

# Predicting Infant Motor Development Status using Day Long Movement Data from Wearable Sensors

David Goodfellow, Ruoyu Zhi  
University of Southern California  
{rzhi,dgoodfel}@usc.edu

Rebecca Funke  
University of Southern California  
rfunke@usc.edu

José Carlos Pulido  
Universidad Carlos III de Madrid  
jcpulido@inf.uc3m.es

Maja Matarić  
University of Southern California  
mataric@usc.edu

Beth A. Smith  
University of Southern California  
beth.smith@pt.usc.edu

## ABSTRACT

Infants with a variety of complications at or before birth are classified as being at risk for developmental delays (AR). As they grow older, they are followed by healthcare providers in an effort to discern whether they are on a typical or impaired developmental trajectory. Often, it is difficult to make an accurate determination early in infancy as infants with typical development (TD) display high variability in their developmental trajectories both in content and timing. Studies have shown that spontaneous movements have the potential to differentiate typical and atypical trajectories early in life using sensors and kinematic analysis systems. In this study, machine learning classification algorithms are used to take inertial movement from wearable sensors placed on an infant for a day and predict if the infant is AR or TD, thus further establishing the connection between early spontaneous movement and developmental trajectory.

## CCS CONCEPTS

• **Binary Classification**; • **Interpretability**; • **Patient Outcome**;

## KEYWORDS

Infant Motor Development, Classification, Sensor Data

### ACM Reference Format:

David Goodfellow, Ruoyu Zhi, Rebecca Funke, José Carlos Pulido, Maja Matarić, and Beth A. Smith. 2018. Predicting Infant Motor Development Status using Day Long Movement Data from Wearable Sensors. In *Proceedings of KDD Workshop on Machine Learning for Healthcare and Medicine (KDD'18)*. ACM, New York, NY, USA, 5 pages. [https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

## 1 INTRODUCTION

An increasing number of infants with neuromuscular impairments survive past birth due to improvements in obstetric and neonatal medicine [1]. These infants with complications at birth are classified as being at risk for developmental delay (AR), and some are later diagnosed with cognitive and/or physical developmental delays [8]. Early motor delays are often the initial signs of later developmental

impairments [8]. A current challenge in the field of physical therapy is identifying neuromotor impairment accurately and early enough so that interventions can take place before the developmental delay is pronounced. One approach is tracking early spontaneous movements which have been shown to be correlated to future motor control [17]. Specifically, researchers propose that neurological deficits could be identified by collecting high quality spontaneous movement patterns from wearable sensors and kinematic analysis systems [9, 10, 19].

Using sensors and kinematic analysis systems across short sessions (5-10 minutes), studies have demonstrated that kinematic variables, such as kicking frequency, spatiotemporal organization, and interjoint and interlimb coordination are different between infants with typical development (TD) and infants at risk (AR) including infants with intellectual disability [13], myelomeningocele [21, 23], Down syndrome [15], as well as infants born preterm [7].

The goal of our work is to use day-long kinematic data collected from wearable inertial sensors to predict the infant's developmental classification as TD or AR. A prediction algorithm that is able to classify between TD and AR would further establish and expand on the initial results of demonstrated differences in spontaneous movements of infants with TD and AR, supporting the potential of full-day wearable sensor data in future research and prediction algorithms.

The remainder of the paper is organized as follows. Section 2 discusses the methodology adopted, data preprocessing, and classification. Section 3 provides the results of various methods and compares them to more interpretable methods. Section 4 discusses the future steps of this project.

## 2 METHODOLOGY

### 2.1 Data

The dataset was provided by the Infant Neuromotor Control Laboratory in the Division of Biokinesiology and Physical Therapy at the University of Southern California. It contains day-long inertial movement data from Opal sensors affixed on each infant ankle using leg warmers with pockets. Past work has validated that these sensors can accurately record quantity of infant limb movements and do not impede or promote infant movement [11, 24]. The dataset contains 19 movement features that summarize the general movement characteristics of each infant across the time the infant was awake during that day. The dataset consists of data from 12 TD infants and 24 AR infants. Each infant was tested 3 times from

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

*KDD'18, August 2018, London, England*

© 2018 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06.

[https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

1 to 16 months (with the exception of 1 AR infant with only two sessions) at various developmental stages.

## 2.2 Data Preprocessing

The sensors collected raw movement data from an accelerometer, gyroscope, and magnetometer. In post-analysis, we used an established algorithm that detects each leg movement from the raw sensor data [24, 26]. Next, the movement's duration, average acceleration, peak acceleration, and type (unilateral, bilateral synchronous, or bilateral asynchronous) were determined by the algorithm.

The dataset contains 19 movement features which can be broken down into 4 categories: movement count, duration, average acceleration, and peak acceleration. Specifically, the dataset includes the overall movement count and type of each movement for both the right and left legs. Movement count is presented as movements per hour of awake time and type is presented as a percentage of each type of movement made out of the total number of movements made. The dataset also includes the average and standard deviation of movement duration, average acceleration, and peak acceleration over all the infant's movements across the full day. Additionally, the movement features, age, developmental scale score (specifically the Alberta Infant Motor Scale [18]), and the AR/TD label were also included. Due to our small sample size (12 TDs and 24 ARs), we treated each new visit as a unique data record allowing for 107 samples of day-long leg movement, with 37 TDs and 100 ARs.

First, we normalized the data to correctly reflect the infants' movements due to each sensor being active for a different amount of time during the data collection. We then split the dataset into two groups, 0 to 6 months and 6 to 12 months, because there are noticeable changes in the movement features after 6 months and important developmental milestones emerge after this time for infants [6]. Also, the few samples over 12 months were each from an AR infant, so we removed these samples to reduce overfitting around age. Thus, we had 16 TD and 15 AR samples for 0 to 6 months and 23 TD and 38 AR samples for 6 to 12 months. The class weights for TD and AR were balanced for each algorithm used.

## 2.3 Feature Extraction

We extracted the relevant features out of the 22 to avoid unnecessary and intercorrelated inputs. We applied three feature extraction methods (univariate feature selection [22], recursive feature elimination (RFE) [14], and stepwise feature selection [16]) separately on the 0-6 month and 6-12 month dataset. For each feature subset selected by the algorithms, we removed highly correlated features to reduce the dimensionality of the feature set.

## 2.4 Classification

We used a binary classification approach for our models. To classify infants as TD or AR, various classification algorithms were spot checked using their default parameters including Random Forest [3], Logistic Regression, Support Vector Machines (SVM) [2], K-Nearest Neighbors (KNN) [12], Decision Trees [20], and AdaBoost [5].

We improved the performance by using grid search for tuning hyper parameters for each method, and the three top performing algorithms were put into a voting ensemble with a majority vote wins

principle. The reason for ensembling is to level individual algorithm bias and to reduce overfitting by smoothing the results [25].

## 3 RESULTS

### 3.1 Evaluation

We used a leave-one-out cross validation for model selection, a common approach for small datasets [4]. Standard metrics such as accuracy, precision, recall, and F1 score were used to evaluate the classification algorithms. We focused on choosing algorithms with higher true positives (TP) and false negatives (FN) because predicting AR infants correctly is more important for initial diagnosis purposes.

### 3.2 Experiments and Outcomes

**3.2.1 0 to 6 Month Prediction Algorithms.** Based on the feature extraction methods described above, we used eight movement features: mean acceleration, peak acceleration, unilateral movements, and total movements for the right and left leg. The results from the tuned algorithms varied from 90% average accuracy for SVM to 67.7% for AdaBoost. To reduce overfitting, the SVM was ensembled with two other well performing algorithms: AdaBoost and Logistic Regression. The resulting ensemble received an 83.9% accuracy with a 84% F1 score. Every AR infant was predicted correctly. Tables 1-6 displays the outcomes of the three individual algorithms, a decision tree, the ensemble, and the baseline algorithm which is the average of ten separate weighted random predictions respectively.

**Table 1: Weighted Average Baseline**

Class	Accuracy	Precision	Recall	F1 Score
TD	.625	.53	.62	.57
AR	.4	.5	.4	.44
Average	.516	.51	.52	.51

**Table 2: Decision Tree**

Class	Accuracy	Precision	Recall	F1 Score
TD	.75	.75	.75	.75
AR	.733	.733	.733	.733
Average	.742	.742	.742	.742

**Table 3: SVM**

Class	Accuracy	Precision	Recall	F1 Score
TD	.813	1	.81	.9
AR	1	.83	1	.91
Average	.903	.92	.9	.9

**Table 4: AdaBoost**

Class	Accuracy	Precision	Recall	F1 Score
TD	.625	.71	.62	.67
AR	.733	.65	.73	.69
Average	.677	.68	.68	.68

**Table 5: Logistic Regression**

Class	Accuracy	Precision	Recall	F1 Score
TD	.625	.77	.62	.69
AR	.8	.67	.8	.73
Average	.71	.72	.71	.71

**Table 6: Ensemble:SVM, AdaBoost, and LogReg**

Class	Accuracy	Precision	Recall	F1 Score
TD	.688	1	.69	.81
AR	1	.5	1	.86
Average	.839	.88	.84	.84

3.2.2 *6 to 12 Month Prediction Algorithms.* Using the features above, the same tests were conducted. The individual algorithms' accuracies ranged from 83% for AdaBoost to 67.2% for SVM and KNN. An ensemble was made to reduce overfitting using SVM, Adaboost, and Random Forest with majority vote as the winning prediction. The ensemble ended with a 77% accuracy and 77% F1 Score. Tables 7-12 list the results for the weighted average baseline, a decision tree, each individual algorithm in the ensemble, and the ensemble.

**Table 7: Weighted Average**

Class	Accuracy	Precision	Recall	F1 Score
TD	.391	.35	.39	.37
AR	.553	.6	.55	.58
Average	.492	.5	.49	.5

**Table 8: Decision Tree**

Class	Accuracy	Precision	Recall	F1 Score
TD	.478	.611	.478	.537
AR	.816	.721	.816	.765
Average	.689	.68	.689	.679

**Table 9: SVM**

Class	Accuracy	Precision	Recall	F1 Score
TD	.609	.56	.61	.58
AR	.711	.75	.71	.73
Average	.672	.68	.67	.67

**Table 10: AdaBoost**

Class	Accuracy	Precision	Recall	F1 Score
TD	.739	.81	.74	.77
AR	.895	.5	.89	.87
Average	.83	.83	.84	.83

**Table 11: Random Forest**

Class	Accuracy	Precision	Recall	F1 Score
TD	.609	.64	.61	.62
AR	.789	.77	.79	.78
Average	.721	.72	.72	.72

**Table 12: Ensemble:SVM,AdaBoost,Random Forest**

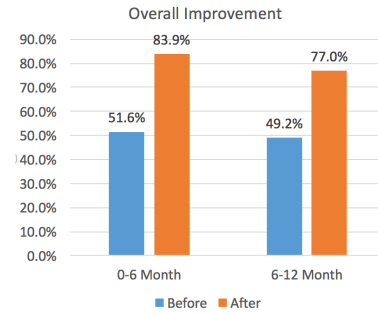
Class	Accuracy	Precision	Recall	F1 Score
TD	.608	.74	.61	.67
AR	.868	.79	.87	.82
Average	.77	.77	.77	.77

3.2.3 *Overall Outcome.* Overall, our ensemble approach in the 0-6 and 6-12 model resulted in a 32.3% and 27.8% increase in accuracy, respectively, when compared to the baseline prediction. Fig. 1 graphs the overall improvement in accuracy.

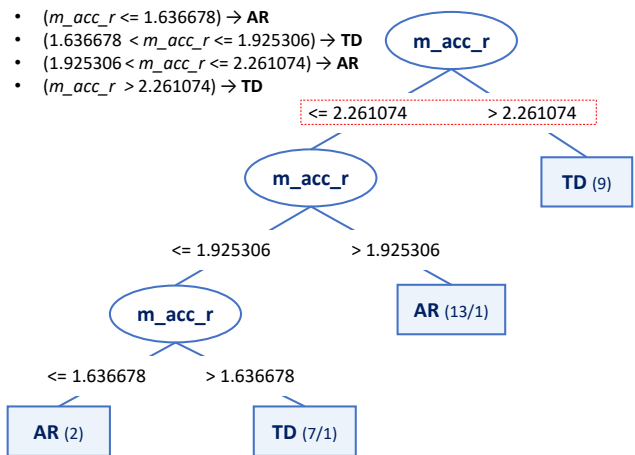
### 3.3 Accuracy versus Interpretability

In selecting an ensemble model, there is an inherent trade-off between increasing accuracy and decreasing interpretability. We need our model to be as accurate as possible, but we also need interpretable results for medical professionals where interpretability is necessary for professional insight and patient and parent trust [27]. Simpler algorithms, such as decision trees, are easily interpretable compared to ensemble models. We generated several decision trees to provide explanations to medical specialist on the importance of specific features in predicting TD and AR.

Specifically, Fig 2 shows the model decision tree resulting from executing algorithm J48 (C4.5) for the first age range (0 to 6 months). Unlike the second age range (6 to 12 months), this model obtained the accuracy of 70.75% with an improvement significantly better than the baseline. The peculiarity of this model is that the generated tree uses a single feature of the data set that is the mean acceleration of the right leg ( $m\_acc\_r$ ). Specifically, a mean acceleration value higher than 2.261074 was always TD implying some connection with higher accelerations and typical development.



**Figure 1: Comparison of accuracies between the baseline and ensemble models for 0-6 and 6-12 months.**



**Figure 2: The top-down model generated with the 0 to 6 months dataset using J48 (C 4.5).**

## 4 DISCUSSIONS

In this work, we investigated the use of binary classification algorithms applied to day-long sensor movement data to predict typical or at risk developmental status in infants. Our analysis was split for 0-6 months and 6-12 months, and both final algorithms performed significantly better than our baseline predictions. It is clear from our analysis that from 0 to 6 months the kinematic data of AR infants are more distinct than TD, but from 6 to 12 months it becomes hard to separate the two classes. This could be due to the small sample size. Overall, these results further establish the relationship between the kinematic features used in the classification algorithms and infant development.

Working with interpretable and understandable models is very important in order to obtain knowledge that can be shared with the medical experts. The top-down model generated in this work shows the importance of the mean acceleration in the prediction of AR infants. In fact, the model has found a possible threshold value (2.261074) above which typical development can be considered for an infant between 0 and 6 months. Future studies are necessary to determine the robustness of this classifier and test for relationships to developmental outcomes.

The nature of this dataset, similar to all small datasets, leads to the possibility of overfitting as well as lacking representation of the general population. When using a leave one out cross validation, the trained models may be picking up noise beyond the true signal of the pattern which could cause the model to perform poorly on new unseen instances. We paired our final ensemble model with the decision tree to provide more interpretable results for healthcare professionals.

There are several next steps we are pursuing. We would like to recruit more TD and AR infants for this study. This requires significant longitudinal interaction with each infant such as traditional developmental assessment tests, movement sensor tests, and a check-up at two years old to determine if the infant has a diagnosed developmental delay.

Currently, each infant visit is treated independently despite having three samples per infant. A larger dataset would allow us to divide the samples by infant rather than by visit, providing more data about each infant and potentially improving our predictions. We could also create an algorithm that predicts some movement feature of an infant at visit three based on the infant's first two visits. This prediction algorithm would be helpful if an infant is participating in a new intensive intervention after their second visit, and the healthcare professional wants to compare the infant's movement results after the intervention to the infant's predicted movement pattern from the algorithm that assumes no intervention.

The eventual goal is to use this methodology to predict if an at risk infant will be diagnosed with a developmental delay. Currently, developmental delays are often not diagnosed until an infant is two years old, preventing many infants from receiving early targeted interventions. Current tools, such as the AIMS score, detect early signs of atypical development, but they can only detect the extreme cases. The prediction algorithm we aim to develop would confirm that developmental delays are reflected in the movement of infants in the first few months of infancy and thus allow for more infants to receive earlier, directed interventions.

## ACKNOWLEDGMENTS

This work was funded in part by the American Physical Therapy Association Academy of Pediatric Physical Therapy Research Grant 1 and 2 Awards (PI: Smith) and in part by NSF award 1706964 (PI: Smith, Co-PI: Matarić). Study procedures were approved by the Institutional Review Boards of the Oregon Health & Science University and the University of Southern California.

## REFERENCES

- [1] G. R. Alexander, M. Kogan, D. Bader, W. Carlo, M. Allen, and J. Mor. 2003. US birth weight/gestational age-specific neonatal mortality: 1995-1997 rates for whites, hispanics, and blacks. *Pediatrics* 111, 1 (2003), 61A-66.
- [2] B. E. Boser, I. M. Guyon, and V. N. Vapnik. 1992. A training algorithm for optimal margin classifiers. In *Computational Learning Theory*. 144.
- [3] L. Breiman. 2001. Random Forests. *Machine Learning*, 45, 1 (2001).
- [4] G. C. Cawley. 2006. Leave-One-Out Cross-Validation Based Model Selection Criteria for Weighted LS-SVMs. *International Joint Conference on Neural Networks* (2006).
- [5] Y. Freund and R. E. Schapire. 1999. A Short Introduction to Boosting. *Journal of Japanese Society for Artificial Intelligence* 14, 5 (1999), 771-780.
- [6] E. Gajewska, M. Sobieszka, E. Kaczmarek, A. Suwalska, and B. Steinborn. 2013. Achieving Motor Development Milestones at the Age of Three Months May Determine, but Does Not Guarantee, Proper Further Development. *Scientific World Journal* (2013).
- [7] J. J. Geerdink, B. Hopkins, W. J. Beek, and C. B. Heriza. 1996. The organization of leg movements in preterm and full-term infants after term age. *Developmental Psychobiology* 29, 4 (1996), 335-351.
- [8] A. Ghassabian, R. Sundaram, E. Bell, S. C. Bello, C. Kus, and E. Yeung. 2016. Gross Motor Milestones and Subsequent Development. *Pediatrics* 138, 1 (2016).
- [9] S. E. Groen, A. C. de Blecourt, K. Postema, and M. Hadders-Algra. 2005. General movements in early infancy predict neuromotor development at 9 to 12 years of age. *Developmental Medicine and Child Neurology* 47, 11 (2005), 731A-738.
- [10] M. Hadders-Algra and A. M. Groothuis. 1999. Quality of general movements in infancy is related to neurological dysfunction, ADHD, and aggressive behaviour. *Developmental Medicine and Child Neurology* 41, 6 (1999), 381-391.
- [11] C. Jiang, C. J. Lane, E. Perkins, D. Schiesel, and B. A. Smith. 2017. Determining if wearable sensors affect infant leg movement frequency. *Developmental Neurorehabilitation* 5, 771 (2017), 1-4.
- [12] W. R. Klecka. 1980. Discriminant Analysis. 19, 1 (1980).
- [13] M. Kouwaki, M. Yokochi, T. Kamiya, and K. Yokochi. 2014. Spontaneous movements in the supine position of preterm infants with intellectual disability. *Brain and Development* 36, 7 (2014), 572-577.
- [14] Y. Mao, D. Pi, Y. Liu, and Y. Sun. 2006. Accelerated Recursive Feature Elimination Based on Support Vector Machine for Key Variable Identification. *Chinese Journal of Chemical Engineering* 14, 1 (2006), 65-72.
- [15] S. M. McKay and R. M. Angulo-Barroso. 2006. Longitudinal assessment of leg motor activity and sleep patterns in infants with and without Down syndrome. *Infant Behavior and Development* 29, 2 (2006), 153-168.
- [16] M. McShane, B. Cameron, G. Cotea, M. Motamedib, and C. Spiegelman. 1999. A novel peak-hopping stepwise feature selection method with application to Raman spectroscopy. *Analytica Chimica Acta* 388, 3 (1999), 251-264.
- [17] J. P. Piek and D. J. 1998. The influence of preterm birth on early motor development. (1998).
- [18] M. C. Piper and J. Darrah. 1994. Motor Assessment of the Developing Infant. *WB Saunders* (1994).
- [19] H. F. Prechtl. 1997. State of the art of a new functional assessment of the young nervous system. An early predictor of cerebral palsy. *Early Human Development* 50, 1 (1997), 1-11.
- [20] J. R. Quinlan. 1986. Induction of Decision Trees. *Machine Learning* 1, 1 (1986).
- [21] N. Rademacher, D. P. Black, and B. D. Ulrich. 2008. Early spontaneous leg movements in infants born with and without myelomeningocele. *Pediatric Physical Therapy* 20, 2 (2008), 137-145.
- [22] Inza I. Larranaga P. Says, Y. 2007. A review of feature selection techniques in bioinformatics. *Bioinformatics* 23, 19 (2007), 2507-2517.
- [23] B. A. Smith, C. Teulier, J. Sansom, N. Stergiou, and B. D. Ulrich. 2008. Approximate entropy values demonstrate impaired neuromotor control of spontaneous leg activity in infants with myelomeningocele. *Pediatric Physical Therapy* 23, 3 (2008), 241-247.
- [24] Beth A. Smith, Ivan A. Trujillo-Priego, Christianne J. Lane, James M. Finley, and Fay B. Horak. 2015. Daily quantity of infant leg movement: Wearable sensor algorithm and relationship to walking onset. *Sensors* 15, 8 (2015), 19006-19020.
- [25] P. Sollich and A. Keogh. 1995. Learning with ensembles: How over-fitting can be useful. *Advances in Neural Information Processing Systems* (1995).

- [26] Ivan A Trujillo-Priego and Beth A Smith. 2017. Kinematic characteristics of infant leg movements produced across a full day. *Journal of Rehabilitation and Assistive Technologies Engineering* 4, 1 (2017), 1–10.
- [27] A. Vellido, J. D. Martin-Guerrero, and P. J.G. Lisbon. 2012. Making machine learning models interpretable. *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*. (2012).