

Predicting Preterm Birth in Maternity Care by means of Data Mining

Sónia Pereira, Filipe Portela, Manuel F. Santos, José Machado, António Abelha
Algoritmi Centre, University of Minho, Portugal
b7004@dps.uminho.pt; {cfp, mfs}@dsi.uminho.pt
{jmac, abelha}@di.uminho.pt

Abstract. Worldwide, around 9% of the children are born with less than 37 weeks of labour, causing risk to the premature child, whom it is not prepared to develop a number of basic functions that begin soon after the birth. In order to ensure that those risk pregnancies are being properly monitored by the obstetricians in time to avoid those problems, Data Mining (DM) models were induced in this study to predict preterm births in a real environment using data from 3376 patients (women) admitted in the maternal and perinatal care unit of Centro Hospitalar of Oporto. A sensitive metric to predict preterm deliveries was developed, assisting physicians in the decision-making process regarding the patients' observation. It was possible to obtain promising results, achieving sensitivity and specificity values of 96% and 98%, respectively.

Keywords: Data Mining, Preterm Birth, Real data, Obstetrics Care, Maternity Care

1 Introduction

Preterm birth portrays a major challenge for maternal and perinatal care and it is a leading cause of neonatal morbidity. The medical, education, psychological and social costs associated with preterm birth indicate the urgent need of developing preventive strategies and diagnostic measures to improve the access to effective obstetric and neonatal care [1]. This may be achieved by exploring the information provided from the information systems and technologies increasingly used in healthcare services.

In Centro Hospitalar of Oporto (CHP), a Support Nursing Practice System focused on nursing practices (SAPE) is implemented, producing clinical information. In addition, patient data plus their admission form are recorded through EHR (Electronic Health Record) presented in Archive and Diffusion of Medical Information (AIDA) platform. Both SAPE and EHR are also used by the CHP maternal and perinatal care unit, Centro Materno Infantil do Norte (CMIN). CMIN is prepared to provide medical care / services for women and child. Therefore, using obstetrics and prenatal information recorded from SAPE and EHR, it is possible to extract new knowledge in the context of preterm birth. This knowledge is achieved by means of Data Mining (DM) techniques, enabling predictive models based on evidence. This study

accomplished DM models with sensitivity and specificity values of approximately 96% and 98%, which are going to support the making of preventive strategies and diagnostic measures to handle preterm birth.

Besides the introduction, this article includes a presentation of the concepts and related work in Section 2, followed by the data mining process, described in Section 3. Furthermore, the results are discussed and a set of considerations are made in Section 4. Section 5 presents the conclusions and directions of future work.

2 Background and Related Work

2.1 Preterm Birth

Preterm birth refers to a delivery prior to 37 completed weeks (259 days) of labour. Symptoms of preterm labour include uterine contractions occurring more often than every ten minutes, or the leaking of fluids. Preterm birth is the leading cause of long-term disability in children, since many organs, including the brain; lungs and liver are still developing in the final weeks of pregnancy [2]. Preterm Birth has not decreased in the last 30 years, due to the failure identifying the high-risk group during routine prenatal care [3]. Many studies were conducted to identify a way to predict preterm deliveries, focusing on physiologic measures, ultrasonography, obstetrics history and socioeconomic status [4]. For instance, in 2011 a model was developed for predicting spontaneous delivery before 34 weeks based on maternal factors, placental perfusion and function at 11-13 weeks' gestation, through screening maternal characteristics and regression analysis. They detected 38.2% of the preterm deliveries in women with previous pregnancies beyond 16 weeks and 18.4% in those without [3]. Most of the efforts to predict preterm birth face limited provision of population based data, since registration of births is incomplete and information is lacking on gestational age [6].

2.2 Interoperability Systems and Data Mining in Healthcare

As mentioned in the previous section, this study is based on real data acquired from CMIN. The knowledge extraction depends substantially on the interoperability between SAPE and EHR systems assured through AIDA. This multi-agent platform enables the standardization of clinical systems and overcomes the medical and administrative complexity of the different sources of information from the hospital [5].

In healthcare systems, there is a wealth of data available, although there is a lack of effective analysis tools to extract useful information. Thus, data mining have found numerous applications in scientific and clinical domain [8]. Successful mining applications have been implemented in the healthcare. In obstetrics and maternal care, some of these studies were employed to predict the risk pregnancy in women performing voluntary interruption of pregnancy (VIP) [9] and manage VIP by predicting the most suitable drug administration [7].

3 Study Description

This study was conducted by following the Knowledge Discovery in Database (KDD), allowing the extraction of implicit and potentially useful information, through algorithms, taking account the magnitudes of data increasing [10].

The DM methodology employed was the Cross Industry Standard Process for Data Mining (CRIP-DM), a non-rigid sequence of six phases, carried out in this section, which allow the implementation of DM models to be used in real environments [11]. To induce the DM models, four different algorithms were implemented: Decision Trees (DT), Generalized Linear Models (GLM), Support Vector Machine (SVM) and Naïve Bayes (NB). This study used data collected from 3376 patients (women) admitted in the maternal and perinatal care unit (CMIN) of CHP comprising a period between 2012-07-01 and 2015-01-31, in a total of 1120 days.

3.1 Business Understanding

The Business aim of this project is to identify the risk group of preterm delivery, to ensure the proper monitoring and to avoid its associated problems. The DM goal is to develop accurate models able to support the decision-making process by predicting whether or not a woman will be subjected to a preterm delivery, based on data from clinic cases.

3.2 Data Understanding

The initial dataset extracted from SAPE and EHR admission records was analysed and processed in order to be used in the DM process. A set of 13 variables were selected: age (corresponds to the age of the pregnant patient), programmed (indicates whether or not a delivery is programmed), gestation (singular or multiple pregnancies), PG1 and PG2 (first echography measures), motive (reason of intervention - normal delivery or unexpected events), patients' weight and height, BMI (body mass index), blood type, cardiotocography (CTG) (biophysics exam that evaluates the fetal wellbeing), streptococcus (presence of the bacterium streptococcus in the pregnant system) and finally, marital status of the pregnant patient. The target variable *Group Risk* denotes the preterm birth risk and it is presented in Table I.

Table I: Representation of the target variable *Group Risk*.

Description	Value	Target Distribution	Percentage
≥ 37 weeks of gestation (Term)	0	3137	92.92%
< 37 weeks of gestation (Preterm)	1	239	7.08%

In Table II are shown statistics measures related to the numerical variables age, gestation, PG1, PG2 and BMI, while in Table III it is represented the percentage of occurrences for some used variables.

Table II: Statistics measures of age, PG1, PG2, weight, height, BMI variables.

	Minimum	Maximum	Average	Standard Deviation
<i>Age</i>	14	46	29.88	5.81
<i>PG1</i>	5	40	12.81	2.96
<i>PG2</i>	0	8	3.09	1.96
<i>BMI</i>	14.33	54.36	29.40	4.57

Table III: Percentage of occurrences of some variables.

<i>Variable</i>	<i>Class</i>	<i>Cases</i>
<i>Programmed</i>	True	12.53%
<i>Gestation</i>	Singular	89.90%
<i>Motive</i>	Normal	81.33%
<i>Streptococcus</i>	Positive	13.27%
<i>Cardiotocography</i>	Suspect	2.19%

3.3 Data Preparation

After understanding the data collected, the variables were prepared to be used by the DM models. The data pre-processing phase started with the identification of null and noise values. These values were eliminated from the dataset. To ensure the data normalization, all the values, such as weight and height, were transformed to International System measures, using the point to separate decimal values.

As shown in Table 1, there is a disparity in the distribution of values of the target variable Risk Group (low percentage of preterm birth cases). In order to balance the target, the oversampling technique was implemented by replicating the preterm birth cases until it reached approximately 50% of the dataset, obtaining 6244 entries.

3.4 Modelling

A set of Data Mining models (DMM) were induced using the four DM techniques (DMT) mentioned in Section 3: GLM, SVM, DT and NB. The developed models used two sampling methods Holdout sampling (30% of data for testing) and Cross Validation (all data for testing). Additionally there were implemented two different approaches, one using the raw dataset (3376 entries) and another with oversampling. Different combinations of variables were used, obtaining 5 different scenarios:

- S1: {Age (A), Gestation (G), Programmed (P), PG1, PG2, Motive (M), Height (H), Weight (W), BMI, Blood Type (B), Marital Status (MS), CTG, Streptococcus (S)}
- S2: {A, H, W, BMI, B, MS, CTG, S}
- S3: {G, P, PG1, PG2, M, CTG, S}
- S4: {A, G, PG1, PG2, M, H, W, BMI, B, CTG, S}
- S5: {A, G, P, M, H, W, BMI, B}

Therefore, a total of 80 Data Mining models (DMM) were induced:

$$DMM = \{5 \text{ Scenarios}, 4 \text{ Techniques}, 2 \text{ Sampling Methods}, 2 \text{ Approaches}\}$$

All the models were induced using the Oracle Data Miner with its default configurations. For instance, GLM was induced with automatic preparation, with a confidence level of 0.95 and a reference value of 1.

3.5 Evaluation

The study used the confusion matrix (CMX) to assess the induced DM models. Using the CMX, the study estimated some statistical metrics: sensitivity, specificity and accuracy. Table IV presents the best results achieved by each technique, sampling method and approach. The best accuracy (93.00%) was accomplished with scenario 3 by both DT and NB techniques using oversampling and 30% of data for testing. The best sensitivity (95.71%) was achieved by scenario 4 with oversampling using SVM technique and all the data for testing. Regarding specificity, scenario 2 reached 97.52% using SVM with oversampling and all the data for testing.

Table IV: Sensitivity, specificity and accuracy values for the best scenarios for each DMT, approach and sampling method. Below, the best metric values highlighted for each DMT.

<i>DMT</i>	<i>Oversampling</i>	<i>Sampling</i>	<i>Scenario</i>	<i>Sensitivity</i>	<i>Specificity</i>	<i>Accuracy</i>
<i>DT</i>	No	30%	3	0.8889	0.9303	0.9300
	No	All	1	0.2896	0.9723	0.8599
<i>GML</i>	No	All	4	0.2896	0.9723	0.8599
	Yes	All	4	0.8674	0.7126	0.7687
<i>NB</i>	No	30%	3	0.8889	0.9303	0.9300
	No	All	1	0.4868	0.9646	0.9271
<i>SVM</i>	No	All	2	0.1023	0.9752	0.4570
	Yes	All	4	0.9571	0.6647	0.7410

In order to choose the best models a threshold was established, considering sensitivity, accuracy and specificity values upper than to 85%. Table V shows the models that fulfil the threshold.

Table V: Best model achieving the established threshold.

<i>Scenario</i>	<i>Model</i>	<i>Oversampling</i>	<i>Sampling</i>	<i>Sensitivity</i>	<i>Specificity</i>	<i>Accuracy</i>
3	NB,DT	No	30%	0.8889	0.9303	0.9300

4 Discussion

Should be noted that the best sensitivity (95.71%) and specificity (97.52%) are reached by models that did not achieve the threshold defined, showing low values in the remaining statistical measures used to evaluate the models. It can be settled that scenario 3 meets the defined threshold, presenting good results in terms of specificity and sensitivity, as seen in Table V. Thus, it appears that the most relevant factors that affect the term of birth are: pregnancy variables, Gestation and physical conditions of the pregnant woman. In a clinic perspective, the achieved results will enable the prediction of preterm birth, with low uncertainty, allowing those responsible better monitoring and resource management. In a real time environment, physicians can rely on the model to send a warning informing that a specific patient has a risk pregnancy and it is in danger of preterm delivery. Consequently, the physician can be observant and alert to these cases and can put the patients on special watch, saving resources and time to the healthcare institution.

5 Conclusions and Future Work

At the end of this work it is possible to assess the viability of using these variables and classification DM models to predict Preterm Birth. The study was conducted using real data. Promising results were achieved by inducing DT and NB, with oversampling and 30% of the data for testing, in scenario 3, achieving approximately 89% of sensitivity and 93% of specificity, suited to predict preterm births. The developed model support the decision-making process in maternity care by identifying the pregnant patients in danger of preterm delivery, alerting to their monitoring and close observation, preventing possible complications, and ultimately, avoiding preterm birth.

In the future new variables will be incorporated in the predictive models and other types of data mining techniques will be applied. For instance, inducing Clustering techniques would create clusters with the most influential variables to preterm birth.

Acknowledgments

This work has been supported by FCT - Fundação para a Ciência e Tecnologia within the Project Scope UID/CEC/00319/2013.

References

1. Berghella, V. (Ed.): Preterm birth: prevention and management. John Wiley & Sons (2010)
2. Spong, C. Y.: Defining "term" pregnancy: recommendations from the Defining "Term" Pregnancy Workgroup. *Jama*, 309(23), 2445-2446 (2013)
3. Beta, J., Akolekar, R., Ventura, W., Syngelaki, A., Nicolaides, K. H.: Prediction of spontaneous preterm delivery from maternal factors, obstetric history and placental perfusion and function at 11–13 weeks. *Prenatal diagnosis*, 31(1), 75-83 (2011)
4. Andersen, H. F., Nugent, C. E., Wanty, S. D., Hayashi, R. H.: Prediction of risk for preterm delivery by ultrasonographic measurement of cervical length. *AJOG*, 163(3), 859-867 (1990)
5. Abelha, A., Analide, C., Machado, J., Neves, J., Santos, M., Novais, P.: Ambient intelligence and simulation in health care virtual scenarios. In *Establishing the Foundation of Collaborative Networks* (pp. 461-468). Springer US (2007)
6. McGuire, W., Fowlie, P. W. (Eds.): *ABC of preterm birth* (Vol. 95). John Wiley & Sons (2009)
7. Brandão, A., Pereira, E., Portela, F., Santos, M. F., Abelha, A., Machado, J.: Managing voluntary interruption of pregnancy using data mining. *Procedia Technology*, 16, 1297-1306 (2014)
8. Kaur, H., Wasan, S. K.: Empirical study on applications of data mining techniques in healthcare. *Journal of Computer Science*, 2(2), 194 (2006)
9. Brandão, A., Pereira, E., Portela, F., Santos, M. F., Abelha, A., Machado, J.: Predicting the risk associated to pregnancy using data mining. *ICAART 2015 Portugal*. SciTePress (2015)
10. Maimon, O., Rokach, L.: Introduction to knowledge discovery in databases. In *Data Mining and Knowledge Discovery Handbook* (pp. 1-17). Springer US (2005)
11. Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C., Wirth, R.: *CRISP-DM 1.0 Step-by-step data mining guide* (2000)