We are IntechOpen,
the world's leading publisher of
Open Access books
Built by scientists, for scientists

## 5,300
Open access books available

## 130,000
International authors and editors

## 155M
Downloads

Our authors are among the

## 154
Countries delivered to

## TOP 1%
most cited scientists

## 12.2%
Contributors from top 500 universities

**CLARIVATE ANALYTICS**
**BOOK CITATION INDEX**
**INDEXED**

**WEB OF SCIENCE™**

Selection of our books indexed in the Book Citation Index
in Web of Science™ Core Collection (BKCI)

## Interested in publishing with us?
## Contact book.department@intechopen.com

Numbers displayed above are based on latest data collected.
For more information visit www.intechopen.com

**Chapter**

# Predicting Type 2 Diabetes Complications and Personalising Patient Using Artificial Intelligence Methodology

*Leila Yousefi and Allan Tucker*

## Abstract

The prediction of the onset of different complications of disease, in general, is challenging due to the existence of unmeasured risk factors, imbalanced data, time-varying data due to dynamics, and various interventions to the disease over time. Scholars share a common argument that many Artificial Intelligence techniques that successfully model disease are often in the form of a "black box" where the internal workings and complexities are extremely difficult to understand, both from practitioners' and patients' perspective. There is a need for appropriate Artificial Intelligence techniques to build predictive models that not only capture unmeasured effects to improve prediction, but are also transparent in how they model data so that knowledge about disease processes can be extracted and trust in the model can be maintained by clinicians. The proposed strategy builds probabilistic graphical models for prediction with the inclusion of informative hidden variables. These are added in a stepwise manner to improve predictive performance whilst maintaining as simple a model as possible, which is regarded as crucial for the interpretation of the prediction results. This chapter explores this key issue with a specific focus on diabetes data. According to the literature on disease modelling, especially on major diseases such as diabetes, a patient's mortality often occurs due to the associated complications caused by the disease over time and not the disease itself. This is often patient-specific and will depend on what type of cohort a patient belongs to. Another main focus of this study is patient personalisation via precision medicine by discovering meaningful subgroups of patients which are characterised as phenotypes. These phenotypes are explained further using Bayesian network analysis methods and temporal association rules. Overall, this chapter discussed the earlier research of the chapter's author. It explores Artificial Intelligence (IDA) techniques for modelling the progression of disease whilst simultaneously stratifying patients and doing so in a transparent manner as possible. To this end, it reviews the current literature on some of the most common Artificial Intelligent (AI) methodologies, including probabilistic modelling, association rule mining, phenotype discovery and latent variable discovery by using diabetes as a case study.

**Keywords:** diabetes, complex disease progression, artificial intelligence in medicine, patient model, Bayesian statistics, causal networks, data mining, hidden risk factors

## 1. Summary

Intelligent systems, whether biological or artificial, require the ability to make decisions under uncertainty using the available evidence. Several computational models exhibit some of the required functionality to handle uncertainty. These computational models in Artificial Intelligent (AI) and Machine Learning are judged by two main criteria: ease of creation and effectiveness in decision making. For example, Neural Networks (NNs) which represent complex input/output relations using combinations of simple nonlinear processing elements, are a familiar tool in AI and computational neuroscience. Alternatively, probabilistic networks (also called Bayesian Networks) are a more explicit representation of a domain through modelling the joint probability distribution (the probability of all possible outcomes in a domain). This paper provides a short summary of the previous methods in Intelligent Data Analysis (IDA) in disease progression, decision making and probabilistic modelling of patients. It then describes some existing key methods that can be updated or combined to model multiple diabetes complications in the presence of unmeasured factors. There is considerable research on predicting diabetes, especially Type 2 Diabetes Mellitus (T2DM), complications. Nevertheless, the previous research of the author is discussed in order to address these issues. In particular, these previously proposed methods have contributed to the diabetes literature by explaining unknown risk factors and identifying temporal phenotypes employing hybrid methods (including descriptive and predictive). These suggested methodology includes rule-based methods for an explanation of patient subgroups and a probabilistic framework for modelling data explicitly.

## 2. Literature review: intelligent data analysis in complex disease progression modelling

This article reviews the current literature on some of the most common AI methodologies, including probabilistic modelling, association rule mining, and latent variable discovery. Intelligent Data Analysis (IDA) is a subcategory of AI that is focused on data analysis and modelling. These methods are known to be highly successful in combining advantages of modern data analytics, classical statistics and the expertise of scientists and experts [1–3]. IDA techniques have already proved successful in clinical modelling [4]. A large and growing body of literature has investigated IDA approaches that have shown excellent results modelling cross-sectional clinical data for classification. There has also been substantial modelling on longitudinal data using IDA techniques. However, there is still an urgent need to improve these models to take account of the variability of disease progression from person to person, and explicitly model the time-varying nature of the disease. Many studies have attempted to find automated ways of helping clinicians predict disease progression [5].
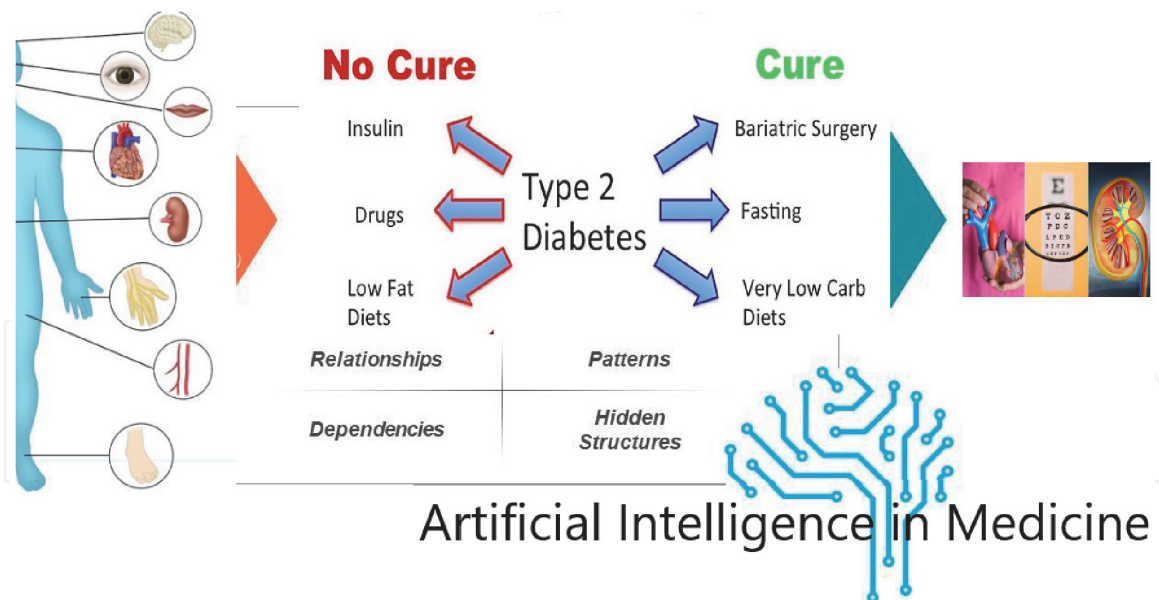
For many clinical problems, the underlying structure of unmeasured variables may play an essential role in the progress of the disease. However, it is still a relatively unexplored area. Identifying these unmeasured variables as hidden or latent variables is key. What is more, understanding the semantics behind these unmeasured risk factors can improve the understanding of the disease mechanisms and thus better improve clinical decision making. Interpreting these latent variables is complicated; however, as they may represent different many types of unmeasured information such as social deprivation, missing clinical data, environmental factors, time-based information or some combination of these. To gain trust in any AI model, it is mandatory to understand/explain influencing factors of

disease that guide predictions or decisions. This is because clinicians expect to understand AI diagnoses to be able to make decisions. There is a great deal of debate over the importance of explanation in AI models inferred from health data. In particular, there is a balance that needs to be made between the accuracy of complex deep models such as convolutional neural networks (in predictive strategy) and the transparency of models (in descriptive strategy) that aims to model data in a more human way such as expert systems.

A combination of explainable and deep strategies rather than either one of them alone would have a better prognostic value. Furthermore, in order to obtain a more accurate and explainable prediction of progression, the predictive models need to be personalised based on how an individual patient matches historical data by identifying patient subgroups.

## 3. Probabilistic model for time-series analysis

Understanding the pattern of complications associated with the disease has been used significantly in the clinical domain [6]. It provides an insight into the prediction and relative prevention of the associated complications which are expected to occur in a patient follow-up [7]. It generally can lead to less suffering time for patients while saving time and cost to healthcare. However, that is highly dependent on the stage of disease along with the prior occurring complications, which is associated with time-series analysis. In time-series analysis, every disease risk factor and complication is determined by various features in previous patient visits (time interval). At every medical visit, all diabetic patients have a unique profile of symptoms and complications that change over time, regardless of the phase of the disease. This non-stationary characteristic of clinical data collected as part of the monitoring of T2DM creates a difficult context for effective forecasting [8]. Clinical data needs to be considered as time-series data in order to provide a description of the progression of a disease over time. Nevertheless, dealing with time-series patient records is known to be a major issue in the prognosis of comorbidities [9], particularly when time-series data is imbalanced and contains few examples of patients without comorbidities that are common to all patients. In Type 2 Diabetes, for example, once patients are diagnosed with T2DM, half of them show signs of complications [10]. Unfortunately, these life-threatening complications remain undiagnosed for a long time because of the hidden patterns of their associated risk factors [11]. If T2DM is not appropriately managed, the development of serious complications, such as neuropathy, retinopathy, and hypertension lead to disability, premature mortality and financial cost [12]. The prediction process is complex due to the interactions between these complications and other features, as well as between complications themselves. More importantly, each patient has a unique profile of complication occurrence and the status of T2DM risk factors during a patients time-series is subject to change, as their levels may rise and fall over time. Early diagnosis and prevention techniques are needed to reduce the associated mortality and morbidity caused by T2DM complications [13]. Although there are various methodologies for T2DM prediction, for example, risk-prediction equation and Markov models [14], studies that enable early predictions of diabetes using predictive models are limited [15]. The risk-prediction equations suffer from uncertainty as well as performing only one-step-ahead predictions, while Markov models are limited to a small number of discrete risk factors. Other existing literature on investigating the prognosis of T2DM complications [16, 17] focuses particularly on logistic regression and Naive Bayes. Such studies are unsatisfactory for modelling the complex T2DM complications/risk factors. Logistic regression does

**Figure 1.**
*The organs/muscles affected by the common complications associated with type 2 Diabetes.*

not perform well when there are multiple or non-linear decision limitations. In Naive Bayes, there is an assumption of independence among the risk factors whereas all features are independent of one another.

The major limitation of the previous work in T2DM literature derives from time discretisation in temporal time slices per year. Therefore, in this study, we consider all T2DM patient's follow-up visits regardless of year basis while precisely monitoring the location of change within the unequal number of visits. This chapter suggests that AI in Medicine can provide useful techniques to analyse patient data to be able to find cure for the disease or reduce patient's suffering time (see **Figure 1**).

## 3.1 Dynamic Bayesian networks

In the field of medical informatics, probabilistic IDA techniques are exploited to obtain different clinical solutions. To improve patients' quality of life, there is an urgent need to extend and explore probabilistic IDA methods to investigate the disease complications from a clinical point of view. Thus, a Bayesian Network (BN) decision model was exploited in [18] for supporting the diagnosis of dementia, Alzheimer disease and mild cognitive impairment. Bayesian Network models appear to be well suited T2DM progression modelling, because of their flexibility in modelling spatial and temporal relationships as well as their ease of interpretation [19]. It has been reported that Dynamic Bayesian Networks (DBNs) are simple BNs for modelling time-series data and popular for modelling uncertain noisy time-series clinical data [20]. More importantly, DBNs are probabilistic graphical models that can handle missing data and hidden variables.

Previous work on learning DBNs have inferred both network structures and parameters from (sometimes incomplete) clinical datasets [20]. For example, a recent study presented a DBN method but to analyse fisheries data [21]. Authors in [22] proposed a Bayes Network to predict diabetes on the Pima Indian Diabetes dataset. However, the study failed to consider the time-series analysis. Similarly, authors in another study [23] simulated the health state and complications of type 1 diabetes patients by using partially and entirely learned Bayesian models. Apart from using a different type of Diabetes, this chapter is utilising a different approach from the above studies for the representation of the relationship between T2DM

risk factors. Many diseases involved structural changes based upon key stages in the progression, but many models did not appear to take this into account. There has been some work in extending DBNs to model underlying processes that are non-stationary [24]. In [24], clinical features were modelled using a second-order time-series model while time-invariant temporal dependencies were assumed. Among this, some studies, for example, Marini and co-authors conducted research [23] that variables were connected within two-time-series and within the same time slice assumed that the temporal dependencies were time-invariant. In addition, in Marinis paper for learning the network structures, a Tabu search was used based on the Hill climbing algorithm for Bayesian Networks but with no use of latent variables. However, the approach was useful for stratifying patients according to the probability of developing complications, the major limitation of the Marinis work derived from time discretisation in time slices of one year.

Another work in [25] retained the stationary nature of the structure in favour of parameter flexibility, arguing that structure changes lead almost certainly to over-flexibility of the model in short time-series. Alternatively, a paper [26] formalised non-stationary DBN models and suggested MCMC sampling algorithm for learning the structure of the model from time-series biological data. Similarly, authors in [27] estimated the variance in the data structure parameter with an MCMC approach, but the search space was limited to a fixed number of segments and indirect edges only, which is not suitable for T2DM data. Such studies remained narrow and limited by constraints on one or more degrees of freedom: the segmentation points of the time-series, the parameters of the variables, the dependencies between the variables and the number of segments and the ignorance of the incomplete data and latent variable.

## 3.2 Dealing with time-series imbalanced data

Another common problem with classifying complications in longitudinal data is that there may be many more cases where the complication does not manifest compared to those where it does. Early prediction of T2DM complications while discovering the behaviour of associated aggressive risk factors can help to improve a patients quality of life [28]. This study suggests that while there is an association between the latent variable and joint complications in the prognosis of T2DM patients, this relationship is complex. In T2DM data analysis, another challenge can be to classify/group patients in imbalanced clinical data with several binary complications. Models of the time-series data are needed to manage diabetic complications and deal with their imbalanced and complex interactions. In particular, mining time-series is one of the challenging problems in the prognosis of disease. In addition, it has received considerable critical attention in data mining especially when there are rare positive results [29]. It has been reported that a class imbalance in the training data caused by one class (here positive cases) massively outnumbers the examples in another class (negative class) [30]. This situation may occur where the number of positive clinical test results for a complication is not equal or even close to the number of negatives. That can be solved by applying an appropriate balancing strategy in a multi-class classification problem. Different learning techniques deal with imbalanced data, such as oversampling, undersampling, boosting, bagging, bootstrapping, and repeated random sub-sampling [31]. Therefore, this chapter in order to prepare T2DM data for the prediction has utilised these strategies and customised them based on dataset nature (time-series patients records with the unequal number of visits). As a result, various balancing strategies such as pair-sampling, bootstrapping undersampling and over-sampling have been proposed in [32, 33].

The bootstrap approach can be used to identify the significant statistics from classifiers learnt from such data. For example, in a study [34], Li and co-authors provide an extension to the temporal bootstrap approach while applied on cross-sectional data. Similarly, a study conducted in [21], the bootstrap strategy is extended to longitudinal data by sampling pairs of time points, thus enabling the (first-order) temporal nature of the data to be inferred. However, these solutions only can be applicable when the imbalance ratio for all binary complications is similar. Otherwise, it can be more difficult if we need to over-sample one class value and under- sample others in order to reduce bias from data. Overall, the observed balancing strategies from the prior studies have not been sufficient for analysing more than one complication at a time, whereas it was almost impossible to obtain a satisfactory prediction performance enhancement for all complications. As well as modelling unmeasured factors, hidden variables can also be used to model non-stationary processes. This chapter attempts to address this issue by using hidden variables discovery approaches based upon T2DM risk factors/complications dependencies. Before explaining these strategies, it is necessary to understand unmeasured variables and analyse their dependencies that are generated by causal structures.

### 3.3 Causal structure learning and latent variable discovery

Various studies on longitudinal data sets have suggested an association between complications and risk factors of the disease. To discover probabilistic dependencies given clinical data, it is necessary to search the space of belief networks or casual models, which is called casual discovery of BNs [35]. These patterns of dependency with no model based solely upon the observed variables can be explained by using a latent variable. The casual discovery indicates dependencies that are generated by casual structures with unmeasured factors, i.e., hidden variables. Hidden variable modelling, introduced in [36], has a long tradition in casual discovery. One of the research gap in the previous literature of disease prediction is the existence of the unmeasured or latent variables. This is because clinicians cannot measure all risk factors and carry out all kinds of tests, so there are some unmeasured factors that clinicians fail to measure, which need to be discovered at the early stage of diabetes.

Furthermore, Factor Learning (FL) was introduced in [37], which has been known as one method for learning a probabilistic model from data. It can also be helpful to understand latent variables and measure their hypothetical impacts. FL contrasts with most other BN learning methods in that it learns a factor structure. As Martin and co-authors in [37] stated that FL for hidden variables could identify the most probable structures of factors have given the data and suitable priors. However, with a large number of variables, FL methods might be prohibitively expensive. Again in the same research these authors provided a factor structure for learning methods that efficiently utilised hidden variables. Factor structure indicates the joint probability distribution among discrete observed variables. It also contributes an explanation across a small number of variables. Although factor structures are suitable for polynomial time inference, they can cause a reduction in the prediction accuracy and precision; they contribute an explanation across a small number of variables. Nevertheless, these techniques failed to consider prior belief in the factor structure, and therefore, it could be hard to rely on the final structure.

Factor structure indicates the joint probability distribution among discrete observed variables. Interestingly, each factor in a factor structure corresponds to a completely connected dependency graph. Although they are suitable for polynomial time inference, caused reducing accuracy and precision. By contrast, they are not able to decide precisely whether or not latent variables are present, and in

consequence there has been some controversy about that status of exploratory versus confirmatory factor analysis. In this regard, casual discovery methods in AI have the advantages as they can discover the actual dependencies and independencies in the data.

The causal discovery of BNs is a critical research territory, which depends on looking through the space of causal models for those which can best clarify a pattern of probabilistic conditions appeared in the data [35]. As a result, [38] showed the integration of structure-search algorithm with a latent variable in a DBNs model. However, the method did not consider the discovery of the long-range dependencies with an equal number of time slices. Similarly, in [39], Bayesian belief networks was used to find the most probable structure, using the K2 algorithm, while adding a hidden variable. Nevertheless, Cooper in [39] applied the K2 method that needs an ordering on the nodes. Witting focused on using hidden variables in a known structure [40]. Cooper in [39] used Bayesian techniques to find the most probable structure and can use this technique to add hidden variables. In principle, exact Bayesian methods for hidden variables could identify the most probable structures of factors given the data and suitable priors. However, with a large number of variables, exact methods are prohibitively expensive. Furthermore, in [41] Silva highlighted the weakness of DAG (Directed Acyclic Graph) models in the marginalisation of Hidden factors and representing the independencies over a subset of features in a DAG with more links. They suggested that Directed mixed graphs (DMGs) are a solution to this drawback. Therefore, they represented how to perform Bayesian inference on two DMGs, such as Gaussian and Probit, which is not the focus of this chapter. Nevertheless, such studies remained narrow and limited by constraints on one or more degrees of freedom: the segmentation points of the time-series, the parameters of the variables, the dependencies between the variables and the number of hidden factors. As a result, Chicharro in [42] analysed causal influences to find the relationship among different brain regions in several disorders. Similar to this chapter, Chicharros research made use of Inductive Causation (IC*) algorithm in the latent process to analyse Granger causality and Dynamic Causal Modelling. However, Chicharros study did not consider DBNs to understand causal influences.

Difficulties arise, however, when an attempt is made to implement a Bayesian Network structure as authors in [43] have argued that the number of potential DAGs over the disease risk factors is super-exponential. Additionally, the real cause-effect relationship DAG is not distinguishable while from equivalent structures when learning only using from observational data. This issue will be worse, especially when each expert has a unique probability of correctly labelling the inclusion or exclusion of edges in the disease structure. As noted by Amirkhani [43], some scoring functions are provided with that score each suitable graph based on the data and experts knowledge. Another research in [44] shows that networks with the fixed structure containing hidden variables can be learned automatically from data using a gradient-descent mechanism similar to that used in neural networks. A few algorithms have been created to understand the structure for Bayesian Networks from both fully observed models and those with hidden variables. Structure Expectation–Maximisation (SEM) has been produced for learning Probabilistic system structure from information with latent factors and missing data. A structure learning algorithm has been created for non-stationary dynamic probabilistic models. For example, REVEAL (REVerse Engineering ALgorithm) has been utilised as a structure learning algorithm, that learns the optimal set of parents for each node of a network independently, based on the information-theoretic concepts of mutual information analysis. However, the two-stage temporal Bayes network (2TBN) cannot be well recovered by the application of REVEAL. A normally

utilised structure learning algorithm depends on REVEAL which takes in the ideal arrangement of guardians for every hub of a system autonomously, in light of the theoretical data ideas of common data examination. Be that as it may, the two-arrange fleeting Bayes organise as the 2TBN which cannot be all around recuperated by use of REVEAL. Rijmen in [45] exploited an HMM to study the temporal pattern of symptoms burden in brain tumour patients. He showed that the discovery of symptom experience over time is necessary for treatment and follow-up of patients with symptom-specific intervention. In general, Bayesian learning methods could determine network structure and how the networks variables should be represented along with the causal links among them. Moreover, it addressed the difficulty of qualifying causal relationships in terms of Conditional Probability Tables (CPTs). Witting focused on using hidden variables in a known structure [40] as the knowledge of the latent variable in predictive modelling is important for an understanding of the complex AI models. Discovering latent variables can potentially capture unmeasured effects from clinical data, simplifying complex networks of interactions and giving us a better understanding of disease processes. In addition, it can improve classification accuracy and boost user confidence in the classification models [46]. Elidan and co-authors in [47] emphasised the importance of the presence of hidden variables. In addition, they determined a hidden variable that interacted with observed variables and located them within the Bayesian Network structure. They also showed that networks without hidden variables are clearly less useful because of the increased number of edges needed to model all interactions, which caused overfitting. Despite the productivity of exploring trees of hidden variables to render all observable variables independently [48], these hidden variables were non-optimal with independencies among observable variables. Overall, previous works on learning DBNs have presented both network structures and parameters from clinical data sets and learning parameters for a fixed network of incomplete data, in the presence of missing data and latent variables [20]. Much of the current literature on disease prediction have argued that a complex AI model, with many unexplainable hidden variables, also has several serious drawbacks. Therefore, this chapter has chosen AI DBNs model to learn parameters and latent variables to predict complications. The next section intends to emphasise the explainability of the proposed methodology in order to uncover the meaning behind the latent AI model.

## 4. Black box models and AI in medicine

Investigating unmeasured risk factors can improve the modelling of disease progression and thus enable clinicians to focus on early diagnosis and treatment of unexpected conditions. However, the overuse of hidden variables and lack of explainability can lead to complex models, which are not well understood (being black box in nature). Models need to be understood by clinicians to facilitate transparency and trust.

### 4.1 Explainability in deep learning

This stage outlines and discusses the limitations of Deep Learning approaches that have been proposed, so far, to gain deeper insights into the understanding of black box AI models. AI medical machine such as Deep Learning has become ubiquitous to provide a high-performance prediction. Nevertheless, understanding their mechanisms has become a significant concern worldwide whereby the goal is to gain clinicians and patients trust. The reason behind this is due to several

obstacles that arise to interpret the findings, such as the scale of big data, complex interactions, and high-dimensional internal state.

### 4.1.1 Google's novel approach

Most medical algorithms proposed by [49], such as AI Doctor designed to reproduce current problem-solving methods (e.g., the detection of cancers). In addition, the concept assignment can help people to strengthen their skills and talents for a computer system that showcased superhuman effectiveness and efficiency.

Google's AI Doctor can be demonstrated how they could be used to provide an explanation further into predictions generated by local classifiers, first from conventional image classification networks to a focused clinical application. The concept attribution approach in AI Doctor offers several promising avenues for future work. In addition to this, the concept assignment can help people to strengthen their skills and talents for a computer system that showcases superhuman effectiveness and efficiency. The concepts of explanatory power are outlined by Google under three principle assumptions/limitations: firstly, comprehension for whatever hidden layer and artificial neurons would offer. This is based on most of the information in a deep neural network consists of hidden layers. Secondly, it recommends that acknowledging the numerous hidden layers and understanding their design on a meta-level would lead to more in-depth modelling. Finally, to comprise how nodes become active, it considers groups of interconnected neurons that trigger at the same time and space. These principles are defined instead of explaining the structural nature of each neuron in each network. This is because the stratification of a network for the categories of interconnected neurons would enable its configurations even more abstractable. This is the main weakness of the black box models.

One of the most highlighted ones is Google's approach to resolve the explainability issues while enabling human-like description of the internal state of a deep network by employing Concept Activation Vectors (CAVs). While medical systems are mostly designed to reproduce current decision-making methods such as the classifier used in the detection of cancers, Google has claimed that its novel strategy can interpret existing clinical data. Although Google has made a claim that the CAVs can directly relate to one's anticipated theories, to draw conclusions about the decision-making process, it needs to consider the human needs of a higher level of understandability.

Nevertheless, cardiac specialists have been critical of the conclusions derived by Google in the clinical domain. With the proper information, AI is optimistic that innovative, unique healthcare insights might be created without human intervention. Unfortunately, this new approach is only established based on extensive and adequate datasets. This is presumably part of the explanation of why Google has established projects as its benchmark research proposal is capturing detailed patients' history of 100,000 population across four years. However, the investigation conducted out by Google did not necessarily indicate that the suggestion was entirely distant. Such as image classifiers that could be applied to low-level structures. The central concept and assumption are to consider a neural network as additional assistance that can cause issues related to the internal representation. As a result, the clinicians commented on the deep explanatory networks. They questioned the hypotheses, by stating that although the AI algorithms and Deep Learning could improve current prediction methods of clinical domain, the research would not be trustworthy unless it had been assessed with caution while a broader range of disease had been explored. Difficulties arose, when an attempt was made in

order to implement the principles and these assumptions. It seemed to be evident that their approach was overconfident and yet to be trusted.

### 4.1.2 Prototyping examples in Artificial Neural Networks

In order to introduce a different perspective on Deep Learning models' interpretability, Zintgra and co-authors [50] conducted a study to simplify the black box structure of Artificial Neural Networks (ANNs). They made use of prototypic examples method that indicate tools in order to diagnose trained ANNs. In general, ANNs analyse discrete decision-making processes and obtain high-performance prediction results.

The prototype examples may be computationally intractable, including a pre-determined normal distribution to prevent the proliferation of unreasonable prototype cases. They provided an explanation of tools to train ANNs based on two datasets. Moreover, it can often be like such a losing battle to describe precisely how ANNs operate mathematically. Therefore, a much more comprehensive pre-processing methodology could also be used in a related development (e.g., generative adversarial network proposed by Goodfellow et al. in [49]). Furthermore, experimental results and hypotheses in ANNs were portrayed and tested only on two datasets. Alternatively, a more detailed analysis is required to rely on the empirical results, which might be achieved by including rich data containing imbalance issue, different types of features. Selection bias was another potential concern because it could involve possible measurement errors. It could be extended through more set of data with various features. Finally, conclusions and interpretations of data were drawn from an inevitably subjective mechanism on the investigator's basis. This was because to examine whether the produced case studies should satisfy the investigator's standards about the phenomena of been modelled (e.g., decisions could be only made by the time it came). This was established based on approaches or standards for collecting and analysing concepts that might be more unbiased. As a result, this could also enable investigators/analysers to understand the implications and weaknesses of the use of ANNs for the discrete decision-making process, which might enhance the strictness of the approach. However, many healthcare methods are required to reconstruct conventional prediction methods (e.g., the identification of cancers), but so far, different ideas to interpret previous clinical records have been discovered.

### 4.1.3 Visualisation in deep learning

For the time being, the possibility of an AI physician planning to roll new prognosis without direct human intervention is a significant distance in which the more presumably in decades rather than a few years later. Recent developments in several technologies in the Deep Learning area have been powered by the steadily declining expense of computing and storage. That being said, realistic apps, including certain integrated smartphone and electronic devices, have intensified explainability issues for Deep Learning in the black box resource-limited environments. Liu et al. in [51] introduced the leading solution to address these issues where a deteriorated image of Binary Convolutionary Networks caused by binarising Filtres. They offered a range of Circulant Filtres (CiFs) and a Circulant Binary Convolution (CBConv) to strengthen efficiency and to tackle those limitations for Binary Convolutionary functionalities through their proposed Circulant Backpropagation (CBP). Then, CiFs effortlessly was integrated into the current deep neural networks (DCNNs). Enormous research has indicated that perhaps the output difference among one-bit and total-precision DCNNs could be reduced by

extending the variety and distributing the filtres. Zintgraf et al. in [52] identified numerous tools to test the model and understand how DCNNs could provide a reliable outcome by using the visualisation method.

Overall, the existing explanatory Deep Learning approaches would need to be adapted for further sophisticated longitudinal modelling strategy (rather than with a multivariate distribution). This would result in better outcomes, for example, in pixel values which could be estimated reliably by everyone's environment while it skewed down much more. By providing the black box models with sufficient data, machine learning seemed to be overconfident that completely different health knowledge could then be generated without user intervention. The black box models of Deep Learning can be simplified in several aspects. For example, if an object is detected, an image detection machine can breakdown back and towards specific attributes including shape, colour and texture of the image, and then reduce the predictions to a mathematical method by checking the classification error and then background diffusion to improve the practises. In particular, in the world that it is possible to fully allocate decision making to computer systems, confidence in AI systems will be hard to achieve. In the future work, one approach that can be applied to the small-sized T2DM dataset can be the use of Bayesian Neural Networks, which will deal with uncertainties in data and model structure by exploiting the advantages of both Neural Networks and Bayesian modelling. To conclude, AI can improve current methods of medical diagnosis in terms of interpretability but cautioned that the technology would need to be more evaluated to be trusted by both patients and practitioners.

In black box models, it can be challenging to determine what is coordinating the visible patterns. Such models are problematic not only for lack of transparency but also for possible biases inherited by the algorithms from clinicians mistakes [53]. This issue is caused based on the human errors and biased sampling of training data as well as the underestimation of the impact of the risk factors underlying behaviour/pattern. In general, as observed from prior studies, it is difficult to obtain performance enhancement while simultaneously trying to explain hidden factors. Lakkaraju in [54] suggested that there is a trade-off between patient personalisation (in a descriptive analysis) and prediction performance (in predictive analysis). Generally speaking, an improvement in explainability is often possible through a less accurate model or at a higher cost of the predictive accuracy (in a Black box model) [6]. There are quite few research studies on predicting T2DM complications and T2DM black box models. However, studies on explaining an unknown risk factor/latent phenotype by using a hybrid data mining methodology (including descriptive and predictive) are rare to find in literature. Therefore, this study attempts to open the AI, black box model by using both predictive and descriptive strategies.

## 5. Patient personalisation and explanation

Most of the previously published studies in diabetes prediction have tended to focus on all patients as one integrated database rather than separating patients [16]. It can be challenging to stratify patients based on their longitudinal data in order to determine what is triggering the visible patterns that may be specific to one cohort of patients. There is some research, such as [55] that assesses the disease prediction performance based upon different IDA techniques. For example, the onset of the disease is modelled in [56] while other studies focus on patient modelling [57]. The approach described in this chapter aims to personalise patients by using unsupervised methodologies to group time-series patient data.

The proposed descriptive strategy in this chapter has been regarded as a useful tool known as association rules to detect interesting relationships among T2DM complications.

## 5.1 Time-series clustering

Time-series clustering is often problematic [58], especially when we need to analyse risk factors from matching patterns across time. The literature on time-series clustering and pattern discovery has highlighted several studies [59]. There have been some qualitative measures for clustering time-series data, which captured similar risk factor patterns in dynamic temporal data, regardless of whether the correlation between them was linear or not [60]. However, they did not seem to be very suitable for a long and an unequal number of time-series data (e.g., T2DM data). For instance, authors in [59] proposed an algorithm to cluster patients based on clinical data whilst utilising the clustering information for identifying distinct patterns. Altiparmak in [59] provided a slope-wise comparison method (SWC) to find the correlation between local distance vectors of patients visits, and group clinical test results into different sub-groups, based upon the related risk factors, by using feature selection. In their method each cluster of patients was considered as a transaction data that included a pattern indicating which cluster belonged to each patient. Authors in [61] used a similar method [59] in clustering, but they clustered fixed length time-series. Ceccon and coauthors [62] exploited a variation of the naive Bayes classifier with a hidden variable for segmenting patients into disease sub-types. Ceccon's study intended to enhance the classification performance of Glaucoma patients based upon visual field data. Nevertheless, they only focused on standard/static BNs (instead of DBNs) to infer the parameter in a cross-sectional dataset. Moreover, they failed to analyse the influences of multiple hidden variables on the prediction results.

## 5.2 Pattern discovery and association rules mining

It has previously been observed that patients with T2DM are also at an increased risk of microvascular comorbidities, including nephropathy, neuropathy, and retinopathy [63]. The underlying pattern of T2DM complications and how their co-occurrence is followed/caused/related by other complications associated with the disease, known as the major source of mortality and morbidity in T2DM [64]. That is because predicting a target complication can be challenging without the consideration of the effects of its associated complications. Similar to Diabetic type 1 patients, although genetic factors impact on developing T2DM, it is believed ignorance of developing complications harms patients' life. What is more, T2DM patients develop a different profile of complications and features, which changes over time per follow-up visit. One of the most important factors in the high number of dependencies among T2DM features and complications is the appearance of unmeasured risk factors. Surprisingly, the effect of understanding unmeasured variables, which play an important role in disease prediction, does not seems that closely examined.

Understanding these associated patterns has a remarkable actual value and can significantly being used in the clinical domain [6]. It provides an insight into the prediction and relative prevention of the associated complications which are expected to occur in patient followups [7]. It also leads to less suffering time for patients while saves time and cost to healthcare. However, that is highly dependent on the stage of disease along with the prior occurring complications, which is associated with time-series analysis. In time-series analysis, every disease risk factor and complication is determined by various features in previous patient visits (time

interval). To better understand the complications of the disease and their effects, this chapter clusters patient the associated rules among the complications. It attempts to address this issue and present an informative rules/ordering pattern of patient behaviour, with an aim to capture the complexities of the associated complications' over time. The proposed descriptive strategy has been regarded as a useful tool known as association rules (ARs) to detect interesting relationships among T2DM complications.
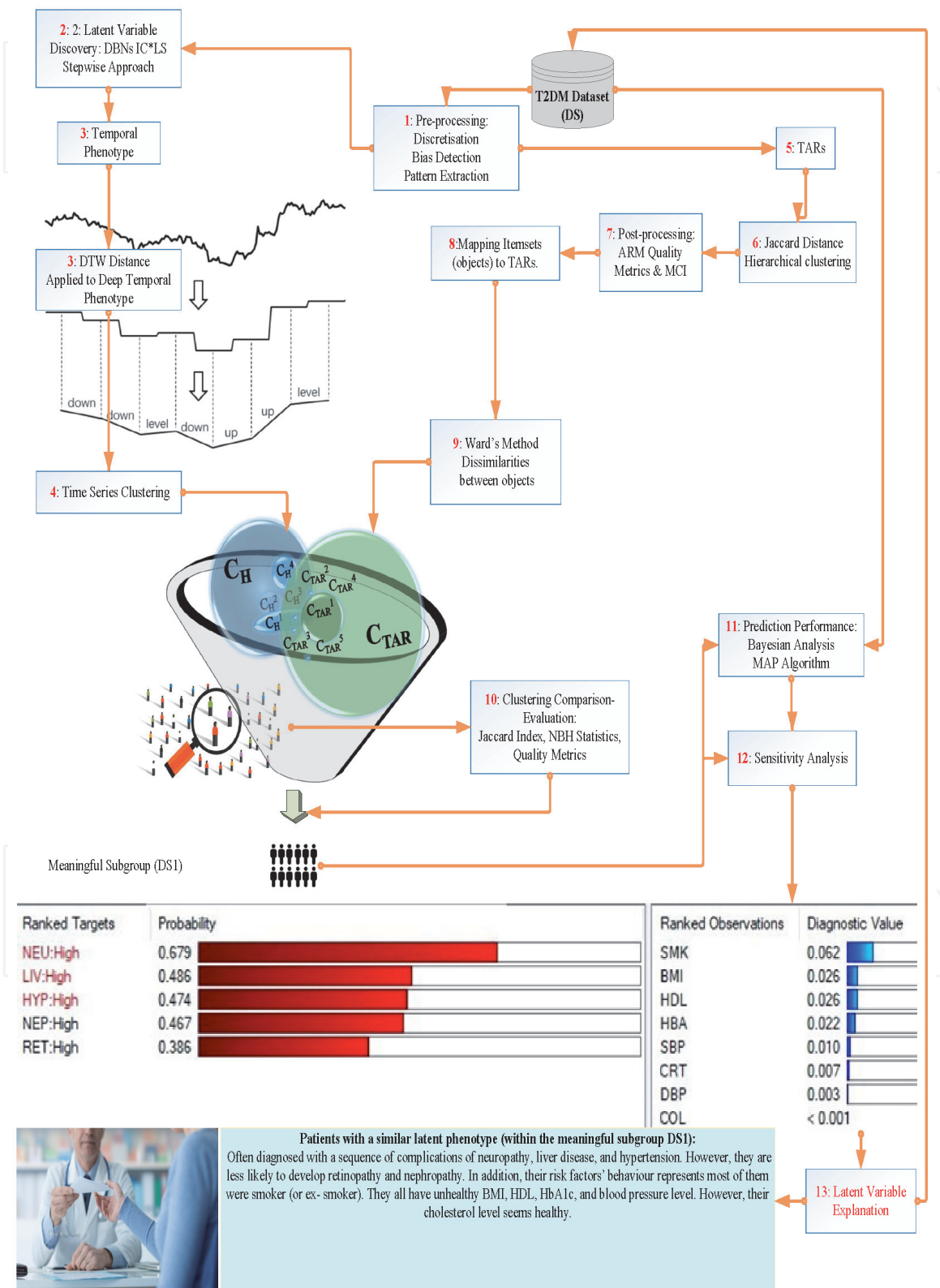
Temporal Association Rules (TARs) [65] is an extension to association rules [66] to analyse basket data that includes a temporal dimension to order related items. Many algorithms with temporal rules work by dividing the temporal transitions database into different partitions based on the time granularity obliged. For example, different mining algorithms were reformulated and presented to reflect the new general temporal association rules. These include Progressive Partition Minder (PPM), Segmented Progressive Filter (SPF), and TAR algorithm [65–67]. Various algorithms have been proposed for the incremental mining of temporal association rules, especially for numerical attributes [68]. Allen's rules [69] generalised abstracted time-series data into a relation (PRECEDES) to find TARs in [70]. Various ways were proposed to explore the problem of temporal association rules discovery [71]. Nevertheless, previous studies performed discovering association rules on a given subset specified by the time [72], whilst not considering the specific exhibition period of the elements.

Association Rule Mining (ARM) finds frequent patterns by mining ARs with the use of two basic parameters of support and confidence [73]. The majority of the previous ARM algorithms worked by dividing the temporal transitions database into different partitions based on the time granularity obliged.

Difficulties arise with TARs when there are some rare rules of particular interest [74]. Many studies have employed the most common filtering metrics rather than support and confidence in order to detect interesting rules [75]. There is a controversy to this, as a study in the literature argued that a conservative ARM methodology only based on a fixed and rigid threshold for the filtering metrics could be problematic. A few studies attempted to mine frequent underlying patterns of diabetic complications [76]. The frequent pattern mining research significantly affects data mining techniques in longitudinal data. A post-processing approach in [77] attempted to extract interesting subsets of temporal rules within T2DM data. However, it only considered characteristic patterns of administrative data without the appearance of latent variables. Other researchers have undertaken association rule mining of clinical data [78, 79]. Lee et al. attempted to address the issue in [67] and have led to the proposal of the concept of general TARs, where the items were allowed to have varying exhibition periods, and their support was made based on that accordingly. Another research conducted by Plasse et al. in [80] looked at finding homogeneous groups of variables. They suggested that a variable clustering method could be applied to the data in order to achieve a better result in pattern discovering methodology. However, their strategy to mine ARs differed from this chapter in which the number of rules was reduced only based on hierarchical clustering applied to items, not to multiple identical binary attributes. Among these, some methods uncovered temporal patterns and relationships among clinical variables, including causal information [81], numeric time-series analysis [82]. Nevertheless, considering all of this evidence, none of the above studies has clustered uneven time-series clinical data based on a hidden variable for extracting temporal phenotype and behaviours of patients.

# 6. The suggested methodology

This chapter, so far, has described the research gap in the modelling and explaining of complex disease processes and thus given the motivation behind the suggested methodology. The previously discussed methods suffer from some limitations in addressing imbalance issues, complex and temporal relationships between



**Figure 2.**
*The proposed hybrid methodology to find explainable subgroup of patients by personalising diabetic patients in precision medicine. This figure is an abstract methodology explained in Figures 1–4 in the previous work in [83].*

(sometimes unmeasured) factors, and the identification of different underlying characteristics of disease for different subgroups of the population. There is considerable research on predicting T2DM complications. Among these, studies on explaining unknown risk factors and identifying temporal phenotypes by using hybrid methods (including descriptive and predictive) are rare to find in literature. It represented the reason of the earlier research conducted by the author in [32, 33, 83–85]. The current work of this chapter's author has attempted to address these issues in the previous research in [32, 33, 84], after describing the case study data as a starting point, the suggested methodology is explored as a framework for modelling real time-series clinical data. In the recent work conducted in [83, 85], the identification of informative hidden factors is investigated followed by methods to cluster patients into meaningful subgroups along with the identification of a latent temporal phenotype and the characterisation of these groups using temporal association rules (as illustrated in **Figure 2**).

## 7. Type 2 Diabetes as a case study

The World Health Organisation (WHO) reported that Type 2 Diabetes Mellitus (T2DM) accounts for at least 90% of all diabetes types. Another study in WHO revealed that T2DM patients are at increased risk of long-term vascular comorbidities, which is known as "underlying cause of death" and severe phenotype of the disease [86]. It has previously been observed that patients with T2DM are also at an increased risk of microvascular comorbidities, including nephropathy, neuropathy, and retinopathy [86]. Similar to Diabetic type 1 patients, although genetic factors impact on developing T2DM, it is believed ignorance of developing complications harms patient life because it may develop a different profile of complications and features, which changes over time per follow-up visit. However, these life-threatening complications remain undiagnosed for a long time because of the hidden patterns of their associated risk factors [11]. The underlying pattern of the complications is known as the major source of mortality and morbidity in T2DM and how their co-occurrence is followed/caused by other complications associated with the disease [64]. That is because predicting a target complication can be challenging without the consideration of the effects of its associated complications.

### 7.1 Data description

The observed dataset in this chapter is similar to the data utilised in the previous study of Diabetes patients in [83] of pre-diagnosed T2DM patients aged twenty five to sixty five years (inclusive) that were recruited from clinical followups at the "IRCCS Instituti Clinic Scientifici" (ICS) Maugeri of Pavia, Italy. The MOSAIC project funds the information based on the seventh Framework Program of the European Commission, Theme ICT201152 Virtual Physiological Human (600914) from 2009 to 2013. These consists of physical examinations and laboratory data for complications and risk factors (predictors) in T2DM which were selected supported existing literature on T2DM [76, 87–90] as well as the recommendations from the clinicians at ICS. These are Retinopathy (RET), Hypertension (HYP), Nephropathy (NEP), Neuropathy (NEU) and LIVer disease (LIV) (see **Table 1**). Here, the predictors are known and selected from the dataset: Body Mass Index (BMI), Systolic Blood Pressure (SBP), High-density Lipoprotein (HDL), Glycated Haemoglobin (HbA1c or HBA), Diastolic Blood pressure (DBP), ChOLesterol (COL), Smoking habit (SMK) and Creatinine (CRT). Control Values for T2DM risk factors are classified in **Table 2** illustrates three clinical level of risk, particularly low (zero),

| Node ID | Target complication | Diagnosis outcome[a] | Clinical risk class[b] |
|---------|--------------------|--------------------|----------------------|
| 2 | Retinopathy (RET) | {Negative,Positive} | {low,high} |
| 3 | Neuropathy (NEU) | {Negative,Positive} | {low,high} |
| 4 | Nephropathy (NEP) | {Negative,Positive} | {low,high} |
| 5 | Liver Disease (LIV) | {Negative,Positive} | {low,high} |
| 6 | Hypertension (HYP) | {Negative,Positive} | {low,high} |

[a]*Negative test result, Positive test result.*
[b]*Low clinical risk, High clinical risk.*

**Table 1.**
*The description of T2DM target complication, clinical node control values, and discretised states [83].*

| Node ID | T2DM risk factors | Control value[a] | Discretised value[b] |
|---------|-------------------|------------------|---------------------|
| 1 | HbA1c (HBA) | $6.6 \pm 1.2$ (%) | {low,medium,high} |
| 7 | Body Mass Index (BMI) | $26.4 \pm 2.4$ (kg/m$^2$) | {low,medium,high} |
| 8 | Creatinine (CRT) | $0.9 \pm 0.2$ (mg/dL) | {low,medium,high} |
| 9 | Cholesterol (COL) | $0.9 \pm 0.2$ (mg/dL) | {low,medium,high} |
| 10 | High-Density Lipoprotein (HDL) | $1.1 \pm 0.3$ (mmol/l) | {low,medium,high} |
| 11 | Diastolic Blood Pressure (DBP) | $91 \pm 12$ (mmHg) | {low,medium,high} |
| 12 | Systolic Blood Pressure (SBP) | $148 \pm 19$(mmHg) | {low,medium,high} |
| 13 | Smoking Habit (SMK) | {0,1,2} | {low,medium,high} |

[a]*(Mean ± SD).*
[b]*low, medium, high.*

**Table 2.**
*The description of the T2DM clinical features, risk factors, control values, and the discretised states [83].*

medium (one) and high (two). In T2DM data, the worsening level of the micro-vascular diseases and HYP is known as a significant cause of death [91]. Even though micro-vascular complications such as RET, NEP, NEU are less frequent comparing to HYP, an inadequate estimation of them causes long-term suffering and life-threatening comorbidities [64]. Fowler and co-authors in [7] researched type 2 Diabetic American patients. This research utilised T2DM key risk factors such as HbA1c, SBP, and DBP to investigate relationships among complications such as HYP, NEP, RET, and NEU. In addition, LIV is a severe phenotype of diabetes and associated with T2DM complications, especially NEU [92]. Litwak analysed Russian diabetic patients in [93] which referred to the influence of macro-vascular and micro-vascular disease on one anther. For example, important features in T2DM dataset such as blood pressure, HDL, lipid, BMI, and HbA1c influence diabetic patients' complications. They also revealed that HDL has a negative effect on HYP, NEP, NEU, and RET, whereas HbA1c negatively associated with HYP. Again, a study conducted by Ramachandran [94] referred to the high prevalence of NEU and RET in Type 2 diabetes in India. Similarly, research in [76] suggested that most of the diabetic patients have objective evidence for some variety of NEU, but only a few of them have identified by symptoms. This research also showed that there is a strong association among NEP, NEU, and RET. This study only concentrates on five binary complications as the predictive target classes in a binary classification problem (with two categories of classes: "high" or "low" risk). Furthermore, a complication class value of low risk (zero) represents a patient visit in

which the complication is not present; otherwise, it is at high risk (one). For instance, a complication class value of zero represents a patient visit in which the complication is not present; otherwise, it is one. Alternatively, other risk factors associated with a patient (symptoms/clinical tests) are abstracted in the multi-class classification problems with more than two targets including high, medium, and low risk patient, according to a diabetes experts definitions [95, 96]. For each patient in T2DM dataset, time-series analysis is described in Appendix A with definition of the related notations.

## 8. Experimental results and conclusions

This section summarises the clinical implications and shows how the obtained experimental findings in the previous works [83, 84, 97] and their significance have led to developing explanatory AI models. For example **Table 3** illustrated the promising results obtained by the proposed Stepwise approach discussed in [33, 84].

In **Table 4**, the prediction performance of the underlying patterns of complications for these patients within the discovered subgroup dataset (was introduced in [83] as DS1 and discovered using the descriptive strategy) was analysed and compared to all patients belonged to DS (the raw T2DM dataset). It also suggested that DS1 (by personalising patients) could be considered as a dataset with less uncertainty than DS. In order to describe the inference problem in this chapter, the causal

| Percentage (%) | Accuracy | Sensitivity | Specificity | Precision |
|---|---|---|---|---|
| No Hidden variable in [32] | 48 | 53 | 48 | 53 |
| Stepwise IC* in [33] (Step1) | 60 | 40 | 80 | 70 |
| Enhanced stepwise in [97] (Step1) | 80 | 51 | 98 | 97 |
| Stepwise IC* in [33] (Step2) | 78 | 98 | 58 | 68 |
| Enhanced stepwise in [97] (Step2) | 95 | 80 | 96 | 86 |
| Stepwise IC* in [33] (Step3) | 78 | 98 | 58 | 68 |
| Enhanced stepwise in [97] (Step3) | 95 | 81 | 96 | 82 |
| Enhanced stepwise in [97] (Step4) | 96 | 81 | 97 | 092 |
| Enhanced stepwise in [97] (Step5) | 95 | 82 | 97 | 85 |

**Table 3.**
*Comparison of our new and enhanced stepwise IC\*LS approach in [97] with its previous version (stepwise IC\*) in [33] and without latent variable in [32].*

| Complication | Accuracy | |
|---|---|---|
| | **DS** | **DS1** |
| NEP | 0.81 | 0.93 |
| LIV | 0.77 | 0.88 |
| HYP | 0.91 | 0.99 |
| NEU | 0.76 | 0.81 |
| RET | 0.81 | 0.79 |
| All | 0.81 | 0.88 |

**Table 4.**
*The overall prediction accuracy of T2DM complications for patients in DS is compared to DS1.*

relationships seemed to be a reliable option to represent static and dynamic correlations between T2DM risk factors. The causal inference has a greater focus on distinguishing causes from other associations than on uncovering detailed temporal relationships. Therefore, in this work ([83]), several predictive strategies in order to test whether the descriptive approaches have contributed to improving the prediction performance of the ordering patterns of complications.

**8.1 Clinical implications**

This study offered several valuable insights into the prediction challenges in diabetes and similar diseases and explained how they could be tackled. First, throughout this chapter, appropriate machine learning techniques were conducted to model complex interactions among the complications, risk factors and unmeasured factors. For instance, the use of probabilistic graphical models provided a significant improvement in the accuracy of predictive models while reducing uncertainty in disease management. Having adopted DBNs to learn hidden risk factors and effectively understand the AI black box model was the key contribution of this research. The temporal phenotype was identified to represent the overall patterns of disease risk factors for each patient based on the discovered hidden variables over time. The descriptive analytics, in [97], provided valuable insights into the hidden variable effects on stratifying patients into different sub-groups, whether or not they developed the same complications. These findings also explained the influence of the latent variable on the bootstrapped data. Phenotype discovery was utilised to categorise and investigate meaningful subgroups of patients based on how an individual matches historical data. The hybrid type methods in discovering meaningful subgroups and explaining temporal phenotype also led to a better understanding of clinical data as well as aiding to interpret the unmeasured factors while demonstrating their risks.

**8.2 Future works**

The generalisability of the results presented in this study is subject to certain limitations as follows: This research was conducted to explain and discover the unmeasured factors with a few patients and relatively few features. Thus, this study focused on time-series complex clinical dataset like T2DM, which was a small-sized dataset with an unequal number of patient's follow-up visits (which is common in clinical data). This study was specific to T2DM concept and Bayesian modelling; hence, one fundamental criticism could be the bias towards this dataset and whether the method could be developed in other fields of clinical data in the future. In order to help overcome the limitations discussed in the previous section, the following recommendations are suggested: The originality of the proposed study consisted in its innovative, analytical, and methodological strategies to predict and explain complex clinical data to improve patients' quality of life. A natural progression of this work for a better generalisability should involve extending the latent DBNs model with more hidden variables to capture a greater variety of unmeasured factors to characterise critical changes and produce interesting findings that account more for better explainability and predictability. In addition, to address the limitation related to the small-sized dataset, this work could be extended to further investigation and experimentation into clinical impacts and environmental factors, such as family history, pollution, and glucose. More research also might be conducted to monitor disease progression effectively and detect the underlying patterns of complications, which could provide clinicians with a better understanding of the obtained findings. For example, a greater focus on phenotype discovery

could enable assessment of the long-term effects of the temporal phenotype on the patient, which might be done by following qualitative approaches to support the obtained findings from the biomedical literature. The generalisability of the findings obtained in this study might be tested on other data with potentially non-stationary, complex, and incomplete data. For instance, the pre-processing approaches, statistical analysis, temporal phenotype, MCI algorithm and the DBNs model could be applied to another complex data (e.g., COVID-19). Therefore, in a new project, a similar patient model to this research was mainly employed, which primarily concentrated on helping healthcare staff in their understanding of how COVID-19 spread and how they could be better prepared.

In the current work as a Post-Doctorate research fellow at Brunel University and University City London (UCL) associated with the BHF Alan Turing Institute jointly funded research project with the collaborators of the project in UCL and GSK. This project aims to develop a computational tool to investigate the action of drug compounds for the treatment of cardiovascular disease and type 2 diabetes which involves: firstly, the construction of a cardiovascular disease (CVD) and Type-II diabetes (T2D) relevant metabolic measures networks, using repeated measures. Secondly, the combination of different causal networks on the same set of metabolic measures. Lastly, the integration to the system of available drug targets and disease information for testing CVD and T2D drugs.

## Acknowledgements

## Author details

Leila Yousefi* and Allan Tucker
Life Science Department, College of Health, Medicine and Life Sciences, Brunel University London, United Kingdom

*Address all correspondence to: leila.yousefi@brunel.ac.uk

IntechOpen

## References

[1] R. Bellazzi. Big data and biomedical informatics: a challenging opportunity. *Yearbook of medical informatics*, 9(1):8, 2014.

[2] M. R. Berthold, C. Borgelt, F. Höppner, and F. Klawonn. *Guide to intelligent data analysis: how to intelligently make sense of real data*. Springer Science & Business Media, 2010.

[3] D. J. Hand. Intelligent data analysis: Issues and opportunities. In *International Symposium on Intelligent Data Analysis*, pages 1–14. Springer, 1997.

[4] R. Bellazzi, F. Ferrazzi, and L. Sacchi. Predictive data mining in clinical medicine: a focus on selected methods and applications. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1(5):416–430, 2011.

[5] M. B. Sesen, T. Kadir, R.-B. Alcantara, J. Fox, and M. Brady. Survival prediction and treatment recommendation with bayesian techniques in lung cancer. In *AMIA Annual Symposium Proceedings*, volume 2012, page 838. American Medical Informatics Association, 2012.

[6] T. Wang and Q. Lin. Hybrid predictive model: When an interpretable model collaborates with a black-box model. *arXiv preprint arXiv: 1905.04241*, 2019.

[7] M. J. Fowler. Microvascular and macrovascular complications of diabetes. *Clinical diabetes*, 26(2):77–82, 2008.

[8] M. A. Van Gerven, B. G. Taal, and P. J. Lucas. Dynamic bayesian networks as prognostic models for clinical patient management. *Journal of biomedical informatics*, 41(4):515–529, 2008.

[9] M. Van der Heijden, M. Velikova, and P. J. Lucas. Learning bayesian networks for clinical time series analysis. *Journal of biomedical informatics*, 48:94–105, 2014.

[10] R. Turner, R. Holman, D. Matthews, S. Oakes, P. Bassett, I. Stratton, C. Cull, S. Manley, and V. Frighi. Uk prospective diabetes study (ukpds). viii. study design, progress and performance. *Diabetologia*, 34(12):877–890, 1991.

[11] K.-H. Yoon, J.-H. Lee, J.-W. Kim, J. H. Cho, Y.-H. Choi, S.-H. Ko, P. Zimmet, and H.-Y. Son. Epidemic obesity and type 2 diabetes in asia. *The Lancet*, 368(9548):1681–1688, 2006.

[12] U. Diabetes. Diabetes: facts and stats. *Diabetes UK*, 3:1–21, 2014.

[13] S. Mani, Y. Chen, T. Elasy, W. Clayton, and J. Denny. Type 2 diabetes risk forecasting from emr data using machine learning. In *AMIA annual symposium proceedings*, volume 2012, page 606. American Medical Informatics Association, 2012.

[14] E. Mueller, S. Maxion-Bergemann, D. Gultyaev, S. Walzer, N. Freemantle, C. Mathieu, B. Bolinder, R. Gerber, M. Kvasz, and R. Bergemann. Development and validation of the economic assessment of glycemic control and long-term effects of diabetes (eagle) model. *Diabetes technology and therapeutics*, 8(2):219–236, 2006.

[15] S. E. Inzucchi and R. S. Sherwin. The prevention of type 2 diabetes mellitus. *Endocrinology and Metabolism Clinics*, 34(1):199–219, 2005.

[16] A. Dagliati, A. Malovini, P. Decata, G. Cogni, M. Teliti, L. Sacchi, C. Cerra, L. Chiovato, and R. Bellazzi. Hierarchical bayesian logistic regression to forecast metabolic control in type 2 dm patients. In *AMIA Annual*

*Symposium Proceedings*, volume 2016, page 470. American Medical Informatics Association, 2016.

[17] A. Dagliati, A. Marinoni, C. Cerra, P. Decata, L. Chiovato, P. Gamba, and R. Bellazzi. Integration of administrative, clinical, and environmental data to support the management of type 2 diabetes mellitus: From satellites to clinical care. *Journal of diabetes science and technology*, 10(1):19–26, 2016.

[18] F. L. Seixas, B. Zadrozny, J. Laks, A. Conci, and D. C. M. Saade. A bayesian network decision model for supporting the diagnosis of dementia, alzheimer s disease and mild cognitive impairment. Computers in biology and medicine, 51: 140–158, 2014.

[19] J. Pearl. Probabilistic reasoning in intelligent systems. 1988. *San Mateo, CA: Kaufmann*, 23:33–34.

[20] K. P. Murphy and S. Russell. Dynamic bayesian networks: representation, inference and learning. 2002.

[21] N. Trifonova, A. Kenny, D. Maxwell, D. Duplisea, J. Fernandes, and A. Tucker. Spatio-temporal bayesian network models with latent variables for revealing trophic dynamics and functional networks in fisheries ecology. *Ecological Informatics*, 30:142–158, 2015.

[22] Y. Guo, G. Bai, and Y. Hu. Using bayes network for prediction of type-2 diabetes. In *2012 International Conference for Internet Technology and Secured Transactions*, pages 471–472. IEEE, 2012.

[23] S. Marini, E. Trifoglio, N. Barbarini, F. Sambo, B. Di Camillo, A. Malovini, M. Manfrini, C. Cobelli, and R. Bellazzi. A dynamic bayesian network model for long-term simulation of clinical complications in type 1 diabetes. *Journal of biomedical informatics*, 57:369–376, 2015.

[24] A. Tucker, X. Liu, and D. Garway-Heath. Spatial operators for evolving dynamic bayesian networks from spatio-temporal data. In *Genetic and Evolutionary ComputationGECCO 2003*, pages 205–205. Springer, 2003.

[25] M. Grzegorczyk and D. Husmeier. Non-stationary continuous dynamic bayesian networks. In *Advances in Neural Information Processing Systems*, pages 682–690, 2009.

[26] J. W. Robinson and A. J. Hartemink. Learning non-stationary dynamic bayesian networks. *Journal of Machine Learning Research*, 11(Dec):3647–3680, 2010.

[27] M. Talih and N. Hengartner. Structural learning with time-varying components: tracking the cross-section of financial time series. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(3):321–341, 2005.

[28] A. Lloyd, W. Sawyer, and P. Hopkinson. Impact of long-term complications on quality of life in patients with type 2 diabetes not using insulin. *Value in Health*, 4(5):392–400, 2001.

[29] Q. Yang and X. Wu. 10 challenging problems in data mining research. *International Journal of Information Technology & Decision Making*, 5(04): 597–604, 2006.

[30] L. Litwak, S.-Y. Goh, Z. Hussein, R. Malek, V. Prusty, and M. E. Khamseh. Prevalence of diabetes complications in people with type 2 diabetes mellitus and its association with baseline characteristics in the multinational a 1 chieve study. *Diabetology and metabolic syndrome*, 5(1):57, 2013.

[31] N. Japkowicz and S. Stephen. The class imbalance problem: A systematic study. *Intelligent data analysis*, 6(5):429–449, 2002.

[32] L. Yousefi, L. Saachi, R. Bellazzi, L. Chiovato, and A. Tucker. Predicting comorbidities using resampling and dynamic bayesian networks with latent variables. In *Computer-Based Medical Systems (CBMS), 2017 IEEE 30th International Symposium on*, pages 205–206. IEEE, 2017.

[33] L. Yousefi, A. Tucker, M. Al-luhaybi, L. Saachi, R. Bellazzi, and L. Chiovato. Predicting disease complications using a stepwise hidden variable approach for learning dynamic bayesian networks. In *2018 IEEE 31st International Symposium on Computer-Based Medical Systems (CBMS)*, pages 106–111. IEEE, 2018.

[34] G. Liang. An effective method for imbalanced time series classification: Hybrid sampling. In *Australasian Joint Conference on Artificial Intelligence*, pages 374–385. Springer, 2013.

[35] X. Zhang, K. B. Korb, A. E. Nicholson, and S. Mascaro. Latent variable discovery using dependency patterns. *arXiv preprint arXiv: 1607.06617*, 2016.

[36] C. SPEARMAN. " general intelligence," objectively determined and measured. *American Journal of Psychology*, 15:201–293, 1904.

[37] J. Martin and K. VanLehn. Discrete factor analysis: Learning hidden variables in bayesian networks. Technical report, Technical report, Department of Computer Science, University of Pittsburgh, 1995.

[38] X. Boyen, N. Friedman, and D. Koller. Discovering the hidden structure of complex dynamic systems. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pages 91–100. Morgan Kaufmann Publishers Inc., 1999.

[39] L. co Todorovski, B. Cestnik, and M. Kline. Qualitative clustering of short time-series: A case study of firms reputation data. *IDDM-2002*, 141, 2002.

[40] F. Wittig. Learning bayesian networks with hidden variables for user modeling. In *UM99 User Modeling*, pages 343–344. Springer, 1999.

[41] R. Silva and Z. Ghahramani. The hidden life of latent variables: Bayesian learning with mixed graph models. *Journal of Machine Learning Research*, 10 (Jun):1187–1238, 2009.

[42] D. Chicharro and S. Panzeri. Algorithms of causal inference for the analysis of effective connectivity among brain regions. *Frontiers in neuroinformatics*, 8:64, 2014.

[43] H. Amirkhani, M. Rahmati, P. J. Lucas, and A. Hommersom. Exploiting experts knowledge for structure learning of bayesian networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(11):2154–2170, 2017.

[44] S. Russell, J. Binder, D. Koller, and K. Kanazawa. Local learning in probabilistic networks with hidden variables. In *IJCAI*, volume 95, pages 1146–1152, 1995.

[45] F. Rijmen, E. H. Ip, S. Rapp, and E. G. Shaw. Qualitative longitudinal analysis of symptoms in patients with primary and metastatic brain tumours. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 171(3): 739–753, 2008.

[46] N. Friedman, K. Murphy, and S. Russell. Learning the structure of dynamic probabilistic networks. In *Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence*, pages 139–147. Morgan Kaufmann Publishers Inc., 1998.

[47] G. Elidan, N. Lotner, N. Friedman, and D. Koller. Discovering hidden variables: A structure-based approach.

In *Advances in Neural Information Processing Systems*, pages 479–485, 2001.

[48] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.

[49] E. Choi, M. T. Bahadori, A. Schuetz, W. F. Stewart, and J. Sun. Doctor ai: Predicting clinical events via recurrent neural networks. In *Machine Learning for Healthcare Conference*, pages 301–318, 2016.

[50] M. Khashei, M. Bijari, and G. A. R. Ardali. Improvement of auto-regressive integrated moving average models using fuzzy logic and artificial neural networks (anns). *Neurocomputing*, 72(4–6):956–967, 2009.

[51] Y. Li, S. Swift, and A. Tucker. Modelling and analysing the dynamics of disease progression from cross-sectional studies. *Journal of biomedical informatics*, 46(2):266–274, 2013.

[52] L. M. Zintgraf, T. S. Cohen, and M. Welling. A new method to visualize deep neural networks. *arXiv preprint arXiv:1603.02518*, 2016.

[53] D. Pedreschi, F. Giannotti, R. Guidotti, A. Monreale, S. Ruggieri, and F. Turini. Meaningful explanations of black box ai decision systems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9780–9784, 2019.

[54] H. Lakkaraju, E. Kamar, R. Caruana, and J. Leskovec. Faithful and customizable explanations of black box models. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 131–138, 2019.

[55] S. Ceccon, D. F. Garway-Heath, D. P. Crabb, and A. Tucker. Exploring early glaucoma and the visual field test: Classification and clustering using bayesian networks. *IEEE journal of biomedical and health informatics*, 18(3): 1008–1014, 2014.

[56] P. J. Lucas, L. C. Van der Gaag, and A. Abu-Hanna. Bayesian networks in biomedicine and health-care, 2004.

[57] L. Peelen, N. F. de Keizer, E. de Jonge, R.-J. Bosman, A. Abu-Hanna, and N. Peek. Using hierarchical dynamic bayesian networks to investigate dynamics of organ failure in patients in the intensive care unit. *Journal of biomedical informatics*, 43(2):273–286, 2010.

[58] S. Aghabozorgi, A. S. Shirkhorshidi, and T. Y. Wah. Time-series clustering–a decade review. *Information Systems*, 53: 16–38, 2015.

[59] F. Altiparmak, H. Ferhatosmanoglu, S. Erdal, and D. C. Trost. Information mining over heterogeneous and high-dimensional time-series data in clinical trials databases. *IEEE Transactions on Information Technology in Biomedicine*, 10(2):254–263, 2006.

[60] S. Colagiuri. Glycated haemoglobin (hba1c) for the diagnosis of diabetes mellitus–practical implications. *Diabetes research and clinical practice*, 93(3):312, 2011.

[61] G. Das, K.-I. Lin, H. Mannila, G. Renganathan, and P. Smyth. Rule discovery from time series. In *KDD*, volume 98, pages 16–22, 1998.

[62] S. Ceccon, D. Garway-Heath, D. Crabb, and A. Tucker. The dynamic stage bayesian network: identifying and modelling key stages in a temporal process. *Advances in Intelligent Data Analysis X*, pages 101–112, 2011.

[63] R. Raman, A. Gupta, S. Krishna, V. Kulothungan, and T. Sharma. Prevalence and risk factors for diabetic microvascular complications in newly diagnosed type ii diabetes mellitus. sankara nethralaya diabetic retinopathy

epidemiology and molecular genetic study (sn-dreams, report 27). *Journal of Diabetes and its Complications*, 26(2): 123–128, 2012.

[64] K. R. Munana. Long-term complications of diabetes mellitus, part i: Retinopathy, nephropathy, neuropathy. *Veterinary Clinics: Small Animal Practice*, 25(3):715–730, 1995.

[65] W. Wang, J. Yang, and R. Muntz. Tar: Temporal association rules on evolving numerical attributes. In *Data Engineering, 2001. Proceedings. 17th International Conference on*, pages 283–292. IEEE, 2001.

[66] R. Agrawal, T. Imieliński, and A. Swami. Mining association rules between sets of items in large databases. In *Acm sigmod record*, volume 22, pages 207–216. ACM, 1993.

[67] C.-H. Lee, M.-S. Chen, and C.-R. Lin. Progressive partition miner: an efficient algorithm for mining general temporal association rules. *IEEE Transactions on Knowledge and Data Engineering*, (4):1004–1017, 2003.

[68] T. F. Gharib, H. Nassar, M. Taha, and A. Abraham. An efficient algorithm for incremental mining of temporal association rules. *Data & Knowledge Engineering*, 69(8):800–815, 2010.

[69] J. F. Allen et al. Towards a general theory of action and time. *Artificial intelligence*, 23(2):123–154, 1984.

[70] L. Sacchi, C. Larizza, C. Combi, and R. Bellazzi. Data mining with temporal abstractions: learning rules from time series. *Data Mining and Knowledge Discovery*, 15(2):217–247, 2007.

[71] J. M. Ale and G. H. Rossi. An approach to discovering temporal association rules. In *Proceedings of the 2000 ACM symposium on Applied computing-Volume 1*, pages 294–300, 2000.

[72] H. Jen-Wei and C. M. S. Dai Bi-Ru. Twain: Two-end association miner with precise frequent exhibition periods. *ACM Transactions on Knowledge Discovery from Data*, 8(2):800–815, 2007.

[73] Q. Zhao and S. S. Bhowmick. Association rule mining: A survey. *Nanyang Technological University, Singapore*, page 135, 2003.

[74] J. Mennis and J. W. Liu. Mining association rules in spatio-temporal data: An analysis of urban socioeconomic and land cover change. *Transactions in GIS*, 9(1):5–17, 2005.

[75] P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the right objective measure for association analysis. *Information Systems*, 29(4):293–313, 2004.

[76] P. J. Dyck, K. Kratz, J. Karnes, W. J. Litchy, R. Klein, J. Pach, D. Wilson, P. O'brien, and L. Melton. The prevalence by staged severity of various types of diabetic neuropathy, retinopathy, and nephropathy in a population-based cohort: the rochester diabetic neuropathy study. *Neurology*, 43(4):817–817, 1993.

[77] G. F. Cooper and E. Herskovits. A bayesian method for the induction of probabilistic networks from data. *Machine learning*, 9(4):309–347, 1992.

[78] A. Doddi, S. Marathe, D. C. Ravi, and S. Torney. Discovery of association rules in medical data. *Medical informatics and the Internet in medicine*, 26(1):25–33, 2001.

[79] C. Ordonez, C. A. Santana, and L. De Braal. Discovering interesting association rules in medical data. In *ACM SIGMOD workshop on research issues in data mining and knowledge discovery*, pages 78–85. Citeseer, 2000.

[80] M. Plasse, N. Niang, G. Saporta, A. Villeminot, and L. Leblond. Combined

use of association rules mining and clustering methods to find relevant links between binary rare attributes in a large data set. *Computational Statistics and Data Analysis*, 52(1):596–613, 2007.

[81] S. Mani and G. F. Cooper. Causal discovery using a bayesian local causal discovery algorithm. In *Medinfo*, pages 731–735, 2004.

[82] G. Sparacino, A. Facchinetti, A. Maran, and C. Cobelli. Continuous glucose monitoring time series and hypo/hyperglycemia prevention: requirements, methods, open problems. *Current diabetes reviews*, 4(3):181–192, 2008.

[83] L. Yousefi, S. Swift, M. Arzoky, L. Saachi, L. Chiovato, and A. Tucker. Opening the black box: Personalizing type 2 diabetes patients based on their latent phenotype and temporal associated complication rules. *Computational Intelligence*, 2020.

[84] L. Yousefi, S. Swift, M. Arzoky, L. Saachi, L. Chiovato, and A. Tucker. Opening the black box: Discovering and explaining hidden variables in type 2 diabetic patient modelling. In *2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1040–1044. IEEE, 2018.

[85] L. Yousefi, S. Swift, M. Arzoky, L. Sacchi, L. Chiovato, and A. Tucker. Opening the black box: Exploring temporal pattern of type 2 diabetes complications in patient clustering using association rules and hidden variable discovery. In *2019 IEEE 32nd International Symposium on Computer-Based Medical Systems (CBMS)*, pages 198–203. IEEE, 2019.

[86] U. P. D. S. Group et al. Uk prospective diabetes study 16: overview of 6 years' therapy of type ii diabetes: a progressive disease. *Diabetes*, 44(11): 1249–1258, 1995.

[87] A. Z. Ali, M. Hossain, R. Pugh, et al. Diabetes, obesity and hypertension in urban and rural people of bedouin origin in the united arab emirates. *The Journal of tropical medicine and hygiene*, 98(6): 407–415, 1995.

[88] S. Concaro, L. Sacchi, C. Cerra, P. Fratino, and R. Bellazzi. Mining health care administrative data with temporal association rules on hybrid events. *Methods of information in medicine*, 50 (02):166–179, 2011.

[89] K. G. Tolman, V. Fonseca, A. Dalpiaz, and M. H. Tan. Spectrum of liver disease in type 2 diabetes and management of patients with diabetes and liver disease. *Diabetes care*, 30(3): 734–743, 2007.

[90] R. Turner, H. Millns, H. Neil, I. Stratton, S. Manley, D. Matthews, and R. Holman. Risk factors for coronary artery disease in non-insulin dependent diabetes mellitus: United kingdom prospective diabetes study (ukpds: 23). *Bmj*, 316(7134):823–828, 1998.

[91] M. Cusick, A. D. Meleth, E. Agron, M. R. Fisher, G. F. Reed, G. L. Knatterud, F. B. Barton, M. D. Davis, F. L. Ferris, E. Y. Chew, et al. Associations of mortality and diabetes complications in patients with type 1 and type 2 diabetes: early treatment diabetic retinopathy study report no. 27. *Diabetes Care*, 28(3):617–625, 2005.

[92] P. Thuluvath and D. Triger. Autonomic neuropathy and chronic liver disease. *QJM: An International Journal of Medicine*, 72(2):737–747, 1989.

[93] C. Liu, W. Ding, Y. Hu, X. Xia, B. Zhang, J. Liu, and D. Doermann. Circulant binary convolutional networks for object recognition. *IEEE Journal of Selected Topics in Signal Processing*, 2020.

[94] A. Ramachandran, C. Snehalatha, K. Satyavani, E. Latha, R. Sasikala, and

V. Vijay. Prevalence of vascular complications and their risk factors in type 2 diabetes. *The Journal of the Association of Physicians of India*, 47(12): 1152–1156, 1999.

[95] R. Bellazzi, L. Sacchi, and S. Concaro. Methods and tools for mining multivariate temporal data in clinical and biomedical applications. In *2009 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 5629–5632. IEEE, 2009.

[96] R. Bellazzi and B. Zupan. Predictive data mining in clinical medicine: current issues and guidelines. *International journal of medical informatics*, 77(2): 81–97, 2008.

[97] L. Yousefi, M. Al-Luhaybi, L. Sacchi, L. Chiovato, and A. Tucker. Identifying latent variables in dynamic bayesian networks with bootstrapping applied to type 2 diabetes complication prediction. 2020.