
Detecting early signs of depressive and manic episodes in patients with bipolar disorder using the signature-based model

Andrey Kormilitzin¹ Kate E.A. Saunders^{2,3} Paul J. Harrison^{2,3} John R. Geddes^{2,3} Terry Lyons¹

Abstract

Background: Recurrent major mood episodes and subsyndromal mood instability cause substantial disability in patients with bipolar disorder. Early identification of mood episodes enabling timely mood stabilisation is an important clinical goal. Recent technological advances allow the prospective reporting of mood in real time enabling more accurate, efficient data capture. The complex nature of these data streams in combination with challenge of deriving meaning from missing data mean pose a significant analytic challenge. The signature method is derived from stochastic analysis and has the ability to capture important properties of complex ordered time series data.

Objective: To explore whether the onset of episodes of mania and depression can be identified using self-reported mood data.

Methods: Self-reported mood data were collected from 261 participants with bipolar disorder using the True Colours monitoring system. Manic and depressive episodes were defined as an ASRM score 6 or a QIDS score 10, respectively. The signature method was used to extract features from a rolling window of k -weeks, where $k = \{4, 6, 8, 12, 20, 50\}$. Two independent models for prediction of depressive and manic mood episodes were trained using logistic regression with the elastic net regularisation scheme to reduce overfitting. The stability of the generalisation error metrics of the predictive models was estimated by pooling 100 repetitions of 10-fold cross validation.

Results: The signature method on average accurately predicted 79.2% of precursors to depres-

sive episodes with a sensitivity of 76.9% and specificity of 79.5% (PPV=78.6%, AUC=0.86) and 71.9% of precursors to manic episodes with a sensitivity of 73.3% and specificity 79.2% (PPV=77.1%, AUC=0.83). This was more accurate than predictions (accuracy 77.4% for depression and 69.5% for mania, p -value < 0.001) based up on a model comprises three features: the mean value, variability and a number of missing responses within the window.

Conclusions: The signature method offers a systematic approach to the analysis of longitudinal self-reported mood data. The accuracies of the signature-based linear models are considerably better than linear models based on manually crafted features and the models have the potential to significantly enhance self-management and clinical care in bipolar disorder.

1. Introduction

1.1. Related Work

Predicting the future outcomes on the basis of historical observations is a long standing challenge in many scientific areas. The accurate prediction of future symptoms in patients with bipolar disorder could facilitate timely clinical interventions enhancing patients quality of life, altering the course of disease and could have economic benefits for the national health service (Manning, 2005). Several predictive models have been proposed using the self-reported longitudinal mood data from psychometric questionnaires. These models included: autoregressive linear models (Moore et al., 2012; 2014), relaxation oscillator framework (Bonsall et al., 2015), Kalman filter with the data from geographic location (Tsanas et al., 2016; Palmius et al., 2017) and an effective subgroup selection among diverse cohort of patients (Palmius & De Vos, 2016). However, these methods focused only on predicting the numerical value of the rating scale. In this work, we took a different route in building a predictive model. In contrast to previous approaches, we aim to identify the precursors of upcoming episodes of depression and mania, using the self-

¹Mathematical Institute, University of Oxford, Andrew Wiles Building, Woodstock Rd, Oxford OX2 6GG, UK ²Department of Psychiatry, University of Oxford ³Oxford Health NHS Foundation Trust, Warneford Hospital, Oxford OX3 7JX, UK . Correspondence to: <andrey.kormilitzin@maths.ox.ac.uk>.

reported rating scale. In fact, we solve a binary classification problem, where our developed model, using an interval of k -consecutive self-reported observation, estimates a probability of being a precursor to an episode.

2. Methods

2.1. Data

The data were collected as part of the OXTEXT-1 study (Bilderbeck et al., 2017). The participants completed standardised questionnaires on a weekly basis using the True Colours mood monitoring system after receiving a text or email prompt. The data were collected in an observational manner and independent from the clinical care.

Self-reported mood data using the Quick Inventory of Depressive Symptoms (QIDS-SR16) (Rush et al., 2003) and Altman Self-Rating Mania scale (ASRM) (Altman et al., 1997). A depressive episode was defined as a QIDS score of above 11 for at least two consecutive weeks. A manic episode was defined as an AMRS score above a threshold 6 for at least one week (Kessler et al., 2005).

2.2. Demographics and Patient Selection

The original cohort contained 286 subjects. Data collected from 01-Jan-2012 until 31-Dec-2016 was included in this analysis. We excluded 25 participants from the analysis: 22 participants had unconfirmed diagnoses and 3 participants withdrew consent. Of the 261 included participants 148 and 113 subjects were diagnosed with bipolar type I and type II, (denoted by BP-I and BP-II) respectively. All identical duplicate values were removed. If there were multiple responses within a week, only the first response was considered.

Additionally, we excluded patients who stayed less than 5 weeks in self-monitoring and all of those who had no recorded mood episodes during the study period. From the exploratory analysis, we found that only 59% of patients (155/261) in this cohort had episodes of depression (spent at least two weeks in depression) and 69% (181/261) had manic episodes (at least one week). The reported ethnicity and education level in accordance with the standard classification. One male from the BP-II group did not report his ethnicity. One male from the BP-I group, one male from the BP-II group, two females, one BP-I and one BP-II did not report their education levels. Percents of relevant groups may not add up to 100 due to rounding errors. The adjusted adherence corresponds to the period of observations between the first and the last actual responses. The demographic data summarised in Table 1.

2.3. Signature-based Predictive Model

We used a rolling window analysis to learn temporal dependencies in the data. This is a commonly used method to estimate changes and reveal patterns in sequential data over time (Zivot & Wang, 2007). The rolling window analysis allows to estimate the variability of data within certain time intervals and to compare between them.

The proposed predictive model is based on logistic regression, which learns patterns in data within the window of size k -weeks to predict the future occurrence of an episode at one-week time horizon. Two linear logistic models were trained independently to predict depressive and manic episodes. Using the definition of depressive and manic episodes, we developed an algorithm which identifies episodes in patients data and marks the weekly observations with a binary label (0 or 1) corresponding to the absence or presence of an episode.

The features used in logistic regression to predict outcomes, were extracted using the signature method (Chevyrev & Kormilitzin, 2016). The novel signature method, from stochastic analysis, has the ability to capture important properties of complex time-ordered data. It maps N unique data streams into a single piece-wise continuous path in N dimensions with the iterated integrals of such path representing a feature set (signature features) of the data. The signature is an infinite sequence of ordered iterated integrals and in practice. For machine learning applications and data analyses, we truncate the infinite sequence at some level L by taking only first terms and use them as features.

The signature method allows missing responses to be incorporated into the analysis by introducing a new indicator binary variable, which takes values 1 or 0, indicates whether a response is missing or present respectively (Little & Rubin, 2014). The indicator variable stream is then combined with N data streams of interest and mapped into a lifted path in $N+1$ dimensions. Practically, we used the self-reported data streams (QIDS and Altman questionnaires) together with their missing responses indicator data streams. The detailed description, rigorous mathematical foundations and algorithms of the signature methods are beyond the scope of the current paper and have been extensively covered in (Kormilitzin et al., 2016; Levin et al., 2013; Gyurkó et al., 2013; Yang et al., 2015; Lyons, 2014; Király & Oberhauser, 2016; Yin et al., 2013; Graham, 2013; Chen, 1957; Lyons & Xu, 2011; Lyons et al., 2007; Hambly & Lyons, 2010; Hairer, 2014; Yang et al., 2016; Xie et al., 2016; Lai et al., 2017) and references therein.

In order to estimate the generalisation error of predictive models, we split the entire data set into two non-overlapping groups of patients for training (67%) and test-

ing (33%) data sets, while preserving the ratio of weeks in episodes to weeks without episodes in both sets. The stratified splitting strategy allowed the variance and bias of a model to be minimised (Kohavi et al., 1995; Cawley & Talbot, 2010). A model was trained and cross-validated using the data only from one group and then the data from the withheld second group of patients was used to assess the classification metrics. The process of random splitting of patients into two disjoint groups was repeated 100 times to estimate the uncertainty of the metrics.

The elastic net regularisation scheme (Zou & Hastie, 2005) was used to prevent the models from overfitting the data. The elastic net regularisation introduces the convex combination of L1 and L2 penalties, mitigating the problem of multicollinearity of features and improving the stability of a classifier. As an additional preprocessing step, the signature features matrix was standardised along columns to have zero mean and unit variance. To account for the variability in longitudinal data, we used the lead-lag transformation, where the lead transform of the signal is paired with its lagged transform of the signal. To assess the performance of the classification procedure we computed accuracy, sensitivity, specificity, and positive predicted value (PPV). Additionally, we used the area under the receiver operating characteristic curve (AUC) to assess the performance of the classification models at different values of a threshold.

We compared the signature-based model (Sig) to four baseline models, where three of them based only on one of the mean value (Mean), variability (Rmssd) and a number of missing responses (MissRes) within the window predictors and the fourth model comprises all these three predictors simultaneously in a linear combination (MRM). The variability was computed, ignoring the missing responses, using the root mean square successive distance (RMSSD) (Electrophysiology, 1996), which captures the temporal dependencies in the data.

2.4. Parameters of the Model

The longitudinal self-reported data were partitioned into k-week intervals, transformed into a set of features using the signature method and the resulting features were used as inputs to logistic regression to predict a binary outcome, presence or absence of an episodes, at the week following the interval. The intervals that precede the episodes (precursors to episodes) are labelled as positives (1). The data from QIDS and AMSR questionnaires were used to train models predicting respectively depressive and manic episodes. Intervals that precede missing responses or contain only one non-missing response were excluded. Additionally, we sought to optimise the classification model using the hyperparameters of the penalised logistic regres-

sion during the 10-fold cross-validation phase with the parameter $\lambda \in \{0, 0.1, 0.2, 0.3, \dots, 1.0\}$ that mixes L1 and L2 penalties.

We used the Python Pandas package (version 0.20.1) (McKinney et al., 2010) for statistical analysis, data manipulations and processing, Python Scikit-learn package (version 0.18.1) (Pedregosa et al., 2011) for implementing machine learning tasks and Matplotlib for plotting and graphics (version 2.0.1) (Hunter, 2007).

3. Results

3.1. Predicting Depressive and Manic Episodes

The performance of the signature-based model and comparison to the baseline models for depression and mania trained on corresponding test sets are summarised in Tables 2-3 respectively. Classification results evaluated on the corresponding test sets of the models for depression and mania using 100 repetition of random splitting to train/test (67%/33%) sets. The features used in the model are elements of a truncated signature at level 2. The repeated results of metrics were nearly normal distributed and following the suggested procedure (Salzberg, 1997; Dietterich, 1998), we used the corrected paired t-test to pair-wise compare the performance of all metrics of five models. We found a strong statistical significance of their differences (p -value < 0.001), also for depression and mania models. The stability of the classification metrics was estimated by repeating the training and testing the models 100 times and the results presented as mean and standard deviation. The classification metrics of models for depression and mania as function of the size of the rolling window are presented in Figures 1-2.

3.2. Experimental Clinical Applications

The proposed models for classification of intervals are potentially useful for clinical insights and for identifying particular states, where a patient is in transition between the episodes. First, we are interested in examining whether it is possible to distinguish between the intervals which are close to an episode and those which are further away. For the sake of simplicity and demonstration of the conceptual approach, we consider intervals of fixed length of size 6 weeks. Two types of intervals are: a wellness and a precursor interval, which are spaced 14 and n weeks prior to the beginning of an episode respectively. By varying $n = \{0, 1, 2, 3, 4, 6\}$, and applying the developed models, we will measure the area under the ROC curve (AUC) to assess the classification performance. We restrict the wellness and the precursor intervals from overlapping. For classification procedure we considered only the signature-based model. To estimate the variability due to split of

the data into training and testing sets (67% / 33%) we repeated the process 100 times. The results are presented as mean and standard deviation. The AUC estimations are summarised in Table 4.

4. Discussion

4.1. Principal Results

In this study we have presented an application of the signature method for modelling and prediction of mood episodes in bipolar disorder. We demonstrated that the results of the signature-based model outperform four other models based on the manually selected predictors. The main advantage of the signature method is that it allows to systematically combine multimodal longitudinal data streams, including the distribution of missing responses, and to extract features used in predictive models, avoiding a complicated feature engineering process. The signature terms faithfully represent the underlying data, have an interpretable geometrical meaning as functions of data and should be used as canonical features in machine learning and data analysis tasks. Intuitively, one can think of the signature terms as of an ordered collection of the sample statistical moments. The lead-lag transformation naturally captures the successive variability in the data, avoiding the need for designing a special feature for that task as we did with the RMSSD.

The plot of the classification metrics (Figures 1-2) demonstrate the dependence of the predictive performance of the models on the size of the rolling window. For small values of the window size ($k < 12$) all models perform similar, while the signature-based model, MRM and MissRes, continue improving as the number of historical observations grows and the Mean and Rmssd models decline or saturate. The predictive performance does ultimately depend on the ability of models to account for the temporal dependencies. The signature model is based on the unique set of predictors derived from a trajectory (a path) of the longitudinal self-reported responses, where the predictors account for the sequential dependencies through the path integrals and thus more observations lead to a higher accuracy of predictions. The Mean model, in contrast with the others, strongly declines when used with a large number of historical observations. That manifests the fact, that the average score over a long interval of observations is a poor predictor of future episodes, because it does not take into account the temporal dependencies. Interestingly, it appears that the variability of the rating scale along is a very poor predictor of future episodes, even though it does account for the sequential dependencies in the data through the RMSSD metric. Another interesting observation is that the number of missing responses along within the window can predict with up to 80% accuracy the future episodes of depression and mania. Finally, we observed that the composite model

(MRM) that comprises all three predictors linearly, underperforms the accuracy of the signature-based model on average by 2.36%.

The developed models aim to discover the temporal patterns of mood episodes and accurately predict the possible future outcomes. The proposed algorithms can be potentially deployed in modern health-monitoring platforms and serve as a patient self-management tool. Detection of precursors to upcoming mood episodes and early intervention of the healthcare professionals can help to reduce the severity of symptoms in patients with bipolar disorder. However, the developed predictive models do not take into account the unique personal traits of patients, for example, the distribution of missing values and the item responses. The main problem of training an individual model is due to relatively small number of episodes each participant experiences during the self-monitoring.

While the results of the current work have demonstrated the feasibility of building the predictive models based on the signature transformation of the self-reported data alone, more in-depth research is needed to refine the model. The current analysis considered only a relatively small number of patients who satisfied the inclusion criteria and in order to develop robust models, in the future we are planning to replicate the proposed method using a larger cohort. It has been shown [40] that training a model on a selected subgroup of patients who share common traits leads to a significant improvement of accuracy of the predictive models. We are planning to address the problem of clustering patient in our future works and retrain the predictive models for each cluster.

Apart from the identifying early signs of deterioration, it is also clinically interesting to identify early signs of improvement. The approach to find precursors to wellness is conceptually similar to the one we presented in this work, but requires different labelling of intervals and retraining the models.

The developed approach to interval classification allows us potentially to identify the early signs of transitions between the mood states in patients. We tested whether we can distinguish between an interval which is not in an episode and an interval which is n -weeks close to the beginning of an episode. The results presented in Table 4 indicate that intervals closer to an impending are distinct than those which are far apart. However, the choice of the interval size (6-weeks) and the subjective definition of the wellness intervals (14 weeks prior to an episode) may not be optimal to conclude the results. Further investigation needed.

In this work we used only linear models for prediction of future episodes. Linear models preferable over more complex and non-linear ones, especially with medical data, as

they allow to understand clearly which variables used as predictors and their influence on the dependent variable. However, the important clinical condition which has not yet been addressed through our research is the case of mixed episodes, where both depression and mania may occur simultaneously (Vieta & Valentí, 2013). The signature method allows to combine and develop a model using both QIDS and AMSR data streams to predict the future states including the missing ones. We are planning to address this problem in a future work.

4.2. Limitations

Notwithstanding the signature approach outperforms other models based on the manually created predictors, the full potential of the signature method, which has been shown to be superior to state-of-the-art methods in other fields (Yang et al., 2017) has not yet been demonstrated on the longitudinal mood data. We refer the limited power of the proposed method to the challenging nature of the self-reported longitudinal data, e.g. presence of non-randomly distributed missing values and more importantly, to the diversity of the cohort and lack of regularly time-stamped records of exogenous factors (for example, medications and hospitalisations). The moderate accuracies of the models are also influenced by the assumption that all patients share the same patterns of early signs of deterioration, which might not be true. That assumption allowed us to train and test the models on randomly selected non-overlapping sets of patients from the entire cohort. With more data collected over longer period of the self-monitoring, we are interested in construction of individual models, to examine the diversity in patients. Additionally, the proposed method for incorporating of missing responses might not be optimal. In this work we laid a significant groundwork with a lot of open questions for further research and are planning to address the aforementioned problems in subsequent publications.

4.3. Conclusions

In this work we developed and compared several linear models for predicting the episodes of depression and mania in patients with bipolar disorders. We found that the accuracy of predictions depends on the length of historical observations (the size of the rolling window). The accuracy of the signature-based model is higher than all others by the amount of at least as 2.36% and 3.37% for depression and mania respectively. The signature method represents a systematic approach for feature extraction from and modelling functions on streams of data. We found that the linear combination of three predictors (Mean, Rmssd and MissRev) outperforms each of them individually for predicting the future episodes in both models for depression and mania.

Ethics

The study was approved by Oxfordshire Research Ethics Committee A (reference no: 10/H0604/13). All participants in gave written, informed consent.

Data Accessibility

Data can be requested from JRG

Competing Interests

PJH, AK, TJL, KEAS and JRG declare no competing interests.

Authors Contributions

JRG designed the trial. JRG coordinated the study. AK and TJL conducted the analysis. AK, TJL and KEAS drafted the paper, which was reviewed by all authors.

Acknowledgements

We thank Maarten De Vos, Guy Goodwin, Keltie McDonald, Nick Palmius and Athanasios Tsanas for valuable discussions.

Funding

AK, TJL, KEAS, JRG and PJH are supported by a Wellcome Trust Strategic Award CONBRIO: Collaborative Oxford Network for Bipolar Research to Improve Outcomes, Reference number 102616/Z. JRG, PH and KEAS are supported by the NIHR Oxford Health Biomedical Research Centre. TJL acknowledges the support of NCEO project NERC, ERC grant number 291244, EPSRC grant number EP/H000100/1 and by the Alan Turing Institute under the EPSRC grant EP/N510129/1. No funder had any role in the study design; data collection, analysis, or interpretation of data; writing of the report; or in the decision to submit the paper for publication. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR or the Department of Health.

References

- Altman, Edward G, Hedeker, Donald, Peterson, James L, and Davis, John M. The altman self-rating mania scale. *Biological psychiatry*, 42(10):948–955, 1997.
- Bilderbeck, Amy C, Atkinson, Lauren Z, Geddes, John R, Goodwin, Guy M, and Harmer, Catherine J. The effects of medication and current mood upon facial emotion recognition: findings from a large bipolar disorder cohort study. *Journal of Psychopharmacology*, 31(3):

320–326, 2017.

- Bonsall, Michael B, Geddes, John R, Goodwin, Guy M, and Holmes, Emily A. Bipolar disorder dynamics: affective instabilities, relaxation oscillations and noise. *Journal of The Royal Society Interface*, 12(112):20150670, 2015.
- Cawley, Gavin C and Talbot, Nicola LC. On over-fitting in model selection and subsequent selection bias in performance evaluation. *Journal of Machine Learning Research*, 11(Jul):2079–2107, 2010.
- Chen, Kuo-Tsai. Integration of paths, geometric invariants and a generalized baker-hausdorff formula. *Annals of Mathematics*, pp. 163–178, 1957.
- Chevyrev, Ilya and Kormilitzin, Andrey. A primer on the signature method in machine learning. *arXiv preprint arXiv:1603.03788*, 2016.
- Dietterich, Thomas G. Approximate statistical tests for comparing supervised classification learning algorithms. *Neural computation*, 10(7):1895–1923, 1998.
- Electrophysiology, Task Force of the European Society of Cardiology the North American Society of Pacing. Heart rate variability. *Circulation*, 93(5):1043–1065, 1996. ISSN 0009-7322. doi: 10.1161/01.CIR.93.5.1043. URL <http://circ.ahajournals.org/content/93/5/1043>.
- Graham, Benjamin. Sparse arrays of signatures for online character recognition. *arXiv preprint arXiv:1308.0371*, 2013.
- Gyurkó, Lajos Gergely, Lyons, Terry, Kontkowsky, Mark, and Field, Jonathan. Extracting information from the signature of a financial data stream. *arXiv preprint arXiv:1307.7244*, 2013.
- Hairer, Martin. A theory of regularity structures. *Inventiones mathematicae*, 198(2):269–504, 2014.
- Hambly, Ben and Lyons, Terry. Uniqueness for the signature of a path of bounded variation and the reduced path group. *Annals of Mathematics*, pp. 109–167, 2010.
- Hunter, John D. Matplotlib: A 2d graphics environment. *Computing In Science & Engineering*, 9(3):90–95, 2007.
- Kessler, Ronald C, Berglund, Patricia, Demler, Olga, Jin, Robert, Merikangas, Kathleen R, and Walters, Ellen E. Lifetime prevalence and age-of-onset distributions of dsm-iv disorders in the national comorbidity survey replication. *Archives of general psychiatry*, 62(6):593–602, 2005.
- Király, Franz J and Oberhauser, Harald. Kernels for sequentially ordered data. *arXiv preprint arXiv:1601.08169*, 2016.
- Kohavi, Ron et al. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Ijcai*, volume 14, pp. 1137–1145. Stanford, CA, 1995.
- Kormilitzin, AB, Saunders, KEA, Harrison, PJ, Geddes, JR, and Lyons, TJ. Application of the signature method to pattern recognition in the cequel clinical trial. *arXiv preprint arXiv:1606.02074*, 2016.
- Lai, Songxuan, Jin, Lianwen, and Yang, Weixin. Toward high-performance online hccr: A cnn approach with dropdistortion, path signature and spatial stochastic max-pooling. *Pattern Recognition Letters*, 89:60–66, 2017.
- Levin, Daniel, Lyons, Terry, and Ni, Hao. Learning from the past, predicting the statistics for the future, learning an evolving system. *arXiv preprint arXiv:1309.0260*, 2013.
- Little, Roderick JA and Rubin, Donald B. *Statistical analysis with missing data*. John Wiley & Sons, 2014.
- Lyons, Terry. Rough paths, signatures and the modelling of functions on streams. *arXiv preprint arXiv:1405.4537*, 2014.
- Lyons, Terry and Xu, Weijun. Inversion of signature for paths of bounded variation. *arXiv preprint arXiv:1112.0452*, 2011.
- Lyons, Terry J, Caruana, Michael, and Lévy, Thierry. *Differential equations driven by rough paths*. Springer, 2007.
- Manning, J Sloan. Burden of illness in bipolar depression. *Primary care companion to the Journal of clinical psychiatry*, 7(6):259, 2005.
- McKinney, Wes et al. Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference*, volume 445, pp. 51–56. SciPy Austin, TX, 2010.
- Moore, Paul J, Little, Max A, McSharry, Patrick E, Geddes, John R, and Goodwin, Guy M. Forecasting depression in bipolar disorder. *IEEE Transactions on Biomedical Engineering*, 59(10):2801–2807, 2012.
- Moore, Paul J, Little, Max A, McSharry, Patrick E, Goodwin, Guy M, and Geddes, John R. Mood dynamics in bipolar disorder. *International journal of bipolar disorders*, 2(1):11, 2014.

- Palmius, Niclas and De Vos, Maarten. Non-parametric subset selection for personalised time-series regression. In *Practical Bayesian Nonparametrics Workshop at the 30th Conference on Neural Information Processing Systems*, 2016.
- Palmius, Niclas, Tsanas, Athanasios, Saunders, KE, Bilderbeck, Amy C, Geddes, John R, Goodwin, Guy M, and De Vos, Maarten. Detecting bipolar depression from geographic location data. *IEEE Trans Biomed Eng*, 2017.
- Pedregosa, Fabian, Varoquaux, Gaël, Gramfort, Alexandre, Michel, Vincent, Thirion, Bertrand, Grisel, Olivier, Blondel, Mathieu, Prettenhofer, Peter, Weiss, Ron, Dubourg, Vincent, et al. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12(Oct):2825–2830, 2011.
- Rush, A John, Trivedi, Madhukar H, Ibrahim, Hicham M, Carmody, Thomas J, Arnow, Bruce, Klein, Daniel N, Markowitz, John C, Ninan, Philip T, Kornstein, Susan, Manber, Rachel, et al. The 16-item quick inventory of depressive symptomatology (qids), clinician rating (qids-c), and self-report (qids-sr): a psychometric evaluation in patients with chronic major depression. *Biological psychiatry*, 54(5):573–583, 2003.
- Salzberg, Steven L. On comparing classifiers: Pitfalls to avoid and a recommended approach. *Data mining and knowledge discovery*, 1(3):317–328, 1997.
- Tsanas, A, Saunders, KEA, Bilderbeck, AC, Palmius, N, Osipov, M, Clifford, GD, Goodwin, GM, and De Vos, M. Daily longitudinal self-monitoring of mood variability in bipolar disorder and borderline personality disorder. *Journal of affective disorders*, 205:225–233, 2016.
- Vieta, Eduard and Valentí, Marc. Mixed states in dsm-5: implications for clinical care, education, and research. *Journal of affective disorders*, 148(1):28–36, 2013.
- Xie, Zecheng, Sun, Zenghui, Jin, Lianwen, Ni, Hao, and Lyons, Terry. Learning spatial-semantic context with fully convolutional recurrent network for online handwritten chinese text recognition. *arXiv preprint arXiv:1610.02616*, 2016.
- Yang, Weixin, Jin, Lianwen, and Liu, Manfei. Chinese character-level writer identification using path signature feature, dropstroke and deep cnn. In *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*, pp. 546–550. IEEE, 2015.
- Yang, Weixin, Jin, Lianwen, Ni, Hao, and Lyons, Terry. Rotation-free online handwritten character recognition using dyadic path signature features, hanging normalization, and deep neural network. In *Pattern Recognition (ICPR), 2016 23rd International Conference on*, pp. 4083–4088. IEEE, 2016.
- Yang, Weixin, Lyons, Terry, Ni, Hao, Schmid, Cordelia, Jin, Lianwen, and Chang, Jiawei. Leveraging the path signature for skeleton-based human action recognition. *arXiv preprint arXiv:1707.03993*, 2017.
- Yin, Fei, Wang, Qiu-Feng, Zhang, Xu-Yao, and Liu, Cheng-Lin. Icdar 2013 chinese handwriting recognition competition. In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, pp. 1464–1470. IEEE, 2013.
- Zivot, Eric and Wang, Jiahui. *Modeling financial time series with S-Plus®*, volume 191. Springer Science & Business Media, 2007.
- Zou, Hui and Hastie, Trevor. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2): 301–320, 2005.

Detecting early signs of depressive and manic episodes in patients with bipolar disorder using the signature-based model

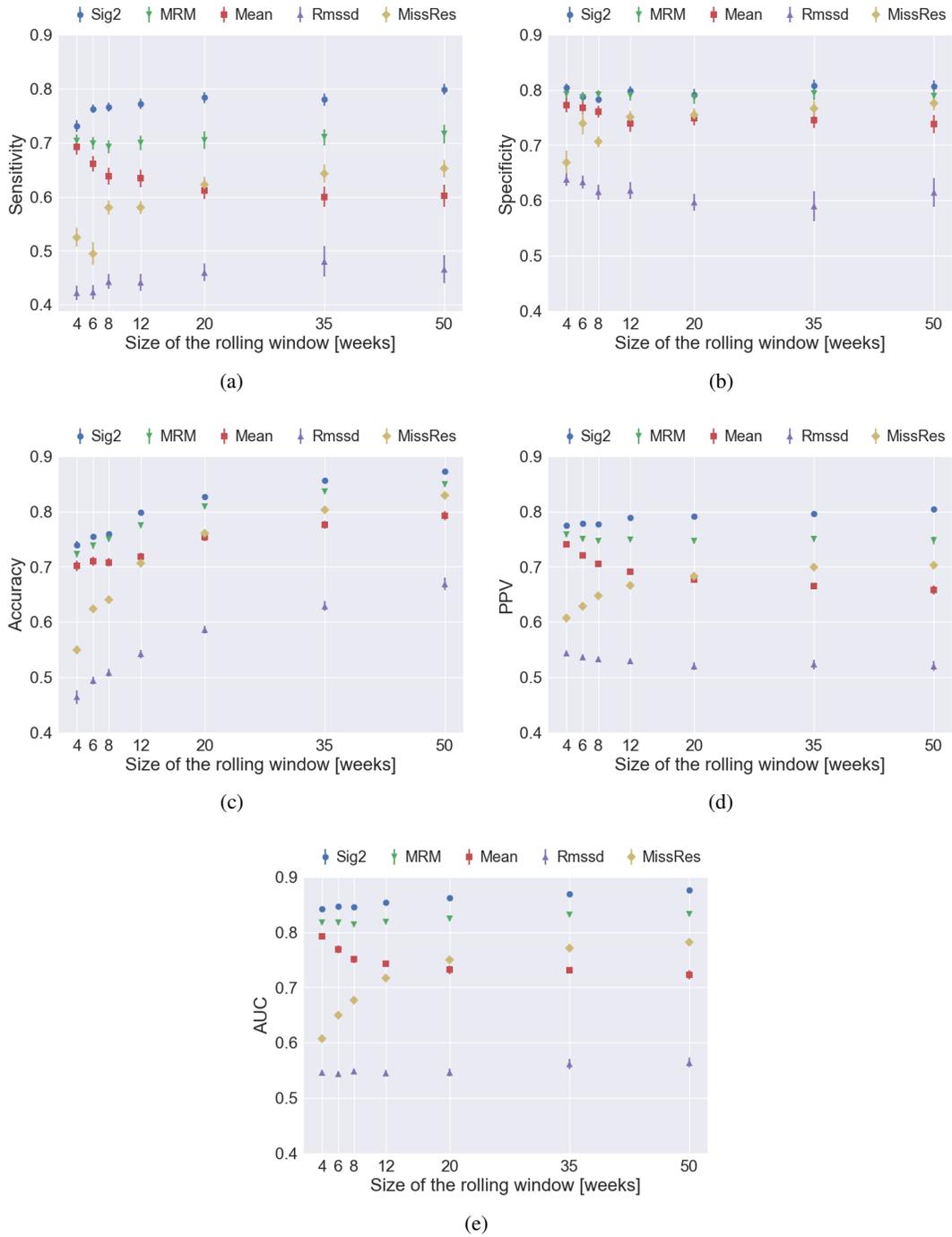


Figure 1. Classification metrics for depression model presented as mean with 95%CI for different size of the rolling window.

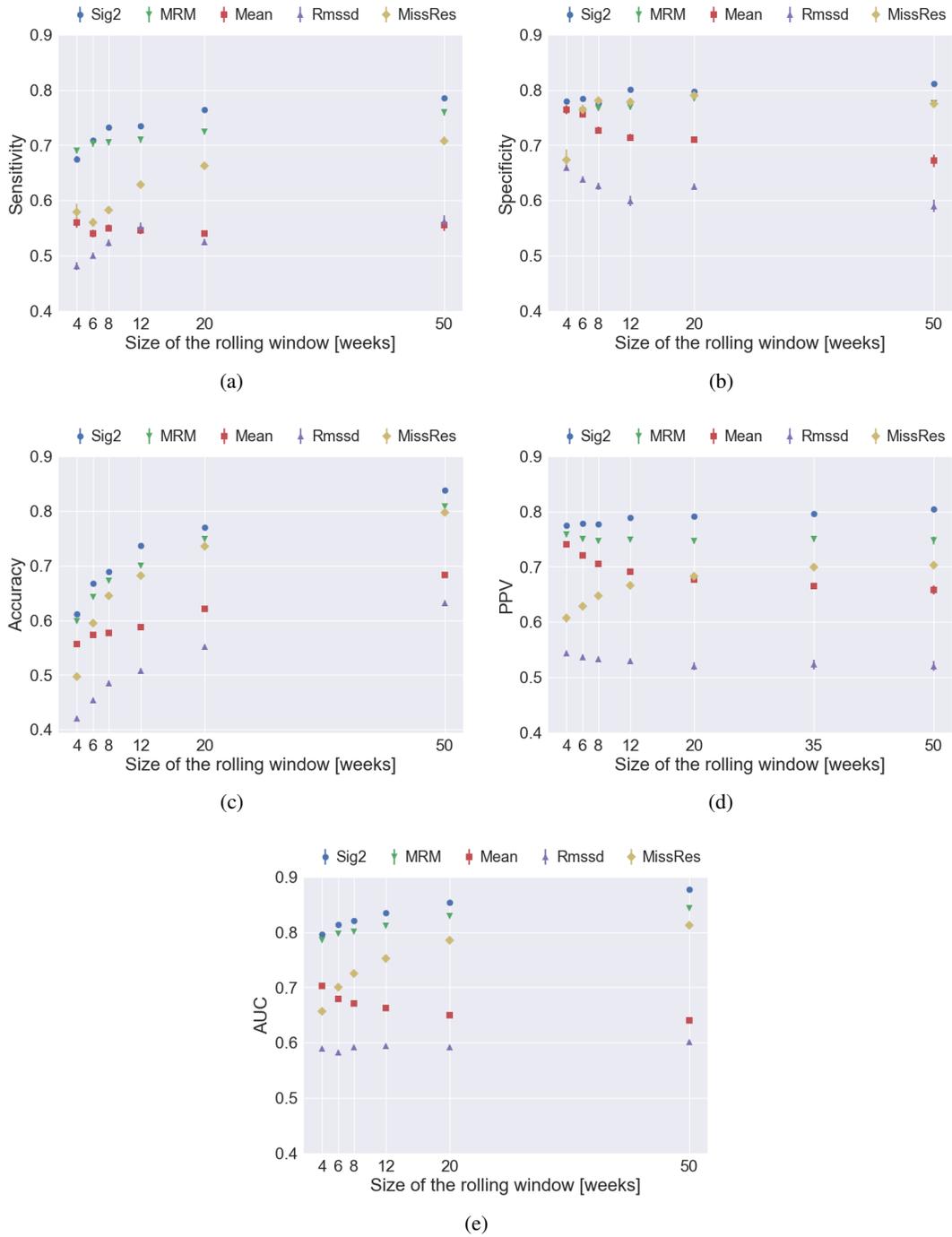


Figure 2. Classification metrics for mania model presented as mean with 95%CI for different size of the rolling window.

Diagnosis		BP-I (N=148)		BP-II (N=113)	
Sex		Male	Female	Male	Female
		51 (34%)	97 (66%)	39 (35%)	74 (65%)
Ethnicity					
	White	50 (98%)	86 (89%)	35 (92%)	66 (92%)
Age	years				
	17 - 25	4 (0.08%)	12 (0.13%)	7 (0.18%)	13 (0.18%)
	25 - 35	13 (0.25%)	25 (0.27%)	7 (0.18%)	23 (0.32%)
	35 - 45	16 (0.31%)	32 (0.34%)	9 (0.23%)	13 (0.18%)
	45 - 55	10 (0.20%)	14 (0.15%)	9 (0.23%)	12 (0.16%)
	55 - 65	6 (0.12%)	6 (0.06%)	7 (0.18%)	7 (0.10%)
	65 - 75	2 (0.04%)	5 (0.05%)	0 (0%)	5 (0.09%)
	Mean (SEM)	41.70 (1.77)	38.70 (1.36)	39.98 (2.23)	38.69 (1.77)
	Range	17 - 70	17 - 74	20 - 64	17 - 72.5
Median (IQR)	42 (16)	38 (18)	37 (22.5)	35 (23.5)	
Education					
	Tertiary	32 (0.63%)	52 (0.53%)	18 (46%)	44 (60%)
Adherence					
	Mean (SEM)	0.45 (0.04)	0.38 (0.05)	0.41 (0.05)	0.41 (0.04)
	Adjusted (Mean (SEM))	0.64 (0.02)	0.59 (0.03)	0.61 (0.02)	0.65 (0.04)
Episodes					
	Depression	23 (45%)	58 (60%)	17 (46%)	57 (77%)
	Mania	33 (65%)	68 (70%)	26 (67%)	54 (73%)

Table 1. Demographic summary of the data used for modelling.

weeks	Model	Sensitivity	Specificity	Accuracy	PPV	AUC
k=4	Sig	0.731 (0.054)	0.804 (0.037)	0.741 (0.032)	0.775 (0.019)	0.842 (0.021)
	MRM	0.703(0.062)	0.793(0.045)	0.723(0.032)	0.758(0.021)	0.818(0.023)
	Mean	0.692(0.071)	0.773(0.067)	0.702(0.048)	0.741(0.024)	0.793(0.033)
	Rmssd	0.422(0.068)	0.638(0.063)	0.464(0.061)	0.544(0.022)	0.546(0.024)
	MissRes	0.525(0.089)	0.669(0.106)	0.55(0.042)	0.607(0.038)	0.608(0.024)
k=6	Sig	0.763 (0.042)	0.788 (0.041)	0.755 (0.025)	0.779 (0.017)	0.847 (0.017)
	MRM	0.699(0.06)	0.787(0.046)	0.738(0.028)	0.75(0.02)	0.817(0.021)
	Mean	0.661(0.074)	0.768(0.059)	0.71(0.043)	0.721(0.024)	0.769(0.034)
	Rmssd	0.423(0.066)	0.633(0.059)	0.494(0.037)	0.537(0.024)	0.543(0.028)
	MissRes	0.495(0.105)	0.74(0.106)	0.624(0.038)	0.628(0.035)	0.65(0.026)
k=8	Sig	0.766 (0.042)	0.783(0.036)	0.76 (0.025)	0.777 (0.017)	0.846 (0.016)
	MRM	0.692(0.062)	0.791 (0.044)	0.75(0.026)	0.747(0.02)	0.814(0.019)
	Mean	0.638(0.078)	0.761(0.054)	0.708(0.039)	0.705(0.022)	0.751(0.033)
	Rmssd	0.443(0.072)	0.615(0.068)	0.508(0.036)	0.533(0.024)	0.548(0.027)
	MissRes	0.58(0.07)	0.706(0.051)	0.64(0.028)	0.647(0.023)	0.677(0.024)
k=12	Sig	0.772 (0.048)	0.798 (0.046)	0.799 (0.019)	0.789 (0.014)	0.854 (0.014)
	MRM	0.7(0.067)	0.79(0.049)	0.775(0.024)	0.749(0.018)	0.819(0.017)
	Mean	0.634(0.085)	0.739(0.078)	0.718(0.039)	0.691(0.022)	0.743(0.032)
	Rmssd	0.441(0.08)	0.618(0.08)	0.542(0.04)	0.529(0.029)	0.545(0.031)
	MissRes	0.58(0.063)	0.751(0.057)	0.707(0.035)	0.665(0.027)	0.717(0.025)
k=20	Sig	0.784 (0.053)	0.792 (0.049)	0.827 (0.018)	0.792 (0.018)	0.862 (0.014)
	MRM	0.705(0.085)	0.787(0.064)	0.809(0.023)	0.745(0.025)	0.825(0.018)
	Mean	0.612(0.085)	0.749(0.068)	0.754(0.038)	0.677(0.031)	0.732(0.037)
	Rmssd	0.46(0.084)	0.597(0.078)	0.586(0.038)	0.52(0.035)	0.546(0.038)
	MissRes	0.622(0.076)	0.755(0.059)	0.761(0.033)	0.683(0.027)	0.75(0.027)
k=50	Sig	0.799 (0.052)	0.807 (0.051)	0.873 (0.018)	0.805 (0.022)	0.877 (0.013)
	MRM	0.716(0.088)	0.789(0.063)	0.849(0.023)	0.748(0.037)	0.833(0.019)
	Mean	0.602(0.104)	0.738(0.084)	0.79(0.042)	0.658(0.041)	0.723(0.04)
	Rmssd	0.466(0.135)	0.614(0.132)	0.669(0.058)	0.52(0.045)	0.564(0.045)
	MissRes	0.652(0.081)	0.776(0.064)	0.829(0.029)	0.703(0.035)	0.782(0.03)

Table 2. Prediction results for depression on the test set of the signature-based and four baseline models for different size of the rolling window. The results are presented as mean(SD).

weeks	Model	Sensitivity	Specificity	Accuracy	PPV	AUC
k=4	Sig	0.675 (0.05)	0.78 (0.035)	0.611 (0.031)	0.747 (0.020)	0.796 (0.021)
	MRM	0.69(0.047)	0.763(0.037)	0.598(0.031)	0.741(0.021)	0.786(0.021)
	Mean	0.560(0.091)	0.764(0.085)	0.556(0.059)	0.697(0.038)	0.703(0.033)
	Rmssd	0.481(0.072)	0.659(0.06)	0.42(0.04)	0.598(0.024)	0.590(0.023)
	MissRes	0.579(0.149)	0.674(0.189)	0.497(0.065)	0.644(0.079)	0.657(0.026)
k=6	Sig	0.709 (0.049)	0.785 (0.043)	0.667 (0.026)	0.759 (0.019)	0.814 (0.018)
	MRM	0.703(0.054)	0.763(0.045)	0.643(0.027)	0.744(0.021)	0.797(0.021)
	Mean	0.54(0.066)	0.756(0.05)	0.573(0.044)	0.677(0.027)	0.679(0.035)
	Rmssd	0.5(0.058)	0.638(0.056)	0.454(0.037)	0.586(0.025)	0.583(0.029)
	MissRes	0.56(0.083)	0.764(0.094)	0.595(0.04)	0.69(0.039)	0.701(0.022)
k=8	Sig	0.732 (0.047)	0.776(0.045)	0.689 (0.024)	0.761 (0.017)	0.821 (0.014)
	MRM	0.705(0.053)	0.767(0.043)	0.672(0.025)	0.745(0.017)	0.801(0.016)
	Mean	0.55(0.069)	0.727(0.061)	0.577(0.042)	0.658(0.024)	0.671(0.03)
	Rmssd	0.523(0.065)	0.626(0.064)	0.485(0.042)	0.584(0.025)	0.592(0.029)
	MissRes	0.583(0.067)	0.781 (0.056)	0.645(0.04)	0.703(0.024)	0.725(0.023)
k=12	Sig	0.735 (0.043)	0.801 (0.045)	0.737 (0.018)	0.776 (0.016)	0.835 (0.014)
	MRM	0.71(0.054)	0.769(0.051)	0.7(0.023)	0.748(0.018)	0.811(0.017)
	Mean	0.546(0.07)	0.714(0.061)	0.588(0.042)	0.644(0.024)	0.663(0.033)
	Rmssd	0.552(0.08)	0.599(0.091)	0.508(0.041)	0.581(0.026)	0.594(0.028)
	MissRes	0.628(0.052)	0.778(0.051)	0.682(0.026)	0.717(0.023)	0.753(0.02)
k=20	Sig	0.764 (0.039)	0.797 (0.043)	0.77 (0.023)	0.784 (0.016)	0.854 (0.013)
	MRM	0.724(0.045)	0.785(0.042)	0.749(0.022)	0.759(0.016)	0.829(0.014)
	Mean	0.54(0.063)	0.710(0.052)	0.621(0.036)	0.632(0.024)	0.65(0.032)
	Rmssd	0.525(0.059)	0.625(0.059)	0.552(0.039)	0.58(0.026)	0.592(0.035)
	MissRes	0.663(0.05)	0.79(0.041)	0.736(0.026)	0.733(0.018)	0.786(0.019)
k=50	Sig	0.786 (0.046)	0.811 (0.04)	0.838 (0.022)	0.8 (0.019)	0.878 (0.013)
	MRM	0.76(0.057)	0.776(0.05)	0.808(0.026)	0.77(0.021)	0.843(0.016)
	Mean	0.555(0.1)	0.672(0.114)	0.683(0.052)	0.608(0.031)	0.64(0.036)
	Rmssd	0.563(0.101)	0.59(0.116)	0.632(0.043)	0.577(0.027)	0.602(0.035)
	MissRes	0.708(0.074)	0.775(0.065)	0.798(0.036)	0.741(0.025)	0.813(0.024)

Table 3. Prediction results for mania on the test set of the signature-based and four baseline models for different size of the rolling window. The results are presented as mean(SD).

weeks before episode	Depression	Mania
n = 0	0.696(0.048)	0.742(0.051)
n = 1	0.676(0.059)	0.675(0.074)
n = 2	0.606(0.063)	0.637(0.038)
n = 3	0.634(0.061)	0.612(0.063)
n = 4	0.614(0.078)	0.558(0.048)
n = 6	0.637(0.066)	0.571(0.071)

Table 4. The results of discrimination between the wellness and the precursor intervals of 6 weeks. Precursor intervals are tested at n week before an impending an episode. The values are the area under the ROC curve (AUC) and presented as mean(SD).