
DERIVING INFORMATION FROM MISSING DATA: IMPLICATIONS FOR MOOD PREDICTION

A PREPRINT

Yue Wu
Mathematical Institute
University of Oxford
Oxford, OX2 6GG, UK;
Alan Turing Institute
London, NW1 2DB, UK
yue.wu@maths.ox.ac.uk

Terry J. Lyons
Mathematical Institute
University of Oxford
Oxford, OX2 6GG, UK;
Alan Turing Institute
London, NW1 2DB, UK
terry.lyons@maths.ox.ac.uk

Kate E.A. Saunders
Department of Psychiatry
University of Oxford;
Oxford Health NHS Foundation Trust
Warneford Hospital
Oxford OX3 7JX, UK
kate.saunders@psych.ox.ac.uk

July 9, 2020

ABSTRACT

The availability of mobile technologies has enabled the efficient collection prospective longitudinal, ecologically valid self-reported mood data from psychiatric patients. These data streams have potential for improving the efficiency and accuracy of psychiatric diagnosis as well predicting future mood states enabling earlier intervention. However, missing responses are common in such datasets and there is little consensus as to how this should be dealt with in practice. A signature-based method was used to capture different elements of self-reported mood alongside missing data to both classify diagnostic group and predict future mood in patients with bipolar disorder, borderline personality disorder and healthy controls. The missing-response-incorporated signature-based method achieves roughly 66% correct diagnosis, with f1 scores for three different clinic groups 59% (bipolar disorder), 75% (healthy control) and 61% (borderline personality disorder) respectively. This was significantly more efficient than the naive model which excluded missing data. Accuracies of predicting subsequent mood states and scores were also improved by inclusion of missing responses. The signature method provided an effective approach to the analysis of prospectively collected mood data where missing data was common and should be considered as an approach in other similar datasets.¹

Keywords Signature method · Missing responses · Bipolar disorder · Borderline personality disorder

1 Introduction

The rapid emergence of mobile technologies has transformed the way in which mental health data can be collected. Until recently clinicians were wholly reliant on anamnestic approaches and thus hampered by the inaccuracy of retrospective recall regarding psychiatric symptoms. Mobile technologies have enabled the efficient capture of self-reported symptoms in an ecologically valid and prospective manner. A number of different approaches to the analysis of longitudinal mood data have been employed [7, 4, 3]. However missing data is ubiquitous and poses a significant methodological challenge. Mood data may be missing unrelated to mood state or in fact be a consequence of current mood state. Standard approaches such as imputation may inadvertently lead to the loss of important information [8].

Signatures from rough path theory [15] are an effective way of analyzing these types of data as they capture the order in which events occur. The approach incorporates missing data into the system in order to capture the underlying patterns and the evolving interactions between missing data and responses, and to predict the future from the past of an evolving system [13]. So far, the signature method has significantly contributed to automated recognition of Chinese

¹MSC2020: 60L10, 60L90, 62D10, 62P10, 92-08

handwriting [9, 24], formulation of appropriate stochastic partial differential equations to model randomly evolving interfaces [11, 12], skeleton-based human action recognition [25], diagnosis of Alzheimer’s disease [17] and speech emotion recognition [23]. In a previous analysis we demonstrated that a signature-feature model could be successfully applied to 6-dimensional self-reported mood data [3], however missing data was not used in this analysis.

In this study, we used a missing-response-incorporated signature-feature-based machine learning model to re-analyse weekly mood data collected from the AMoSS study [22] which used self-reported mood data and wearables to distinguish between individuals with bipolar disorder (BD), borderline personality disorder (BPD) and healthy controls (HC). We sought to test whether this new analytic approach was superior to standard approaches to mood quantification in its ability to distinguish these diagnostic groups and predict future mood states/scores.

2 Methods

2.1 Data

Patients with BD or BPD and healthy volunteers reported their mood using Altman Self-Rating Mania scale (ASRM) [1] and the Quick Inventory of Depressive Symptoms (QIDS-SR16) [20] were collected. ASRM is a short, 5-item self-assessment questionnaire assessing the presence and severity of manic or hypomanic symptoms. QIDS-SR16 contains 16 items covering the 9 DSM-IV symptom criterion domains [2] with the total score ranging from 0 to 27.

The data were collected as part of the AMoSS study [22]. ASRM and QIDS data were collected from 142 individuals as part of the AMoSS study and the participants completed standardised questionnaires on a weekly basis using the True Colors mood monitoring system [10] after receiving a text or email prompt. Two of the 142 participants either withdrew consent or had no clinical diagnosis and were therefore excluded from analysis. We further excluded 14 participants who failed to provide at least 20 weeks data as part of the analysis is based on information in data of at least 20 weeks. Of the remaining 126 participants, 49 were diagnosed as bipolar disorder and 32 were borderline personality disorder. All identical duplicate values were checked and removed, and only the first response of a week was kept if multiple responds happened within that week. Most of the time, assessments were completed one after the other, of which the sum scores can then be paired up. We discarded the excessive observed data which could not be paired. Using this data, 2-dimensional concatenated data of paired observations from ASRM and QIDS was obtained for each participant, where the score is ‘-1’ for a missing response.

Table 1: Demographic characteristics of the three clinical groups (the appropriate distributions are summarised in the form of the median \pm in the interquartile range).

Group	Recruited	for analysis	Weeks in study	Ages	Gender (males)
BD	54	49	51 \pm 3	38 \pm 19	16
HC	52	45	51 \pm 2	37 \pm 18	13
BPD	34	32	51 \pm 2	34 \pm 12	3

2.2 Features extraction

2.2.1 Encoding missing data

Among all the 126 valid participants in our study, 90% missed a response on a least one occasion during their task-active weeks. Signature features allow missing responses to be included in the analysis without the need for imputation. The missing events are translated into a new counting process [14]. An example is illustrated in Figure 1 for the procedure. In the general case, if one work on data with N many time points, the accumulative missing counts can be generated for each of the N time points by calculating the sum of missing observations up to that particular time point; meanwhile each missing observation, i.e., input “-1” in our case, is replaced by the valid value that happened in the nearest past (referred as the *feed forward* method). This does not imply that the missing responses are assumed to take the same value as their nearest valid responses. By doing this, the increments in both observation and missing counts can be preserved and captured, which are indeed the most critical features in the signature method together with their functionals [16].

After transforming missing responses, one then normalises (and accumulate like in [3]) the data to make it scale-free in order to apply signature transformation.

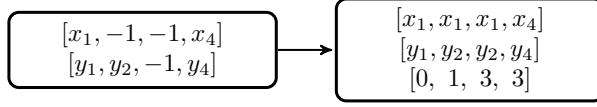


Figure 1: An example for integrating missing responses: the left block contains 2 dimensional data of 4 consecutive observations, where -1 represents one missing observation; in the right block the feed forward method is applied for filling in all missing places with valid values that happened in the corresponding nearest past, while an additional dimension is added to count missing events cumulatively at each time points.

2.2.2 Signature features

In recent year, the signatures of the continuous paths generated from longitudinal data is considered as an efficient feature set for learning purpose because of its nature to capture the nonlinear effect of the evolving systems [16]. Consider \mathbb{R}^d -valued time-dependent, piecewise-differentiable paths of finite length. Such a path X mapping from time domain $[a, b]$ to \mathbb{R}^d is denoted as $X : [a, b] \rightarrow \mathbb{R}^d$. For short we will use X_t for $X(t), t \in [a, b]$. Each coordinate path of X is a real-valued path and denoted as $X^i, i \in [d]$ with $[d] := \{1, \dots, d\}$. The *signature* of a path $X : [a, b] \rightarrow \mathbb{R}^d$, denoted by $S(X)_{a,b}$, is the infinite collection of all iterated integrals of X . That is,

$$S(X)_{a,b} := (1, S(X)_{a,b}^1, \dots, S(X)_{a,b}^d, S(X)_{a,b}^{1,1}, S(X)_{a,b}^{1,2}, \dots), \quad (1)$$

where, the 0th term is 1 by convention, and the superscripts of the terms after the 0th term run along the set of all multi-index $\{(i_1, \dots, i_k) | k \geq 1, i_1, \dots, i_k \in [d]\}$ with the coordinate iterated integral being

$$S(X)_{a,b}^{i_1, \dots, i_k} := \int_{a < t_k < b} \dots \int_{a < t_1 < t_2} dX_{t_1}^{i_1} \dots dX_{t_k}^{i_k}. \quad (2)$$

The finite collection of all terms $S(X)_{a,b}^{i_1, \dots, i_k}$ with the multi-index of fixed length k is termed as the k th level of the *signature*. The truncated signature up to the p th level is denoted by $\lfloor S(X)_{a,b} \rfloor_p$. In machine learning context, truncated signature features are always obtained by truncating the original signature to some finite level.

For discrete data stream $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, where \mathbf{x} contains n observations, and the i th observation $\mathbf{x}_i, i \in [n]$, is assumed to be a d -dimensional column vector at the i th time point, one needs to convert it to a continuous \mathbb{R}^d -valued path via piecewise linear interpolation or other transforms in order to compute signature. The availability of Python packages *iisignature* [19] and *esig* allows easy calculation of signature, where the linear interpolation is implemented automatically by the packages.

For our purpose, we extracted the consecutive paired observations for each participant, incorporated the missing data, and then calculated the corresponding signature features via Python package *iisignature*, where the signature features were truncated to level 2. To distinguish from standard signature features, our features were named the missing-response-incorporated signature features (MRSF).

2.3 Signature-based classification

In order to investigate the role of ASRM and QIDS scores in differentiating between healthy controls and different patient groups, a missing-response-incorporated signature-based classification model was developed to classify the diagnostic group a participant belonged to. For each of 126 participants, a stream of 20 consecutive paired observations, no matter missing or not, was randomly drawn from the 2-dimensional concatenated data for this task. Then the collection of the 20 consecutive paired data were randomly split into a training set (70%) and a testing set (30%). The proposed model was based on a random forest classifier and was trained on the MRSF of the training set.

For comparison, the random forest classifier was also trained on features extracted through a clinic-used metric based on the average score in each category over the valid scores in 20 consecutive observations (the naive classification model). The performance of the missing-response-incorporated signature-based classification model (MRSCM, level 2) and the baseline model (the naive classification model), for classifying the diagnostic groups were tested on bootstrapping and measured in terms of accuracy. Meanwhile the confusion matrices of both methods were generated to allow more detailed analysis, from which f1 scores for different clinic groups were computed. To assess the separation ability of different methods, we created the receiver operating characteristic curves (ROC) at various threshold settings and computed areas under curve (auc).

2.3.1 Spectrum analysis

To further test the performance of the MRSCM (level 2), we investigated the likelihood of each of the three groups being categorised into the correct clinic group. We trained the model on the 20 consecutive streamed data of all but one participant and tested it with the data of this participant. The probability vector of each participant being classified into each group was calculated and then projected onto the equilateral triangle, with each vertex representing one of the three clinic groups. For example, if the inferred probabilities of one participant being classified as BD, HC and BPD are 0.1, 0.5 and 0.4 respectively, then the corresponding probability vector is $[0.1, 0.5, 0.4]$. This vector is indeed on a 3-dimensional triangle surface $[p, q, 1 - p - q]$, with non-negative p, q and $p + q \leq 1$. This triangle is the equilateral triangle that all the inferred 3-dimensional probability vectors will be sitting on.

2.4 Signature-based predictions

We then sought to predict the mania and depression states/scores of each participant in the next week on the last 10 paired observations. Empirically, using the last 10 consecutive data was more effective than using 20, perhaps the future mood state/score is dependent on the most recent states.

2.4.1 State prediction

At this stage, in order to simplify the prediction problem, we reduced the outcome to one of three mutually exclusive responses. We treated ASRM and QIDS similarly but separately. That is, for ASRM, no answer (future response is missing), normal (score of ASRM is no bigger than 5) or manic (score of ASRM is bigger than 5); and for QIDS, no answer (future response is missing), normal (score of QIDS is no bigger than 10) or depressed (score of QIDS is bigger than 10).

The target of the prediction is not quantified, so we approached the problem as a classification problem rather than as a regression. We chose to learn the relation between the pattern of last 10 consecutive observations and the future state through classifiers. For this task, the MRSF truncated at level 2 were extracted from data and a random forest classifier-based predictive model was trained for each clinical group separately. For comparison, the random forest classifier-based predictive model was also trained on features extracted through computing the mean of the last 10 consecutive observations (naive predictive model). The prediction of ASRM/QIDS states of each participant based on diagnostic group using the missing-response-incorporated signature-based predictive model (stateMRSPM, level 2) and the naive predictive model were tested with bootstrapping and measured in terms of accuracy.

2.5 Spectrum analysis

To assist the understanding towards the distribution of the responses to ASRM/QIDS, the bar histograms of the proportions of selected features, namely, no answer/normal/manic for ASRM and no answer/normal/depressed for QIDS, for each clinic group were plotted for an overview.

We also visualised the true proportion vectors as well as the prediction results of each participant using triangle spectrum plots. This is a natural adoption as we also have three different states: no answer, normal and manic/depressed.

The true proportion vector reflects the ground truth. It consists of three elements, that is, the proportion of this participant giving 'no answer', 'normal' and 'manic' (resp. 'depressed') for ASRM (resp. QIDS) during his/her entire study. In order to demonstrate group-dependent characteristics, true proportion vectors for the same questionnaire of patients from the same clinical group were visualised in the same 3-dimensional equilateral triangle surface.

Regarding the predicted result, for each participant of one of the three groups (BD/HC/BPD), ASRM/QIDS states of 5 weeks rather than 1 week were predicted using stateMRSPM (level 2) trained from random 10 consecutive observations of the rest participants in the same group. To avoid overfitting, we excluded the participants that generated at most 5 buckets of 10 consecutive weeks data. Then the predicted proportion of his/her giving each of the three different responses for future ASRM/QIDS could be directly calculated from the predicted ASRM/QIDS states of the five weeks.

2.6 Score prediction

We also sought to predict the next reported ASRM/QIDS score made by a participant based on features extracted from their previous 10 weekly observations, no matter missing or not. To ensure comparability of our results, we predicted raw ASRM/QIDS scores. For our own comparison, we also made more coarse-grained predictions according to the categories for QIDS: none (1-5), mild (6-10), moderate (11-15), severe (16-20) and very severe (21-27) [20], and for ASRM: none (0-5), mild (6-9), moderate (10-13), severe (14-17) and very severe (17-20).

For this task, the MRSF truncated at level 2 were extracted from each sequential data and a random forest regressor-based predictive model was trained for each clinical group separately. The future score prediction of each participant using missing-response-incorporated signature-based predictive model (scoreMRSPM, level 2) was tested with bootstrapping and measured in terms of accuracy and mean absolute error (MAE).

Summary

We used the publicly available Python iisignature package (version 0.23) to calculate signatures of streams of data, Python numpy package (version 1.19.0) for data manipulations and processing, Python scikit-learn package (version 0.23.1) for machine learning tasks and matplotlib for plotting and graphics (version 3.2.1).

The study was approved by the NRES Committee East of England—Norfolk (13/EE/0288) and the Research and Development department of Oxford Health NHS Foundation Trust.

A summary of tasks and models can be found in Table A1 and Table A2.

3 Results

3.1 Classification of the diagnostic group

MRSCM (level 2) categorised 66.3% of participants into the correct with a standard deviation 0.06 while the naive model only classified 54.9% of participants correctly with a higher variability 0.08. The accuracy from MRSCM improved with transformation of missing responses, indicating that missing responses bring additional information and therefore enhance the performance of the model.

Figure 2 gives the confusion matrices for both methods, which illustrates the detailed correct and false classification for each group and allows for computing f1 scores in Table 2. Table 2 shows that the MRSCM (level 2) has higher f1 scores in all three classes. It achieved its lowest f1 score for classifying bipolar group. However, encoding the missing information into the model, the ability of classifying BPD was significantly enhanced from 40%+ (the naive model) to 60%+ (MRSCM, level 2).

The receiver operating characteristic curves for three clinic groups from both methods are plotted in Figure 3 with corresponding areas under curve (auc) recorded in the brackets. The MRSCM (level 2) has better ability in identifying all three diagnostic groups in terms of auc. Both methods have their lowest auc from ROC of bipolar group, which implies it is more likely for bipolar participants to be misplaced into the other two groups.

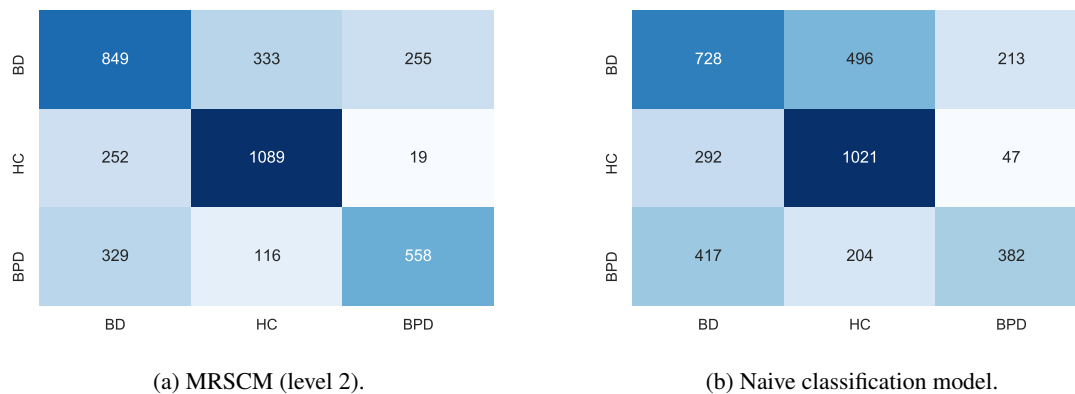


Figure 2: Confusion matrices of the missing-response-incorporated signature-based classification model (MRSCM, level 2) and the naive classification model.

3.1.1 Spectrum analysis

In Figure 4, the triangle spectrum of the predicted diagnosis from MRSCM (level 2) are plotted. In each of the plots, the regions of highest density of participants are located in the correct corner of the triangle. The greatest

Table 2: f1 scores for group classification using the missing-response-incorporated signature-based classification model (MRSCM, level 2) and the naive classification model.

Model	BD	HC	BPD
MRSCM (level 2)	59.2%	75.2%	60.8%
Naive classification model	50.7%	66.3%	46.4%

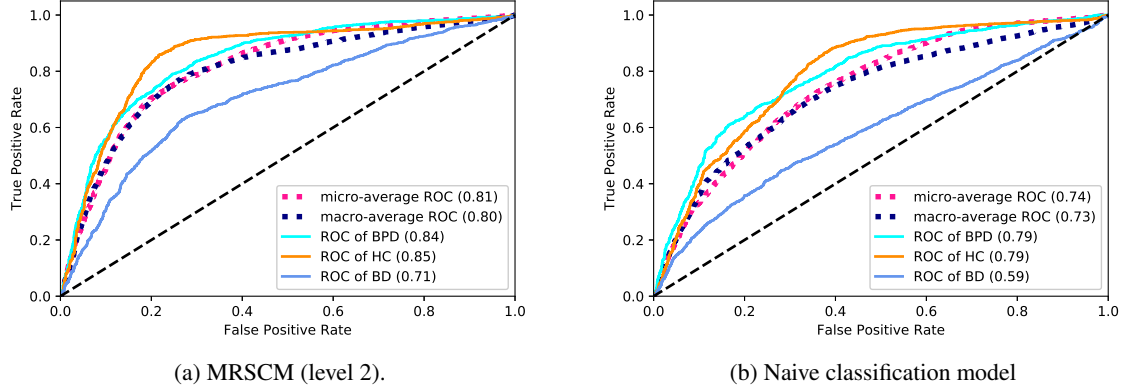


Figure 3: Receiver operating characteristic curves of the missing-response-incorporated signature-based classification model (MRSCM, level 2) and the naive classification model.

consistency is with the healthy participants. Meanwhile, the probabilities of misdiagnosis to other groups can be measured by comparing the distances to the other two vertices to the distance to the right vertex. For instance, one can deduce from Plot (b) that the likelihood of misplacing healthy participants into the borderline group is very low. Plot (c) shows the other way around: BPD participants are unlikely to be misidentified as healthy control. Plot (a) shows that the bipolar patient can be misidentified as healthy control with relatively high probability and as BPD with relatively low probability.

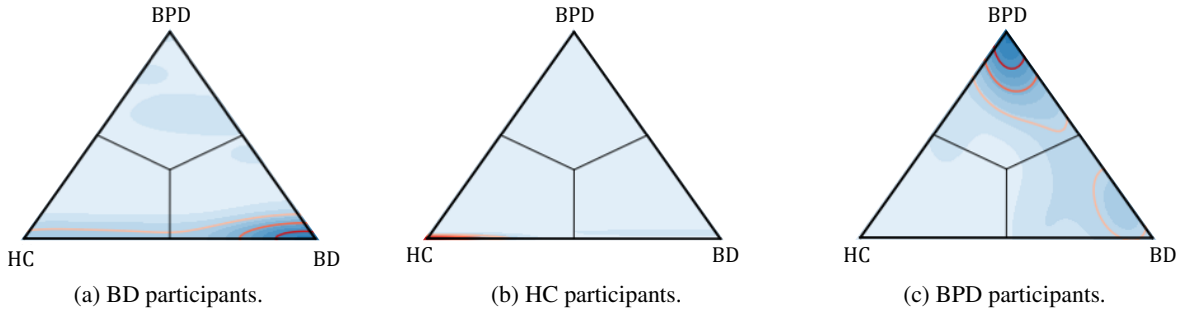


Figure 4: Density plots for the predicted diagnosis from MRSCM (level 2): darker blue areas indicate higher density values, i.e., events that are more likely to happen, and vice versa; red lines indicate the 75% (the lightest red), 50%, 25% (the darkest red) boundaries of density contours, i.e., the events within the area enclosed by the 75% contour line is with probability 75% to happen.

3.2 Prediction of the ASRM/QIDS scores of an individual participant

Prediction of the ASRM/QIDS scores of each participant based on his/her own clinic group, are tested on bootstrapping and summarised in Table 3. Compared to the naive method, the model based on signature features (truncated at level 2) extracted after handling missing data is with the higher accuracy across all different clinic groups.

Among the three clinic groups, Both models have their best performance in predicting the ASRM/QIDS scores for healthy controls and worst performance in predicting the ASRM/QIDS scores for borderline personality disorder.

Table 3: The average accuracy for ASRM/QIDS state prediction using the missing-response-incorporated signature-based predictive model (stateMRSPM, level 2) and the naive predictive model.

Model	BD		HC		BPD	
	ASRM	QIDS	ASRM	QIDS	ASRM	QIDS
stateMRSPM (level 2)	70.6%	64.8%	79.9%	78.9%	65.2%	60.2%
Naive predictive model	61.7%	59.5%	70.1%	71.6%	58.6%	55.6%

Figure 5 demonstrates the proportions of ASRM/QIDS scores of each of the three clinic groups. HC participants rarely experienced either mania or depression in their entire study and were least likely to have missing responses. By contrast, BPD patients were most likely to have missing responses. Figure ?? visualises the densities of the ASRM/QIDS states for three clinic groups using the 3-dimensional equilateral triangle surface.

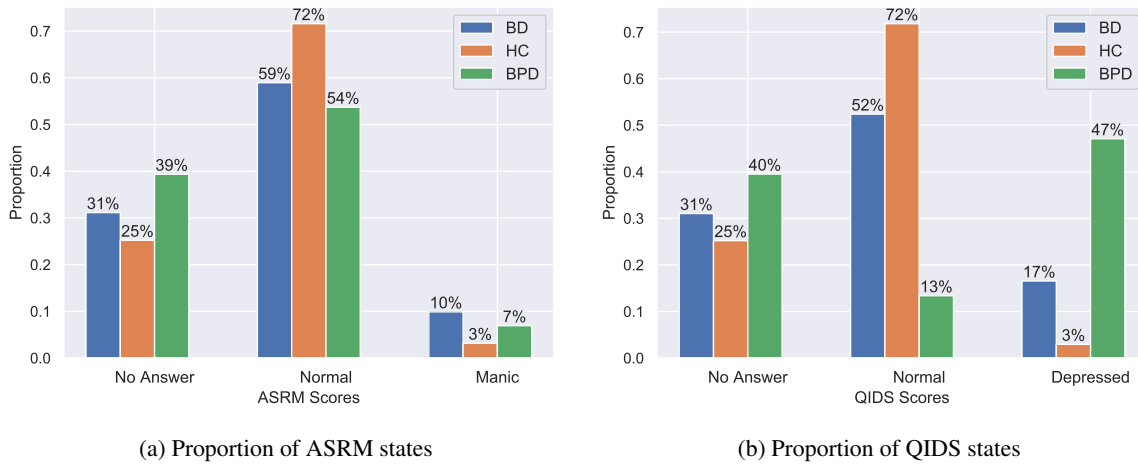


Figure 5: Bar charts: the proportion of ASRM/QIDS states of the three clinic groups, where the total number of ASRM or QIDS questionnaires for BD/HC/BPD are 620, 880, 960.

Figure 7 demonstrates the density spectrum of the predicted ASRM/QIDS states for each clinic group using the missing-response-incorporated signature-based predictive model (level 2), which is roughly consistent with Figure 6. Recall that for ASRM (resp. QIDS), only 5 future observations were predicted for each participant and therefore used for computing the proportion (i.e., density). The difference between Figure 7 and Figure 6 comes from the mixed effects of the prediction error and limited number of observations being predicted for each participant.

3.3 Score predictions

MAEs are recorded in Table 4 for predicting the next reported score using scoreMRSPM (level 2) and the naive model. The worst performances were to predict ASRM scores for bipolar patients and to predict future QIDS scores for patients of borderline personality disorder. Table 5 reports the accuracy and MAE for predicting the future severity of symptoms using signature-based model.

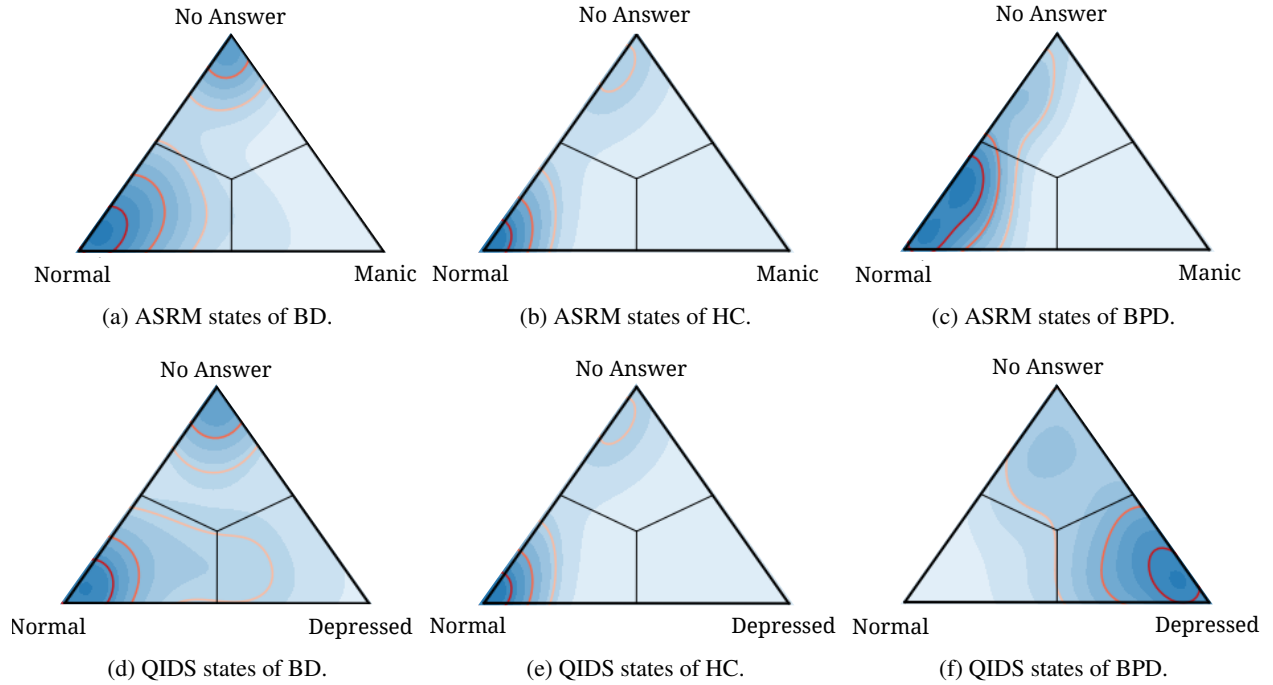


Figure 6: Density plots of true ASRM/QIDS states for three clinic groups: darker blue indicates higher density value and vice versa; red lines indicate the 75% (the lightest red), 50%, 25% (the darkest red) boundaries of density contours.

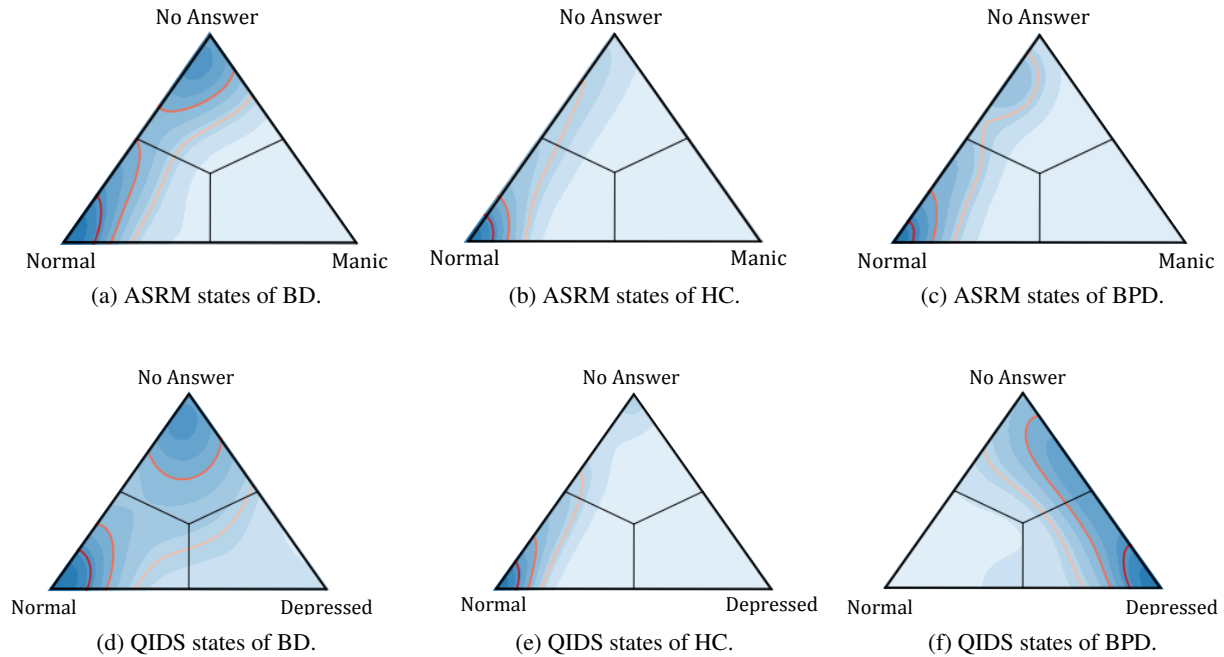


Figure 7: Density plots of the predicted ASRM/QIDS states for three clinic groups using the missing-response-incorporated signature-based predictive method (level 2): darker blue indicates higher density value and vice versa; red lines indicate the 75% (the lightest red), 50%, 25% (the darkest red) boundaries of density contours

Table 4: MAE for ASRM/QIDS score prediction using the missing-response-incorporated signature-based predictive method (scoreMRSPM, level 2) and the naive predictive model.

Model	BD		HC		BPD	
	ASRM	QIDS	ASRM	QIDS	ASRM	QIDS
scoreMRSPM (level 2)	2.386	3.437	0.826	1.532	2.117	3.745
Naive predictive model	3.286	4.600	1.137	1.899	2.571	4.671

Table 5: The average accuracy and MAE for the prediction of severity of symptoms using the missing-response-incorporated signature-based predictive method.

Measure	BD		HC		BPD	
	ASRM	QIDS	ASRM	QIDS	ASRM	QIDS
Accuracy	74.3%	76.4%	95.8%	95.0%	82.4%	69.8%
MAE	1.046	0.685	0.191	0.139	0.625	0.794

4 Discussion

This paper introduces the missing-response-incorporated signature-feature-based random forest models and have them tested on the paired ASRM/QIDS data. Comprising at least 25% of the whole dataset (Figure 5), the missing response is a remarkable and unreplaceable component and cannot be simply ignored. By integrating the informative and non-redundant resource from the missing responses, the empowered signature features can be utilised with different machine learning models on various datasets containing missing information, either to identify the category membership based on the observed characteristics of the data or to predict the future from the past of an evolving system generated by the data.

The missing-response-incorporated signature-based classification model is superior to the commonly used metric (the naive classification model) in differentiating the three clinic groups, i.e., bipolar disorder, healthy control and borderline personality disorder. Although signature-based methods were adopted by Perez et al [3] to improve the existing results that used neuroimaging [21] or verbal fluency [6], our analysis provides a benchmark at the same accuracy level while utilising data of much smaller dimensions and missing responses.

Spectrum analysis showed the overlap between the BD and HC groups in Figure 4, which is consistent with the analysis in [3] and with clinical experience. However, we found a much clearer differentiation between clinical groups than previous work [3] suggesting that the inclusion of missing data added useful information.

The prediction of future ASRM/QIDS states is notable. To our knowledge this is the first time that non-response status has been considered alongside normal status and manic/depressed status. Borderline personality disorder tends to miss the tasks with slightly higher probability for both self-reported assessments than the other two groups, suggesting that some group-dependent feature is concealed in missing responses.

Most of previous literature like [3] and [18] were interested in predicting future scores. They did not use the missing responses as part of their information. In our analysis, the inclusion of the missing response information significantly improved our predictions. The performance (first row in Table 4) in general outperforms the one in [18] in terms of mean absolute error. Unlike models that need to impose population-level distribution and parameters [5], our proposed model offers a unified, missing-response-incorporated and non-parametric approach that is able to capture the nature of the evolving system which generates data streams.

For the naive models, both the performance of ASRM/QIDS-state prediction (Table 3) and the one of diagnostic group classification (Table 2) for BPD patients are the worst among the three clinic groups. The f1 score for classification was even below 0.5. For this task, the confusion matrix (Plot (b) in Figure 2) also shows that about 40% of the BPDs were misclassified as BD patients. The poor performance alerted the unreliability of this naive metric in identifying BPD patients and predicting their moods. On the other hand, by incorporating extra valuable information like missing responses into features, the signature-based model lifted the f1 score for identifying BPD patients to above 60%, with less than one third BPD patients being misclassified as BD. This demonstrates the ability of the missing-response-incorporated signature features to capture and learn the inherent difference in mood instability between BPD and BD. Even though, the performance of ASRM/QIDS-state prediction for BPD patients from signature-based model (the last

two column in Table 3) is still the worst among the three clinic groups. The poor performance is again consistent with the results from [3]. This consistency may be a result of the unpredictable nature in BPD, or due to the bias from the same database that is used by both [3] and us.

4.1 Limitations and implications

The missing-response-incorporated signature-based features offer a systematic approach to the analysis of longitudinal self-reported mood data with the presence of non-randomly distributed missing values. It can be easily utilised with various machine learning methods for classification and prediction tasks on other databases containing missing information. The reasons for the moderate accuracies using MRSF are four-fold: the full potential of signature features is hindered by the small dataset, the proposed feature extraction method might not be the optimal, ASRM/QIDS data was analysed on the overall-score level instead of on the question-score level, and the diversity within the same clinic group was not considered for the prediction task. In the future, we would prioritise on two explorations: assessing our proposed method on different mental health datasets, and adjusting MRSF to the “optimal” signature-based feature by adding reasonable metrics/transformations which account for different attributes. We are also interested in addressing the inter-group difference on a much larger dataset of ASRM and QIDS responses by building independent or dependent missing-response-incorporated signature-based predictive models on clusters within the same group.

Data, code and materials

As the data were collected pre-GDPR and contained sensitive personal data, it cannot be placed into a publicly accessible repository.

The codes can be found through GitHub repository via https://github.com/yuewu57/mental_health_AMoSS.

Acknowledgements

YW and TJL would acknowledge Alan Turing Institute for funding this work through EPSRC grant EP/N510129/1 and EPSRC through the project EP/S2026347/1, titled ‘Unparameterised multi-modal data, high order signature, and the mathematics of data science’. KEAS is supported by the NIHR Oxford Health Biomedical Research Centre. *‘The views expressed are those of the authors and not necessarily those of the NHS, the NIHR or the Department of Health’.*

Appendix

Table A1: A summary of tasks, where conPrediction is short for conditional prediction.

Task	Description	Classes or Values
Classification	To identify the participant’s diagnostic group	BD/HC/ BPD
State prediction	To predict the participant’s future mood state	ASRM: No response/Normal/Manic QIDS: No response/Normal/Depressed
Score prediction	To predict the participant’s future raw score given that it is not missing	ASRM: 0-20 QIDS: 0-27
Severity prediction	The future severity of symptoms obtained by cut-off scores of score prediction	ASRM: 0-4 QIDS: 0-4

References

- [1] Altman, E.G., Hedeker, D., Peterson, J.L. and Davis, J.M., The Altman self-rating mania scale. *Biological psychiatry*, 42.10 (1997): 948-955.

Table A2: A summary of models, where MR is short for missing responses, RF short for random forest and conPrediction short for conditional prediction.

Task	Model	Data length	Feature extraction		Base model
			MR integration	Signatures	
Classification	MRSCM (level 2)	20	Yes	Yes	RF classifier
	Naive classification model	20	No	No	RF classifier
State prediction	stateMRSPM (level 2)	10	Yes	Yes	RF classifier
	Naive predictive model	10	No	No	RF classifier
Score prediction	scoreMRSPM (level 2)	10	Yes	Yes	RF regressor
	Naive predictive model	10	No	No	RF regressor

- [2] American Psychiatric Association, Diagnostic and statistical manual of mental disorders. *BMC Medicine*, 17 (2013): 133-137.
- [3] Arribas, I.P., Goodwin, G.M., Geddes, J.R., Lyons, T. and Saunders, K.E., A signature-based machine learning model for distinguishing bipolar disorder and borderline personality disorder. *Translational psychiatry*, 8.1 (2018): 274.
- [4] Bopp, J.M., Miklowitz, D.J., Goodwin, G.M., Stevens, W., Rendell, J.M. and Geddes, J.R., The longitudinal course of bipolar disorder as revealed through weekly text messaging: a feasibility study. *Bipolar disorders*, 12.3 (2010): 327-334.
- [5] Busk, J., Faurholt-Jepsen, M., Frost, M., Bardram, J.E., Kessing, L.V. and Winther, O., Forecasting Mood in Bipolar Disorder From Smartphone Self-assessments: Hierarchical Bayesian Approach. *JMIR mHealth and uHealth*, 8.4 (2020): e15028.
- [6] Costafreda, S.G., Fu, C.H., Picchioni, M., Touloupoulou, T., McDonald, C., Kravariti, E., Walshe, M., Prata, D., Murray, R.M. and McGuire, P.K., Pattern of neural responses to verbal fluency shows diagnostic specificity for schizophrenia and bipolar disorder. *BMC psychiatry*, 11.1 (2011): 18.
- [7] Faurholt-Jepsen, M., Frost, M., Ritz, C., Christensen, E.M., Jacoby, A.S., Mikkelsen, R.L., Knorr, U., Bardram, J.E., Vinberg, M. and Kessing, L.V., Daily electronic self-monitoring in bipolar disorder using smartphones—the MONARCA I trial: a randomized, placebo-controlled, single-blind, parallel group trial. *Psychological medicine*, 45.13 (2015): 2691-2704.
- [8] Faurholt-Jepsen, M., Geddes, J.R., Goodwin, G.M., Bauer, M., Duffy, A., Kessing, L.V. and Saunders, K., Reporting guidelines on remotely collected electronic mood data in mood disorder (eMOOD)—recommendations. *Translational psychiatry*, 45.13 (2019): 1-10.
- [9] Graham, B., Sparse arrays of signatures for online character recognition. *arXiv preprint*, arXiv:1308.0371.
- [10] Goodday, S.M., Atkinson, L., Goodwin, G., Saunders, K., South, M., Mackay, C., Denis, M., Hinds, C., Attenburrow, M.J., Davies, J. and Welch, J., The true colours remote symptom monitoring system: a decade of evolution. *Journal of Medical Internet Research*, 22.1 (2020): e15188.
- [11] Hairer, M., Solving the KPZ equation. *Annals of Mathematics*, (2013): 559-664.
- [12] Hairer, M., A theory of regularity structures. *Inventiones mathematicae*, 198.2 (2014): 269-504.
- [13] Levin, D., Lyons, T. and Ni, H., Learning from the past, predicting the statistics for the future, learning an evolving system. *arXiv preprint*, arXiv:1309.0260.
- [14] Little, R.J. and Rubin, D.B., *Statistical analysis with missing data*, John Wiley & Sons, 793 (2019).
- [15] Lyons, T. and Qian, Z., *System control and rough paths*, Oxford University Press (2002).
- [16] Lyons, T., Rough paths, signatures and the modelling of functions on streams. *arXiv preprint*, arXiv:1405.4537. 2019. Using path signatures to predict a diagnosis of Alzheimer’s disease. *PloS one*, 14(9).
- [17] Moore, P.J., Lyons, T.J., Gallacher, J. and Alzheimer’s Disease Neuroimaging Initiative, Using path signatures to predict a diagnosis of Alzheimer’s disease. *PloS one*, 14.9 (2019).
- [18] Palmius, N., Tsanas, A., Saunders, K.E.A., Bilderbeck, A.C., Geddes, J.R., Goodwin, G.M. and De Vos, M., Detecting bipolar depression from geographic location data. *IEEE Transactions on Biomedical Engineering*, 64.8 (2016):1761-1771.

- [19] Reizenstein, J. and Graham, B., The iisignature library: efficient calculation of iterated-integral signatures and log signatures. *arXiv preprint*, arXiv:1802.08252.
- [20] Rush, A.J., Trivedi, M.H., Ibrahim, H.M., Carmody, T.J., Arnow, B., Klein, D.N., Markowitz, J.C., Ninan, P.T., Kornstein, S., Manber, R. and Thase, M.E., 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. *JBiological psychiatry*, 54.5 (2003): 573-583.
- [21] Sato, J.R., de Araujo Filho, G.M., de Araujo, T.B., Bressan, R.A., de Oliveira, P.P. and Jackowski, A.P., Can neuroimaging be used as a support to diagnosis of borderline personality disorder? An approach based on computational neuroanatomy and machine learning. *Journal of psychiatric research*, 46.9 (2012): 1126-1132.
- [22] Tsanas, A., Saunders, K.E.A., Bilderbeck, A.C., Palmius, N., Osipov, M., Clifford, G.D., Goodwin, G.M. and De Vos, M., Daily longitudinal self-monitoring of mood variability in bipolar disorder and borderline personality disorder. *Journal of affective disorders*, 205 (2016): 225-233.
- [23] Wang, B., Liakata, M., Ni, H., Lyons, T., Nevado-Holgado, A.J. and Saunders, K., A Path Signature Approach for Speech Emotion Recognition. *Interspeech 2019*, ISCA (2019): 1661-1665.
- [24] Xie, Z., Sun, Z., Jin, L., Ni, H., and Lyons, T., Learning spatial-semantic context with fully convolutional recurrent network for online handwritten Chinese text recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40.8 (2017): 1903-1917.
- [25] Yang, W., Lyons, T., Ni, H., Schmid, C., Jin, L., and Chang, J., Leveraging the Path Signature for Skeleton-based Human Action Recognition. *arXiv preprint*, arXiv:1707.03993.