

Forecasting Potential Diabetes Complications

Yang Yang

Tsinghua University
yangyang@keg.cs.tsinghua.edu.cn

Walter Luyten

Katholieke Universiteit Leuven
Walter.Luyten@med.kuleuven.be

Lu Liu

Northwestern University
liulu26@gmail.com

Marie-Francine Moens

Katholieke Universiteit Leuven
sien.moens@cs.kuleuven.be

Jie Tang

Tsinghua University
jietang@tsinghua.edu.cn

Juanzi Li

Tsinghua University
ljz@keg.cs.tsinghua.edu.cn

Abstract

Diabetes complications often afflict diabetes patients seriously: over 68% of diabetes-related mortality is caused by diabetes complications. In this paper, we study the problem of automatically diagnosing diabetes complications from patients' lab test results. The objective problem has two main challenges: 1) feature sparseness: a patient only undergoes 1.26% lab tests on average, and 65.5% types of lab tests are performed on samples from less than 10 patients; 2) knowledge skewness: it lacks comprehensive detailed domain knowledge of the association between diabetes complications and lab tests. To address these challenges, we propose a novel probabilistic model called Sparse Factor Graph Model (SparseFGM). SparseFGM projects sparse features onto a lower-dimensional latent space, which alleviates the problem of sparseness. SparseFGM is also able to capture the associations between complications and lab tests, which help handle the knowledge skewness. We evaluate the proposed model on a large collections of real medical records. SparseFGM outperforms (+20% by F1) baselines significantly and gives detailed associations between diabetes complications and lab tests.

Introduction

Diabetes mellitus, or simply diabetes, is a very common chronic disease, whose prevalence continues to increase, especially in the "Western" world. People with diabetes normally suffer from high blood glucose levels, which damages tissues over time and leads to life-threatening health complications, e.g., hypertension, coronary heart disease, hyperlipidemia, etc. Many of these complications seriously diminish the quality of life and even lead to premature death. Studies show that diabetes and its complications are major causes of early death in most countries (Roglic and Unwin 2010), and over 68% of diabetes-related mortality is caused by diabetes complications¹.

However, diabetes and its complications are often diagnosed late due to the long subclinical evolution, unclear pathogenesis, and poor medical facilities in some regions over the world. Although more than 471 billion USD were

spent on healthcare for 371 million diabetes patients worldwide in 2012, still half of the people with diabetes remain undiagnosed, and 4.8 million people died in 2012 due to diabetes².

Challenges. In this paper, we study the problem of automatically diagnosing diabetes complications from lab test results. More specifically, given a series of lab test results of a potential diabetes patient, the goal of this study is to find particular complications that the patient may have. This approach can benefit a wide range of applications and studies, such as online diabetes complication diagnosis system, and studies of underlying patterns between diabetes complications and lab tests. The problem has several unique challenges:

Feature sparseness. There are 1945 different lab tests in our data set, while the clinical record of each patient only contains 24.43 different lab tests (1.26%) on average, which means only 1.26% of the features on average are used to represent each instance. At the same time, 65.5% types of lab tests are recoded in less than 10 clinical records among 181,933 records in total (0.0054%). The rare appearance of features will cause ineffectiveness of the model learning and lead to poor performance.

Knowledge skewness. The performance of approaches to solve the objective problem highly relies on the domain knowledge of associations between diabetes complications and different types of lab tests. Previous work studied some of the associations. For example, (Zürbig et al. 2012) studied the association between tests on urinary proteome and diabetic nephropathy. Voulgari et al. (2010) studied the association between echocardiographic methods and diabetic cardiomyopathy. However, to the best of our knowledge, no study has provided a comprehensive detailed association. The skewness of the domain knowledge will cause an unbalanced performance.

Proposed Solution and Contributions. To address these challenges, we propose the Sparse Factor Graph Model (SparseFGM). For handling feature sparseness, SparseFGM projects the sparse features into a lower-dimensional latent space, which alleviates the sparseness issue. For handling knowledge skewness, SparseFGM models the associa-

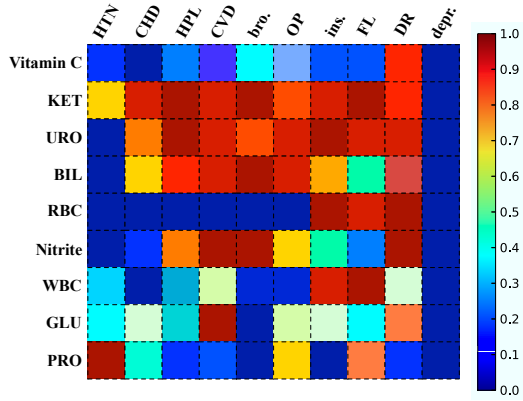


Figure 1: Associations between 10 complications and 9 parameters of routine urine analysis. The parameters include: vitamin C, ketone (KET), urobilinogen (URO), bilirubin (BIL), red blood cell (RBC), nitrite, white blood cell (WBC), glucose (GLU), and protein (PRO). The meanings of the complication abbreviation can be referred to the experimental section.

tions between all diabetes complications and lab tests in our dataset by exponential-linear functions. More importantly, with the model, we can not only forecast a potential diabetes complication, but also discover which types of lab test are most strongly associated with the diagnosed complication.

We evaluate the proposed model on a large collection of real medical records. Figure 1 shows the results obtained in our experiment to demonstrate the associations between ten complications and nine parameters of routine urine analysis. The color in each square represents the association strength, discovered by our proposed model, between each urine test and complication. For example, protein (PRO) has a strong association with hypertension (HTN) as chronic hypertension causes kidney damage, which in turn leads to the appearance of protein in the urine.

In all, our contributions of this paper are summarized in the following.

- We identify and formalize a new problem of diagnosing diabetes complications by a machine learning method. To the best of our knowledge, no previous work has extensively studied this problem.
- We propose a probabilistic model, SparseFGM, to solve the diabetes complication diagnosis problem, which integrates a feature dimensionality reduction process, association mining between complications and lab tests, and complication prediction into a unified model.
- We demonstrate the power of the proposed method using a large collection of real clinical records. Experimental results show that SparseFGM outperforms SVM, traditional factor graph models (FGM), and PCA based methods by +20% on average.

Problem Definition

In this section, we define several related concepts and formulate the diabetes complication forecasting problem.

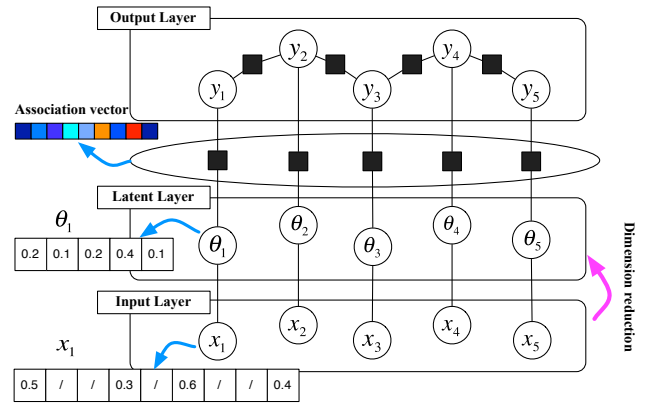


Figure 2: An illustration of the proposed model. Instances in the input layer are represented sparsely. The model projects a 9-dimensional feature space into a 5-dimensional latent space, which alleviates the sparseness. Vectors indicating associations between a particular issue and (groups of) lab tests are derived from the feature factors, which bridge the latent layer and output layer of SparseFGM. Detailed dependencies between input layer and latent layer are omitted.

We first introduce the observed lab tests, which are performed on patients to evaluate the patient’s health condition. Formally, we define a lab test record of a patient as follows:

Definition 1. Lab Test Record. We define a lab test record of a patient n as a set $\mathbf{u}_n = \{u_i\}$, where each u_i denotes a lab test performed on the patient n . We also use a tuple $u = (l, r)$ to represent each lab test, where l is the type of the lab test (e.g., WBC test in urine routine), and r is the result. Lab tests can have numerical or categorical types of results.

A patient can have several different lab test records as he may undergo lab tests more than once. For example, patients with serious diabetes may have to undergo tests once per week. Thus we further define a (lab test) *record sequence* for a particular patient to represent the order of the patient taking lab tests. Specifically, a record sequence of a patient is defined as a sequence $S = (s_i)$, where s_i is a lab test record. We say the test s_i is performed later than s_j if $i > j$.

We finally define the objective problem in this study. Generally, our goal is to diagnose diabetes complications from a set of lab test records of patients. Specifically, the problem can be defined as follows.

Definition 2. Diabetes Complication Forecasting. The inputs of diabetes complication forecasting include a diabetes complication C , a set of patients $V = \{v\}$, and a record sequence S_n for each patient v_n . The goal of diabetes forecasting is, for each patient v_n , and each lab test record $\mathbf{u}_n \in S_n$, determining if the patient v_n has the diabetes complication C , when the lab tests in \mathbf{u}_n are performed on the patient v_n .

Table 1: Notations in the proposed model.

SYMBOL	DESCRIPTION
K	The number of latent variables;
L	The number of different lab tests;
C	A particular diabetes complication the model aims to forecast;
x_{nl}	The l -th value in the n -th instance node;
y_n	A label node indicates whether the patient corresponding to the n -th instance has the diabetes complication C or not;
θ_n	The parameter of the multinomial distribution over latent variables specific to the n -th instance node;
μ_{kl}, δ_k	The parameters of the Gaussian distribution used to sample x_{nl} , which has a numerical value and is specific to the latent variable k ;
$\phi_{klx_{nl}}$	The parameter of the multinomial distribution over x_{nl} , which has a categorical value, specific to the latent variable k ;
α, β	Parameters used to define the feature factor and the correlation factor respectively.

Model Framework

In this section, we explain the model that we have developed for the diabetes complication forecasting problem. Table 1 summarizes the notations used in the proposed model.

For each patient’s lab test record \mathbf{u}_n , we create an instance node \mathbf{x}_n . Assume that there are in total L available (different) lab tests, we define \mathbf{x}_n as a L -dimensional vector. For each lab test $u_{ni} = (l, r) \in \mathbf{u}_n$, we set $x_{nl} = r$. We also set $x_{nl} = /$ to denote that there is no lab test with type l in \mathbf{u}_n . We associate a variable label y_n for \mathbf{x}_n ’s corresponding patient, to denote whether the patient has the given complication C or not. Particularly, we use $y_n = 1$ to denote a positive result, and use $y_n = -1$ to denote a negative result. Based on this formulation, we can build a classification model to map the input instance to the target label, i.e., $f(\mathbf{x}_n) \rightarrow y_n$. With this formulation, more-or-less standard technologies can be employed for learning and inference, for example SVM (Chang and Lin 2011). However, SVM treats all instances independently and cannot capture the dependencies between instances (e.g., the dependencies between records in the record sequence). A factor graph model (Kschischang, Frey, and Loeliger 2001) can be leveraged to model the dependencies. However, the traditional factor graph model still cannot deal with the challenges in our problem: feature sparseness and knowledge skewness.

To this end, we propose a Sparse Factor Graph Model (SparseFGM). The basic idea here is that we first project instances from the original sparse space onto a lower dimensional latent space. The latent space is used to capture the correlation between different lab tests. In this way, the model alleviates the feature sparseness problem. Moreover, the graphical structure of the SparseFGM model is designed to model the correlation between different labels. Figure 2 shows a simple example of a SparseFGM. Instances in the input layer (represented as a 9-dimensional vector) are sparse. The model projects all instances onto a 5-dimensional latent space, which alleviates the sparseness

issue. Furthermore, the factor function (indicated as black rectangles) between labels is defined to model the dependencies between instances. Please notice that SparseFGM can be applied not only in the diabetes complication forecasting problem but also to other similar scenarios.

According to the model, we can write the joint distribution of a given set of instances \mathbf{X} over Y as

$$P(y_n | \theta_n, \mathbf{x}_n) = P(y_n | \theta_n) \prod_l \left(\sum_{k=1}^K \theta_{nk} \cdot \Omega_{x_{nl}k} \right) \quad (1)$$

where l is an index of \mathbf{x}_n which satisfies $x_{nl} \neq /$; K is the number of latent variables in SparseFGM; θ_n is the parameter of a multinomial distribution which assigns latent variables to the n -th instance; $\Omega_{x_{nl}k}$ is the parameter of a distribution which generates values of x_{nl} when it is assigned with latent variable k . In the diabetes complication diagnosis problem, there are two kinds of lab test results: numerical and categorical. Similar with Liu et al. (2013), we assume that the numerical values are drawn from some Gaussian distributions, and categorical values are drawn from multinomial distributions. Thus $\Omega_{x_{nl}k}$ is defined as

$$\Omega_{x_{nl}k} = \begin{cases} N(x_{nl} | \mu_{kl}, \delta_k) & x_{nl} \text{ is numerical} \\ \phi_{klx_{nl}} & x_{nl} \text{ is categorical} \end{cases} \quad (2)$$

where μ_{kl} is the mean of the Gaussian distribution w.r.t. latent variable k and the category of patient symptom l ; we assume that all Gaussian distributions corresponding to the same latent variable k share the same standard deviation δ_k ; $\phi_{klx_{nl}}$ is the probability that the l -th dimension of an instance variable assigned to the latent variable k has the value x_{nl} .

Regarding the feature factor $P(y_n | \theta_n)$, we define it as

$$P(y_n | \theta_n) = \frac{1}{Z_1} \exp\{\alpha \cdot f(\theta_n, y_n)\} \quad (3)$$

where α is a parameter of SparseFGM, Z_1 is a normalization factor, and $f(\theta_n, y_n) = y_n \cdot \theta_n$.

To model the correlations between different labels, we define the correlation factor. Intuitively, a patient’s health condition will not change quickly. Thus, we create a correlation factor between two labels y_n and $y_{n'}$ corresponding to two adjacent lab test records in the same patient’s *record sequence*. Specifically, we define the correlation factor as

$$P(y_n, y_{n'}) = \frac{1}{Z_2} \exp\{\beta \cdot g(y_n, y_{n'})\} \quad (4)$$

where β is a real number, Z_2 is a normalization factor, and $g(y_n, y_{n'})$ is defined as a vector of indicator functions.

By combining all the factors together, we can obtain the log-likelihood objective function for SparseFGM:

$$\begin{aligned} O(\lambda) &= \sum_n \log P(y_n | \theta_n, \mathbf{x}_n) + \sum_c \log P(y_{c_1}, y_{c_2}) \\ &= \sum_n \alpha f(\theta_n, y_n) + \sum_c \beta g(y_{c_1}, y_{c_2}) \\ &\quad + \sum_n \sum_l \log \sum_k \theta_{nk} \Omega_{k,l,x_{nl}} - \log Z \end{aligned} \quad (5)$$

Input: a feature matrix X , learning rate η
Output: estimated parameters λ
Initialize $\alpha, \beta, \theta, \mu, \phi$ randomly;
Initialize $\delta \leftarrow \mathbf{1}$;
repeat
 Calculate $P(k_{nl}|\mathbf{x}_n, \lambda_n)$ according to Eq. 8;
 Update $\theta, \mu, \delta, \phi$ according to Eq. 9-12;
 Call LBP to calculate $E[\sum_n f(\theta_n, y_n)]$ and
 $E[\sum_c g(y_{c_1}, y_{c_2})]$;
 Call LBP to calculate $E_{P_\alpha(\mathbf{y}|\theta)}[\sum_n f(\theta_n, y_n)]$ and
 $E_{P_\beta(\mathbf{y})}[\sum_c g(y_{c_1}, y_{c_2})]$;
 Calculate $\frac{\partial O(\alpha, \beta)}{\partial \alpha}$ and $\frac{\partial O(\alpha, \beta)}{\partial \beta}$ according to Eq. 14;
 $\alpha_{\text{new}} = \alpha_{\text{old}} + \eta \frac{\partial O(\alpha, \beta)}{\partial \alpha}$;
 $\beta_{\text{new}} = \beta_{\text{old}} + \eta \frac{\partial O(\alpha, \beta)}{\partial \beta}$;
until Convergence ;

Algorithm 1: Learning algorithm of SparseFGM.

where c is an indicator for 2-cliques among label nodes \mathbf{Y} , and $Z = Z_1 Z_2$.

Model Learning

We introduce the learning algorithm to estimate the optimal parameter configuration $\lambda = \{\alpha, \beta, \theta, \phi, \mu, \delta\}$ for the proposed SparseFGM, i.e., finding parameters to maximize the log-likelihood, i.e.,

$$\begin{aligned} \lambda^* &= \arg \max_{\lambda} O(\lambda) \\ \text{s.t. } \sum_k \theta_{nk} &= 1, \sum_x \phi_{klx} = 1 \end{aligned} \quad (6)$$

Unfortunately Eq. 6 does not have a closed-form solution. Thus, we propose an EM-like approximation learning algorithm to estimate the parameters.

We first introduce the general idea of the learning algorithm. By regarding labels inferred by the last updated model configuration, α and β as fixed, the remaining part of the model can be learned as a mixture generative model with parameters θ and Ω . On the other hand, by fixing θ and Ω , the remaining part of the model is similar to FGM and we are able to utilize a gradient based method to estimate α and β . Next, we introduce how we update these parameters in detail, and give a framework of the learning algorithm.

Step 1: Update $\{\theta, \phi, \mu, \delta\}$. When fixing $Y, \{\alpha, \beta\}$, SparseFGM can be regarded as a mixture generative model with the generative process as follows: first, each instance variable has a latent variable distribution parameterized by θ_n . Next, for each instance variable x_n , a latent variable k is sampled according to $k \sim \text{Mult}(\theta_n)$. After that, x_{nl} is sampled according to $x_{nl} \sim \Omega_{kl}$. Finally, each label variable y_n is sampled according to $y_n \sim P(y_n|\theta_n)$.

With the generative process above, we can define the mixture generative model's log-likelihood, which can also be treated as the combination of related terms from Eq. 5:

$$\begin{aligned} O(\theta, \phi, \mu, \delta) &= \sum_n \sum_{l_1} \log \sum_k \theta_{nk} \frac{\exp\left\{-\frac{(x_{nl_1} - \mu_{kl})^2}{2\delta_k^2}\right\}}{\delta_k \sqrt{2\pi}} \\ &+ \sum_n \sum_{l_2} \log \sum_k \theta_{nk} \phi_{kl_2} x_{nl_2} \end{aligned} \quad (7)$$

By Jensen's inequality, we obtain the lower bound of Eq. 7. We derive the lower bound with respect to each parameter and set them to zero, we have

$$P(k_{nl}|\mathbf{x}_n, \lambda_n) = \frac{\theta_{nk} \Omega_{klr}}{\sum_k \theta_{nk} \Omega_{klr}} \quad (8)$$

$$\theta_{nk} = \frac{\sum_l P(k_{nl}|\mathbf{x}_n, \lambda_n)}{\sum_l \sum_{k_{nl}} P(k_{nl}|\mathbf{x}_n, \lambda_n)} + \alpha_k y_n \quad (9)$$

$$\mu_{kl} = \frac{\sum_n P(k_{nl}|\mathbf{x}_n, \lambda_n) x_{nl}}{\sum_n P(k_{nl}|\mathbf{x}_n, \lambda_n)} \quad (10)$$

$$\delta_k^2 = \frac{\sum_n \sum_l (x_{nl} - \mu_{kl})^2}{N \times L_1} \quad (11)$$

$$\phi_{klr} = \frac{\sum_n P(k_{nl}|\mathbf{x}_n, \lambda_n)}{\sum_n \sum_r P(k_{nl}|\mathbf{x}_n, \lambda_n)} \quad (12)$$

where L_1 is the number of numerical features, and N is the total number of lab test records.

Step 2: Update $\{\alpha, \beta\}$. When fixing θ and Ω , the remaining part of SparseFGM is similar with FGM, which has the log-likelihood as

$$O(\alpha, \beta) = \sum_n \alpha f(\theta_n, y_n) + \sum_c \beta g(y_{c_1}, y_{c_2}) - \log Z \quad (13)$$

We then adopt a gradient descent method to optimize Eq. 13. The gradient for each parameter is calculated as

$$\begin{aligned} \frac{\partial O(\alpha, \beta)}{\partial \alpha} &= E[\sum_n f(\theta_n, y_n)] - E_{P_\alpha(\mathbf{y}|\theta)}[\sum_n f(\theta_n, y_n)] \\ \frac{\partial O(\alpha, \beta)}{\partial \beta} &= E[\sum_c g(y_{c_1}, y_{c_2})] - E_{P_\beta(\mathbf{y})}[\sum_c g(y_{c_1}, y_{c_2})] \end{aligned} \quad (14)$$

We adopt Loopy Belief Propagation (LBP) (Murphy, Weiss, and Jordan 1999) to compute the marginal probability of \mathbf{y} and thus compute the two expectations. We then update α and β with a learning rate η with the gradient. Algorithm 1 summarizes the learning algorithm for SparseFGM.

Experiments

In this section, we present experimental results to demonstrate the effectiveness of the proposed approach. All codes used in the paper are publicly available³.

Experimental Setup

We use a collection of real medical records from a famous geriatric hospital. The data set spans one year, containing 181,933 medical records corresponding to 35,525 unique patients and 1,945 kinds of lab tests in total. On average each

³<http://arnetminer.org/diabetes>

clinical record contains 24.43 different lab tests (1.26% of all lab tests), which indicates that the feature sparsity problem is serious for this data set.

We view each medical record as an instance, and aim to infer whether the corresponding patient has a particular diabetes complication or not from the lab test results. Nine complications are taken into account in our experiments: hypertension (HTN), coronary heart disease (CHD), hyperlipidemia (HPL), cerebrovascular disease (CVD), bronchitis (bro.), osteoporosis (OP), insomnia (ins.), fatty liver (FL), diabetic retinopathy (DR), and depression (depr.). In the experiments, we randomly picked 60% of the medical records as training data and the rest for testing.

Evaluation Aspects. We evaluate our method on the following two aspects:

- **Forecasting Performance.** We evaluate the proposed model in terms of Precision, Recall, F1-Measure, and compare with baseline methods to validate its effectiveness.
- **Association pattern illustration.** We use the discovered association patterns between diabetes complications and lab tests as anecdotal evidence to further demonstrate the effectiveness of our method.

We compare the following methods for forecasting diabetes complications.

SVM. Lab test results are treated as features and LIB-SVM (Chang and Lin 2011) is employed as the classification model for complication forecasting.

FGM. A traditional factor graph is used as the classification model. We employ a gradient descent algorithm to learn the parameters in FGM (Tang, Zhuang, and Tang 2011; Hopcroft, Lou, and Tang 2011), and set the learning rate parameter as 0.1.

FGM+PCA. To solve the feature sparsity problem in FGM, PCA (Jolliffe 1986) is employed to convert the features into a set of principal components, which is then used as the input of FGM then.

SparseFGM. SparseFGM is our proposed model. In all experiments, we empirically set the number of latent variables in SparseFGM to 100, and set $\eta = 0.1$.

All algorithms were implemented in C++, and all experiments were performed on a Mac running Mac OS X with Intel Core i7 2.66 GHz and 4 GB of memory.

Forecasting Performance

Table 2 shows the performance of the different methods on the diabetes complication forecasting task. Due to the space limitations, we only present the results for the three most common diabetes complications: hypertension (HTN), coronary heart disease (CHD), and hyperlipidemia (HPL). From the table, we can see that SparseFGM achieves the highest F1-score in all three tasks compared with other methods on average. Generally, SparseFGM offers the possibility to model the dependencies and constraints together with modeling latent class information. The results show that recall is increased without hurting precision substantially, or in some cases while even improving precision.

Table 2: Performance of diabetes complication forecasting.

Complication	Method	Precision	Recall	F1
HTN	SVM	0.3804	0.4789	0.4241
	FGM	0.5666	0.4959	0.5075
	FGM+PCA	0.5741	0.3284	0.4178
	SparseFGM	0.4714	0.6319	0.5400
CHD	SVM	0.2132	0.0636	0.0980
	FGM	0.6264	0.1369	0.2247
	FGM+PCA	0.2425	0.8367	0.3761
	SparseFGM	0.2522	0.7972	0.3832
HPL	SVM	0.2208	0.0460	0.0761
	FGM	0.6557	0.0591	0.1084
	FGM+PCA	0.2047	0.8035	0.3262
	SparseFGM	0.2796	0.8396	0.4195

SVM and traditional factor-graph-based methods suffer from the feature sparsity particularly in terms of recall, which is on average 59.9% lower compared with SparseFGM in all tasks. By a careful investigation, we find that most parameters in FGM tend to be zero due to the serious feature sparsity problem, which causes the false diagnose results.

With the effect of dimension reduction, FGM+PCA-based methods improve the performance a lot. For example, FGM+PCA increases 40.3% in terms of recall compared with FGM in the CHD the forecasting task. However, FGM+PCA separates the sparse coding and classification process into two steps, while SparseFGM integrates them into a uniform model. Thus SparseFGM estimates the parameters better and outperforms FGM+PCA on average by 13.5% in terms of F1 on average.

Factor Contribution Analysis. Figure 3 demonstrates how the correlation factor helps in this problem. Generally, models considering the correlations between labels increase 16.29% in terms of F1 compared with their counterparts that do not consider the correlation factors. Intuitively, correlation factors improve the performance by bringing the prior knowledge that “a patient’s health condition will be similar in adjacent time periods”.

Association pattern illustration

Based on SparseFGM, we can discover the underlying patterns between diabetes complications and lab tests at the micro-level and macro-level.

Micro-level. At the micro-level, we calculate the association score $AS(c, e)$ of a complication c with a lab test type e as

$$AS(c, e) = \sum_k \alpha_k^c \theta_{ek} \quad (15)$$

where k is the latent variable in SparseFGM, and α^c is the feature factor parameter of SparseFGM for diagnosing complication c .

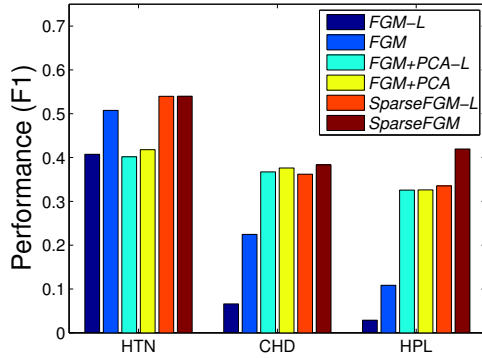


Figure 3: Correlation factor analysis. FGM-L, FGM+PCA-L, and SparseFGM-L denote the methods not considering the correlation factor.

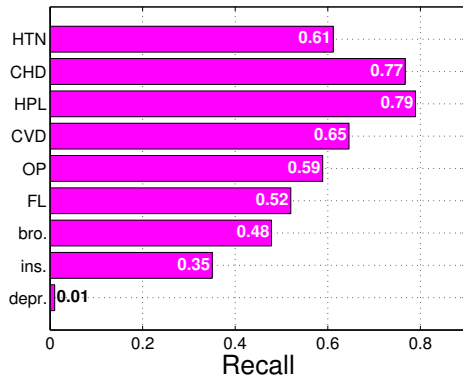


Figure 4: Association analysis at the macro-level. The X-axis denotes the performance of the proposed model for forecasting each diabetes complication in terms of recall.

The example which illustrates the associations between the lab tests of routine urine analysis and 10 diabetes complications is shown in Figure 1. We see that urinary glucose (GLU) is strongly associated with diabetic retinopathy (DR), because DR typically occurs in diabetic patients whose blood glucose levels are not well controlled. To our surprise, we find that insomnia (ins.) is associated with most urine routine measures, such as white blood cells (WBC). This association has an interesting explanation: WBC in the urine are typically found in urinary tract infections which cause frequent voiding, so people who have to go to the toilet frequently can hardly get a good sleep at night.

Macro-level. At the macro-level, we study the different strengths of diabetes complications forecasting depending on lab test results. We say a complication is more diagnosable from lab test results if the proposed model is able to target more positive instances correctly among all positive instances. Thus the strength can be estimated by the measurement recall. Figure 4 shows the results, from which we see that hyperlipidemia (HPL) can be diagnosed based on lab tests precisely, while depression (depr.) is usually rec-

ognized from psychological investigation instead of physiological lab tests such as blood tests. In this experiment, to avoid the effect of the label unbalance issue on the analysis results, we set the ratio of positive to negative instances for each complication as 1 : 5. Thus the recall results are different from the comparison experimental results, in which all the testing samples in the dataset are used.

Related Work

In recent years, utilizing health care data to study diabetes, which is a common chronic disease, has attracted the interest of several researchers. Liu et al. (2013) propose a Symptom-Diagnosis-Treatment model to mine the diabetes complication patterns and symptoms from electronic medical records, which is the same data set employed in this paper. Besides, Neuvirth et al. (2011) studied the personalized care management of diabetes patients at risk. For the medical work, De Groot et al. (2001) examine the strength and consistency of the relationship between depression and diabetes complications in studies of adult patients with type 1 or type 2 diabetes. Katon et al. (2003) and Kanaya et al. (Kanaya, Grady, and Barrett-Connor 2002) investigate two major diabetes complications: depression and coronary heart disease respectively, which are used as sample labels in our experiments. However, to the best of our knowledge, none of the previous work investigated how to predict diabetes complication by machine learning methods.

Sparse modeling is a component in many machine learning tasks. Yoshida et al. (2010) also propose a class of sparse latent factor models and relevant learning algorithms. Krishnapuram et al. (2005) and Zhong et al. (2008) learn models for a sparse Multinomial Logistic Regression and a sparse Conditional Random Field respectively with a Laplacian prior. Tan et al. (2012), Zhang et al. (2011), and Lee et al. (2006) study the sparse coding problem. Mairal et al. (2010) propose a stochastic online learning algorithm for learning dictionaries adapted to sparse coding tasks.

Conclusion

In this paper, we study the problem of diabetes complication forecasting. We propose a novel probabilistic model, Sparse Factor Graph Model (SparseFGM), which integrates sparse modeling and complication diagnosing into a uniform model. By using this model, we cannot only forecast diabetes complications but also discover the underlying association between complications and lab test types. We validate the model on a large collection of real medical records. Experimental results demonstrate that the proposed model outperforms baseline methods significantly. We also demonstrate how to discover the association between complications and lab test types at two granularities.

Acknowledgements. The work is supported by National High-tech R&D Program (No. 2014AA015103), National Basic Research Program of China (No. 2014CB340500, No. 2012CB316006), Natural Science Foundation of China (No. 61222212, 61103065, 61035004), NSFC-ANR (No. 61261130588), and a co-funding of Tsinghua and KU Leuven.

References

- Chang, C.-C., and Lin, C.-J. 2011. Libsvm: a library for support vector machines. *TIST* 2(3):27.
- De Groot, M.; Anderson, R.; Freedland, K. E.; Clouse, R. E.; and Lustman, P. J. 2001. Association of depression and diabetes complications: a meta-analysis. *Psychosom Med* 63(4):619–630.
- Hopcroft, J.; Lou, T.; and Tang, J. 2011. Who will follow you back? reciprocal relationship prediction. In *CIKM'11*, 1137–1146. ACM.
- Jolliffe, I. 1986. Principal component analysis. *Springer Series in Statistics, Berlin: Springer, 1986* 1.
- Kanaya, A. M.; Grady, D.; and Barrett-Connor, E. 2002. Explaining the sex difference in coronary heart disease mortality among patients with type 2 diabetes mellitus: a meta-analysis. *Archives of Internal Medicine* 162(15):1737.
- Katon, W.; Von Korff, M.; Lin, E.; Simon, G.; Ludman, E.; Bush, T.; Walker, E.; Ciechanowski, P.; and Rutter, C. 2003. Improving primary care treatment of depression among patients with diabetes mellitus: the design of the pathways study. *Gen Hosp Psychiatry* 25(3):158–168.
- Krishnapuram, B.; Carin, L.; Figueiredo, M. A.; and Hartemink, A. J. 2005. Sparse multinomial logistic regression: Fast algorithms and generalization bounds. *TPAMI* 27(6):957–968.
- Kschischang, F. R.; Frey, B. J.; and Loeliger, H. A. 2001. Factor graphs and the sum-product algorithm. *IEEE TOIT* 47:498–519.
- Lee, H.; Battle, A.; Raina, R.; and Ng, A. 2006. Efficient sparse coding algorithms. In *NIPS'06*, 801–808.
- Liu, L.; Tang, J.; Cheng, Y.; Agrawal, A.; Liao, W.-k.; and Choudhary, A. 2013. Mining diabetes complication and treatment patterns for clinical decision support. In *CIKM'13*, 279–288. ACM.
- Mairal, J.; Bach, F.; Ponce, J.; and Sapiro, G. 2010. Online learning for matrix factorization and sparse coding. *JMLR* 11:19–60.
- Murphy, K. P.; Weiss, Y.; and Jordan, M. I. 1999. Loopy belief propagation for approximate inference: An empirical study. In *UAI'99*, 467–475.
- Neuvirth, H.; Ozery-Flato, M.; Hu, J.; Laserson, J.; Kohn, M. S.; Ebadollahi, S.; and Rosen-Zvi, M. 2011. Toward personalized care management of patients at risk: the diabetes case study. In *SIGKDD'11*, 395–403. ACM.
- Roglic, G., and Unwin, N. 2010. Mortality attributable to diabetes: estimates for the year 2010. *Diabetes Res Clin PR* 87(1):15–19.
- Tan, M.; Tsang, I. W.; Wang, L.; and Zhang, X. 2012. Convex matching pursuit for large-scale sparse coding and subset selection. In *AAAI'12*.
- Tang, W.; Zhuang, H.; and Tang, J. 2011. Learning to infer social ties in large networks. In *ECML/PKDD'11*, 381–397.
- Voulgari, C.; Papadogiannis, D.; and Tentolouris, N. 2010. Diabetic cardiomyopathy: from the pathophysiology of the cardiac myocytes to current diagnosis and management strategies. *Vasc Health Risk Manag* 6:883.
- Yoshida, R., and West, M. 2010. Bayesian learning in sparse graphical factor models via variational mean-field annealing. *JMLR* 99:1771–1798.
- Zhang, X.; Yu, Y.; White, M.; Huang, R.; and Schuurmans, D. 2011. Convex sparse coding, subspace learning, and semi-supervised extensions. In *AAAI'11*.
- Zhong, P., and Wang, R. 2008. Learning sparse CRFs for feature selection and classification of hyperspectral imagery. *IEEE Trans. Geosci. Remote Sens.* 46(12):4186–4197.
- Zürbig, P.; Jerums, G.; Hovind, P.; MacIsaac, R. J.; Mischak, H.; Nielsen, S. E.; Panagiotopoulos, S.; Persson, F.; and Rossing, P. 2012. Urinary proteomics for early diagnosis in diabetic nephropathy. *Diabetes* 61(12):3304–3313.