machine learning enables systems to learn from data and recognize patterns without much human intervention. Computers using machine learning are shown patterns and data, allowing them to make their own decisions rather than being instructed. Algorithms for machine learning serve a variety of purposes, including assisting in the filtering of email, identifying objects in images, and analyzing large quantities of increasingly complex data sets. Machine learning systems are used by computers to automatically sort through emails, identify spam, recognize objects in pictures, and process large amounts of data. Techniques for machine learning are a growing area of study with numerous potential

applications. For healthcare professionals and health systems, machine learning technology will become increasingly important as patient data becomes more readily available for meaning extraction. Machine learning is especially useful in the healthcare industry because it can help us make sense of the enormous amounts of healthcare data that are generated each day in electronic health records. We can use machine learning in healthcare in a manner that is analogous to how algorithms from machine learning can assist us in discovering patterns and insights that cannot be discovered manually. With the widespread acceptance of machine learning in healthcare, healthcare providers now have the opportunity to implement a more predictive strategy that results in the development of a more unified system with enhanced care delivery and patient-centered processes.

The most common uses of machine learning in the healthcare sector are to automate medical billing, provide clinical decision support, and develop clinical practice guidelines within health systems. Machine learning and healthcare concepts are used in numerous notable high-level applications in science and medicine. Data scientists at MD Anderson developed the first deep learning healthcare algorithm for predicting acute toxicities in patients receiving radiation therapy for head and neck cancers. Deep learning data can automatically identify intricate patterns in healthcare data and provide primary care providers with clinical decision support at the point of care within the electronic health record.

For the purposes of machine learning, large volumes of unstructured healthcare data account for nearly 80% of the data stored or "locked" in electronic health record systems. These are absolutely not data elements; rather, they are pertinent text files or data documents that contain patient data that, in the past, could not be analyzed without a human reviewing the medical records. Human language, also known as "natural language," is extremely complicated, inconsistent, and full of jargon, ambiguity, and ambiguity. In the field of healthcare, machine learning frequently makes use of natural language processing (NLP) software to transform these documents into data that is easier to use and analyze. Machine learning relies on healthcare data for the majority of NLP-based deep learning applications.

Recent trends in machine learning applications in healthcare center on the interpretability of the models. In some healthcare applications, interpretability of the model is more important than accuracy, but there are many situations in which interpretability is preferred despite the loss of accuracy. One way to achieve the ideal situation in which the model has both high accuracy and high interpretability is to start with a simple model like a generalized additive model (GAM) and then make it more complex (and thus more accurate) like a GA2M (GAM with pairwise interactions), or to start with a complex model like XGBoost and try to interpret it locally with methods like LIME. Ann Lecun,

Facebook's head of AI research, and Rich Caruana, Microsoft's lead machine learning researcher on healthcare applications, organized a debate session at NeurIPS (Neural Information Processing Systems), one of the most significant conferences on machine learning. During this session at the "Interpretable ML Symposium's" conclusion, the need for additional case studies on the accuracy versus interpretability trade-off in machine learning for healthcare was highlighted. [9][10] In this paper, we will examine a case study regarding Severe Retinopathy of Prematurity (Severe RoP), which, if left untreated, can result in newborns going completely or partially blind. Stevie Wonder was a well-known musician who had a lot of RoP. Doctors only recently realized that oxygen could save premature infants in the 1950s. However, it took them a few more years to realize that excessive oxygen caused the vessels in the eyes to grow out of place.

Premature babies would go blind for life if this condition was treated with too much oxygen. The World Health Organization (WHO) has estimated that 15 million of the 130 million babies born each year are born prematurely, or before the 37th week of gestation.

Each year, approximately one million children die because of preterm birth complications. Among survivors, learning disabilities as well as visual and auditory impairments are common disabilities. Among these issues, retinopathy of prematurity (RoP) is the leading and most serious cause of disability. Retinopathy of Prematurity (RoP) was first described by Terry [11] as a developmental, vascular, and proliferative retinal disorder that affects the retinas of premature newborns that have not yet undergone vascularization.

Along with cortical blindness, RoP is one of the most common causes of childhood blindness worldwide. RoP has been linked to a number of risk factors, and if severe RoP is not treated, it can cause retinal detachment blindness. With a lower birth weight and shorter gestational weeks, RoP becomes more severe. To stop the disease from progressing, early diagnosis and appropriate treatment are essential. The International Classification for Retinopathy of Prematurity (ICROP) is used to categorize the condition based on its progression and severity. The Newborn Clinic at Istanbul's Zeynep Kamil Woman and Child Diseases Hospital collected the clinical data for our analysis between 2011 and 2014. The diagnosis of retinal detachment, which is the most severe repercussion of RoP, serves as the starting point for classification. Every year, around 150,000 babies in Turkey have birth weights below 1,500 grams. As we improve our prediction model, we will look at how the accuracy of our interpretable model improves as we combine numerical and categorical values and add additional interaction terms. We will also investigate the risk factors that are thought to cause or be correlated with severe RoP. It is more likely that these newborns will be diagnosed with severe RoP.

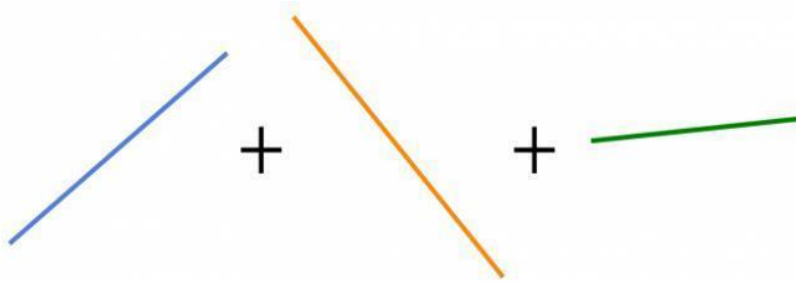
II. PROPOSED METHODOLOGY

Our research begins with data on retinopathy of premature (RoP) from Zeynep Kamil Maternity and Child Diseases Hospital from 2011 to 2014. It is a record of 102 variables across more than 5000 patients. ZK-RoP) For patients with birth weights below 1500 grams, we were interested in the same 20 risk factors examined in the TR-RoP study [4]. We utilized their data for 1066 patients who weighed less than 2000 grams because the ZK-RoP study focused on newborns under 2000 grams. Among the 1066 newborns with birthweights below 2000 grams, 109 cases of severe RoP were examined. Our sample size was reduced to 385 because the important task was to predict who might develop severe RoP from those who had already been diagnosed with any type of RoP. We ran univariate and multivariate logistic regression machine learning algorithms, beginning with the generalized additive model (GAM) to predict severe RoP based on the same risk factors as in the TR-RoP study [4]. We also added interaction terms to our multivariate analysis in order to minimize type II error and minimize wrong diagnosis of a disease that could result in blindness. Out of 385 patients who were diagnosed with any type of RoP, we attempted to construct a model that would predict the 109 patients who we wanted to use the same method as the TR-RoP paper and check their results with a new set of data, ZK-RoP. We also wanted to include accuracy analysis of these predictions, which the TR-RoP study didn't do because it wasn't geared toward ML users. We could use a number of ML algorithms, but we wanted to start with the simplest and move up to the more complex ones. We only focused on the multivariate logistic regression analysis in this paper. We will analyze with more complicated models like decision trees, k-means, and xgboost in further research.

2.1 Generalized Additive Model with Pairwise Interactions (GA2M):

GAM is short for generalized additive model. Before we discuss GAMs, let's first briefly review a common statistical model that you are likely to be familiar with. It is the generalized linear model (GLM). A GLM is a very popular and flexible extension of the classical linear regression model. It enables you to model a target (or response) variable that is not normally distributed.

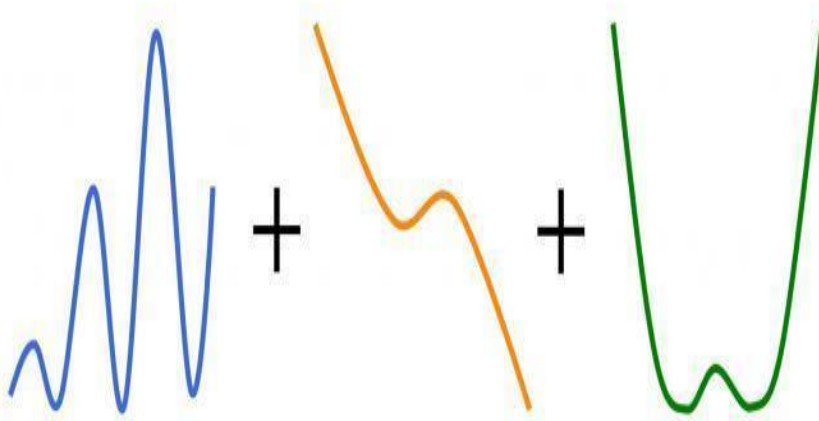
For example, if you are interested in predicting whether an incoming email is spam, or you want to predict the number of people dining in a restaurant, the linear regression model is inappropriate. This is because the target variables for those applications violate the model's normality assumption. However, a GLM can handle both data types by assuming a distribution in the exponential



$$g(\mu_i) = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + c$$

Fig 1: Linear models assume each componentfunction is linear. Source: Adapted from Servén Marín (2018)

For GAMs, on the other hand, the component functions can be arbitrary nonlinear smoothing functions that the model usually estimates by using spline functions. You can see this in Figure 2.



$$g(\mu_i) = f_1(x_1) + f_2(x_2) + f_3(x_3) + c$$

Fig 2: GAMs use smoothing componentfunctions to flexibly model nonlinearities.

Source: Adapted from Servén Marín (2018) Even with this flexibility, GAMs are still interpretable. This is an increasingly desirable feature of machine learning models. This is especially true for high-stakes applications such as health care, finance, and mission-critical systems. Typically, you face a trade-off between model performance and model interpretability.

For the effects in the final model, the node’s results include smoothing component plots for the spline terms and parameter estimates for the parametric terms. For example, the results include a smoothing component plot for the Spline(Debtinc) term. Figure 3 displays this. Interpreting the results, you can see that the probability of default is generally higher for an applicant who has a high debt-to-income ratio and the relationship is nonlinear.

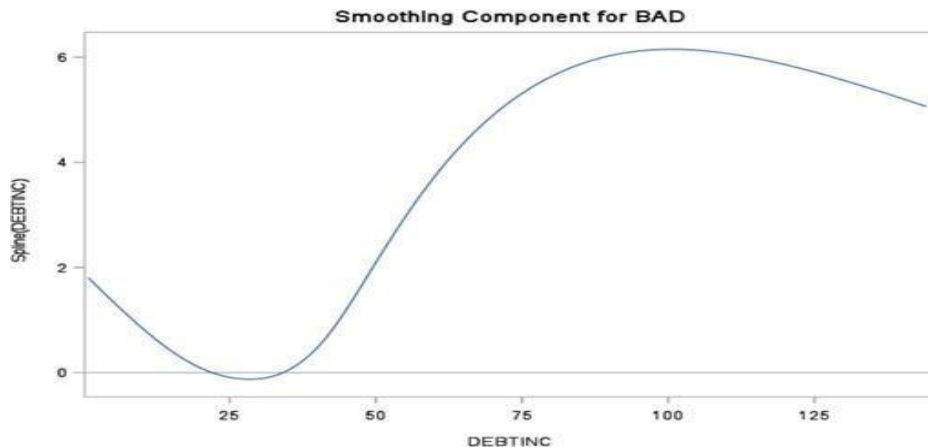


Fig 3: The probability of default is generally higher for higher debt-to-income ratios.

In addition, the ability to interpret the parametric effects like you would with a logistic regression model is another way in which the GAM’s results are relatively easy to understand. For example, an applicant with more delinquent credit lines has a higher predicted probability of default, on average. See Figure 4.

Effect	Parameter ↑	Estimate
DELINQ	DELINQ 0	-4.4301
DELINQ	DELINQ 1	-4.0737
DELINQ	DELINQ 2	-3.5835
DELINQ	DELINQ 3	-3.5960
DELINQ	DELINQ 4	-2.6697
DELINQ	DELINQ 5	-1.0057
DELINQ	DELINQ 6	-0.8645
DELINQ	DELINQ 7	-0.4681
DELINQ	DELINQ 8	-0.2449

Fig 4: More delinquent credit lines typically correspond to a higher probability of default.

. With Model Studio, you can also easily train other supervised learning models for comparison. Let's compare the GAM to gradient boosting, logistic regression, neural network, and support vector machine (SVM) models. If you rank the algorithms by misclassification rate, you see that the GAM ranks second. Figure 5 displays the results.

Algorithm Name	Misclassification Rate ↑
Gradient Boosting	0.0481
GAM	0.0571
Logistic Regression	0.0586
SVM	0.0662
Neural Network	0.0782

Fig 5: The GAM outperforms some of the more complex models but maintains interpretability.

Even though the GAM has a slightly higher misclassification rate, it is more interpretable than the

Fig 4: More delinquent credit lines typically correspond to a higher probability of default.

With Model Studio, you can also easily train other supervised learning models for comparison. Let's compare the GAM to gradient boosting, logistic regression, neural network, and support vector machine (SVM) models. If you rank the algorithms by misclassification rate, you see that the GAM ranks second. Figure 5 displays the results.

Algorithm Name	Misclassification Rate ↑
Gradient Boosting	0.0481
GAM	0.0571
Logistic Regression	0.0586
SVM	0.0662
Neural Network	0.0782

Fig 5: The GAM outperforms some of the more complex models but maintains interpretability.

Even though the GAM has a slightly higher misclassification rate, it is more interpretable than the champion gradient boosting mode. And it still outperforms other complex models such as SVM and neural network.

In this model we are trying to predict a categorical variable, whether the patient is diagnosed with severe RoP or not, we will use logistic regression instead of linear regression. Univariate logistic regression holds only one variable $\ln [Y / (1-Y)] = a + b_1 X_1$, while the formula for multivariate logistic regression is as follows: $\ln [Y / (1-Y)] = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + \dots$. The generalized additive model (GAM) with only 1st order terms is equivalent to the equation above. If we want to add 2nd order terms, we would need to add them as follows: $\ln [Y / (1-Y)] = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_{12} X_1 X_2 + b_{13} X_1 X_3 + b_{23} X_2 X_3 + \dots$. or example, for a GAM with 6 first-order variables ($X_1, X_2, X_3, X_4, X_5, X_6$), there would be $6*5/2=15$ pairwise interaction terms.

It is suggested to add only a few and only start with the most influential pairwise interactions. Thus, as we make our model more complex, we will add the pairwise interactions one by one. However, when we train a machine learning model, we typically face a tradeoff between accuracy and interpretability. GAMs provide a solution to this dilemma because they strike a nice balance between these two competing goals. With Model Studio, you can easily train GAMs alongside other common machine learning models to compare model performance, all without the need to write any code.

II CONCLUSION

In our research 12 risk factors are used in the multivariate linear regression: 10 binary (categorical) values and two numbers. First, we ran the regression with numerical and categorical values separately. The multivariate logistic

regression analysis that would produce a binary prediction of severe RoP was then carried out by combining the numerical and categorical risk factors into a single equation. Last but not least, it is important to note that comparisons between GAM and GA2M and other machine learning methods that can be explained, such as decision trees, require additional research. Perhaps our RoP prediction problem and data were a better fit for other interpretable machine learning techniques given our decision tree approach with cross-validation in SPSS, which produced a higher accuracy rate of 44% and only a few type II errors. Cross-validation on RoP data and various interpretable ML techniques need to be compared in additional research. We hope to use Microsoft's interpretable ML library and the scikit-learn ML library as we expand our explainable ML work based on this paper. Both of these libraries promise high accuracy and high interpretability.

REFERENCES

- [1] R. Caruana, Y. LeCun, The Great AI Debate “Interpretable ML Symposium” as part of NeurIPS- 2017.
<https://www.youtube.com/watch?v=93Xv8vJ2acI>
- [2] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad. Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 1721–1730. ACM, 2015.
- [3] O. Gottesman, F. Johansson, J. Meier, J. Dent, D. Lee, S. Srinivasan, L. Zhang, Y. Ding, D. Wihl, X. Peng, J. Yao, I. Lage, C. Mosch, L. Lehman, M. Komorowski, A. Faisal, L. A. Celi, D. Sontag, and F. Doshi-Velez. Evaluating Reinforcement Learning Algorithms in Observational Health Settings. pp.1-16, 2018.
<https://arxiv.org/pdf/1805.12298.pdf>
- [4] A.Y. Bas, N. Demirel, E. Koc, D. Ulubas Isik, I.M. Hirfanoglu, T. Tunc, and TR-ROP Study Group. Incidence, risk factors and severity of retinopathy of prematurity in Turkey (TR-ROP study): a prospective, multicentre study in 69 neonatal intensive care units. *Br J Ophthalmol.* 102(12):1711-1716, 2018.
- [5] Y. Lou, R. Caruana, J. Gehrke, and G. Hooker. Accurate Intelligible Models with Pairwise Interactions. KDD2013, August 11–14, 2013, Chicago, Illinois, USA.
<http://www.cs.cornell.edu/~yinlou/papers/lou->