

Interpretability on clinical analysis from Pattern Disentanglement Insight

Anonymous ACL submission

Abstract

Diagnosis of a clinical condition can help medical professionals save time in clinical decision-making and prevent overlooking risks. Therefore we explore the problem of clinical text interpretability using free-text medical notes recorded in electronic health records (EHR). MIMIC-III is a de-identified EHR database containing observations from over 40,000 patients in critical care units. Since text corpus is unstructured and in non-database table format, existing machine learning models may have ineffective interpretability; however, interpretability is often desirable for clinical diagnosis. Hence, in this paper, we propose a text mining and pattern discovery solution to discover strong association patterns from patient discharge summaries and the code of international classification of diseases (ICD9 code). The proposed approach offers a straightforward interpretation of the underlying relation of patient characteristics in an unsupervised machine learning setting. The clustering results outperform the baseline clustering algorithm and are comparable to baseline supervised methods.

1 Introduction

If Machine Learning (ML) is to play a significant role in supporting clinical decision making, then it is essential to gain clinician trust (Kim, 2021). Interpretability is frequently defined as the degree to which a human can understand the cause and reason of ML model decisions. The higher the interpretability of a model implies the better the comprehension and explanation of the problem, leading to more accurate and reliable predictions. Most ML algorithms today concentrate on prediction power using general-purpose learning algorithms on large and complex data.

However, even though some ML models can also provide various degrees of interpretability, they generally sacrifice interpretability for predictive power (Ghannam and Techtmann, 2021). Therefore, in this study, we focus on interpreting the

diagnostic characteristics/patterns from the electronic health records (EHR).

An EHR is a digital collection of medical information about a person, which includes information about a patient’s health history, such as diagnoses, medicines, tests, allergies, immunizations, and treatment plans. The MIMIC-III (Medical Information Mart of Intensive Care) is an openly available extensive database comprising de-identified information relating to patients admitted to critical care units at a large tertiary care hospital (Johnson et al., 2016). Data primarily stores both structured (e.g. MIMIC-III medications, laboratory results are stored in the table with columns as features and rows as records) and unstructured data (e.g. MIMIC-III clinical notes, discharge summaries are stored in the format of free text). The patients’ information (e.g., discharge summary) is highly unstructured, thus making interpreting it a challenge.

To address the issue of ML Interpretability, we explore a two-step algorithm, combining text mining and pattern discovery, to discover strong association patterns from patient profiles and discharge summaries to reveal their relationships with the diagnosed disease (ICD9 code, which is the code of international classification of diseases). The first step is transforming free text into a structured dataset formatting as a table with columns as features and rows as records. The second step is discovering patterns and grouping patients’ records based on patterns in an unsupervised manner. The output of the proposed system is an interpretable Knowledge Base, which can link the pattern groups, discovered characteristics of records, and patients’ records together to shows “what” (disease), “who/where” (tracking patient records back) and “why” (discovered patterns) to interpret clinical notes for better clinical decision making.

The contributions of the paper are three-fold: 1) combining NLP and pattern discovery algorithm to interpret free-text clinical notes; 2) Grouping

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

records based on the discovered associations revealing characteristics of records via unsupervised learning ; 3) Generating an all-in-one knowledge base to link knowledge, pattern, and records together for interpretability.

To evaluate the performance of the proposed algorithm, we present both a knowledge base with discovered patterns and clustering results. To verify the effectiveness of discovered patterns, we interpret patterns from a clinical perspective to discuss the interpretability of output. As for the clustering results algorithm, although the process of clustering records does not require class label information, the results can be evaluated by balanced accuracy and weighted F1- score using the presumed class labels (ICD9 code) as ground truth.

2 Related Work

2.1 Clinical Data Analysis with Interpretability

Due to the complex nature of clinical language, clinical texts were hard to interpret. Most of the previous works on clinical data analysis were based on structured data, which lack complementary information such as lab reports or patient history. Clinical expert judgments may thus require information that are available only in unstructured data (e.g. clinical texts) (Culliton et al., 2017).

Latent Dirichlet Allocation (LDA) topic modeling (Blei et al., 2003) has been applied to the unstructured notes of EHRs to predict clinical outcomes (Bright et al., 2021; Huang et al., 2015; Wang et al., 2020). In addition, Ghassemi et al. (2014) showed latent topic features were more predictive than structured features, and a combination of the two performs best.

Topic features cluster terms into a small set of semantically related groups, which is proved useful in text classification and categorizing clinical reports (Chen et al., 2019; Pavlinek and Podgorelec, 2017; Kayi et al., 2013). For example, Horng et al. (2017) combined structured and unstructured data for sepsis prediction using text modeling involving topic models. Further, Gangavarapu et al. (2020) proposed a vector space and topic modeling-based approach applied to structure the raw clinical data by exploiting the data in the nursing notes. Hence, in this study, we use topic modeling to transform free text into a table with features and records.

In addition, with the recent development in neural networks, variants of pre-trained BERT (Devlin

et al., 2018) have widely been applied to clinical domains (e.g. BioBERT (Lee et al., 2020), ClinicalBERT (Alsentzer et al., 2019)). In addition, Feng et al. (2020) used pre-trained BERT-based models as static feature extractors and showed that variants of BERT performed better with Sepsis than Mortality prediction tasks however Wallace et al. (2019) showed that BERT fails to interpret life-threatening important numerical values such as body temperature in the clinical text. In our study, we discretize numerical values into discrete values to make the proposed algorithm can handle mixed-mode dataset. Further, Van Aken et al. (2021) showed that medical-specific negations can be misinterpreted by the pre-trained language models such as BERT (e.g. "abstinence from alcohol" becomes "alcohol dependence syndrome").

2.2 Pattern Discovery

To tackle the interpretability of clinical data analysis, many machine learning algorithms were proposed. For example, the Decision Tree can generate a rule set between features and class labels for interpretable prediction, but the rules need to be trained relying on labeled classes. In addition, Frequent Pattern Mining (Naulaerts et al., 2015) (Han et al., 2007) can discover knowledge in the form of association rules from relational data (Han et al., 2007) (Van Aken et al., 2021) but a manually threshold need to be set for calculated likelihood, support or confidence (Van Aken et al., 2021). And the discovered patterns may be overwhelmed (Wong and Li, 2008) with overlapping/redundant patterns, which requires some post analysis approaches, such pattern pruning and pattern summarizing (Wong and Li, 2008).

Hence, in this study, we applied the pattern discovery and disentanglement (PDD) algorithm, proposed in our previous research work, to discover simple patterns with statistical support to reveal the association between extracted features with class labels without further pattern pruning or pattern summarization. The output patterns are well organized, more clear and easier to be comprehended in a knowledge base.

3 Material:MIMIC-III Data Description

MIMIC-III is a de-identified relational clinical database containing observations from over 40,000 patients in critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012

(Johnson et al., 2016). While MIMIC-III consists of several tabular and time-series datasets, our present study utilizes clinical notes, found in the NOTEEVENTS table, and diagnoses, found in the DIAGNOSES_ICD table.

The former table, NOTEEVENTS, can provide us with the medical notes as text for a detailed description of medical center visits for each patient. The clinical notes contain an internal semi-structured format, which are subdivided into several components, such as: chief complaint, medical history, social history, and discharge information. Each observation refers to a unique hospital stay. The data are related to other tables through unique patient identifiers, hospital stay identifiers, and caregiver identifiers. The latter table, DIAGNOSES_ICD, can provide us with the diagnosis of each patient based on ICD9 codes, which are used as labels to be predicted, and linked with clinical notes.

In summary, our final data contains 11,537 rows/records with the top four classes/diseases represented by ICD9 code. The ICD9 codes are defined as follows: 414 - chronic ischemic heart disease, 038 - septicemia, 410 - acute myocardial infarction, and 424 - diseases of the endocardium. The four classes were slightly imbalanced, with 3502, 3184, 3175, and 1676 observations, respectively.

4 Methodology

In this section, we present the proposed methodology applied to the MIMIC-III dataset. The algorithm proposes tasks in three main steps: preprocessing, feature extraction, and pattern discovery. The overview of the proposed algorithm is shown in Figure 1.

4.1 Preprocessing

We first apply a preprocessing pipeline proposed by (Van Aken et al., 2021) to clean and merge the dataset.

We select the NOTEEVENTS (containing the unstructured text) and DIAGNOSES_ICD tables from the MIMIC-III database. The selected records contain clinical notes of patients, diagnoses, procedures, and ICD9 codes with admission ID columns acting as a link for all the tables. As each admission often had multiple diagnoses, we filter the data by only considering the highest priority diagnosis as the label to be predicted. We then trim the data to

the top four most common ICD9 codes.

After retaining these approximate 11,000 text records, we apply regular expressions to remove invalid characters and common stop words as well as words under three characters. We conform every remaining letter to lowercase and apply lemmatization. Finally, we remove a custom list of stop words that are ubiquitous among all text records.

Then, we process the text into a format suitable to be passed as a corpus (embedded lists). A dictionary, or key-value pair, is created from the tokens that were derived from our corpus of cleaned words.

4.2 Feature Extraction

Topic modelling (Hamed Jelodar and Zhao, 2018) is described as a method for finding a group of words (i.e topic) from a collection of documents that best represents the information in the collection. Hence, we extract features from the clean dataset using topic modelling the value of the features represented by the probabilities of topics occurring in the records. Labels are then merged with the features for unsupervised exploration; in this case, the label is the ICD9 code - the diagnostic code indicating categories of disease. We use LDA (Latent Dirichlet Allocation) for the topic model because it identifies topics best describing distinct subsets of documents within a corpus (Hamed Jelodar and Zhao, 2018).

To determine the ideal number of topics, we choose the optimal number of topics by computing coherence of the topic cluster instance (Röder et al., 2015). We find that the coherence score peaks when the number of topics is 5, 20, and 30 - and therefore we create topic models with those respective parameters. The output of our coherence scores is shown as Figure 2.

4.3 Pattern Discovery and Disentanglement

After preprocessing and extracting features from the text, the dataset has been transformed into a structured table of patients' records in rows and features in columns, which is represented as a $M \times N$ matrix, where M represents the number of patients' records and N represents the number of extracted features ¹.

¹In pattern discovery, we use the term attribute instead of feature.

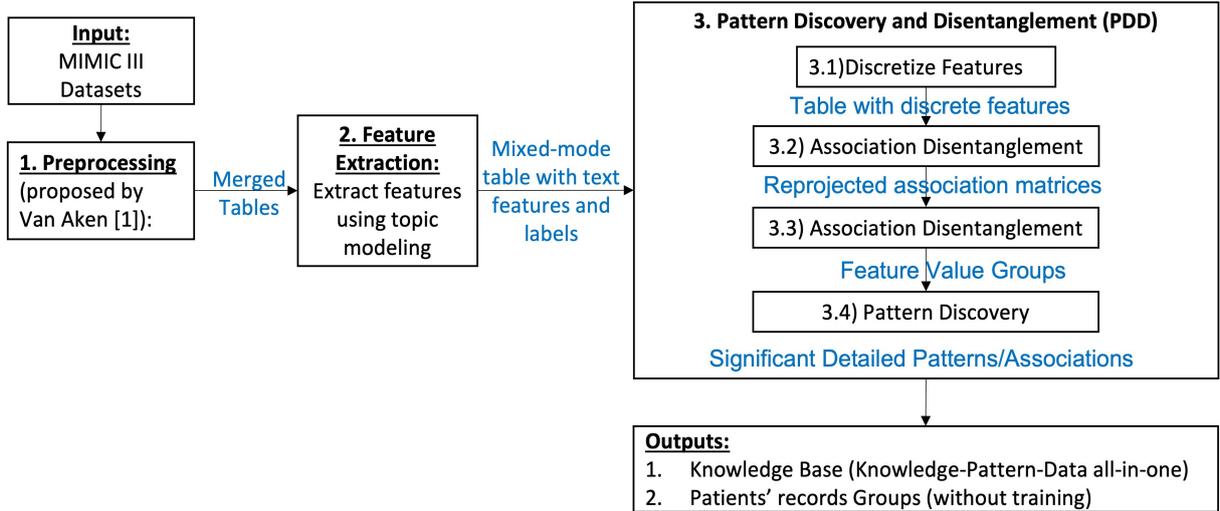


Figure 1: The overview of the proposed algorithm

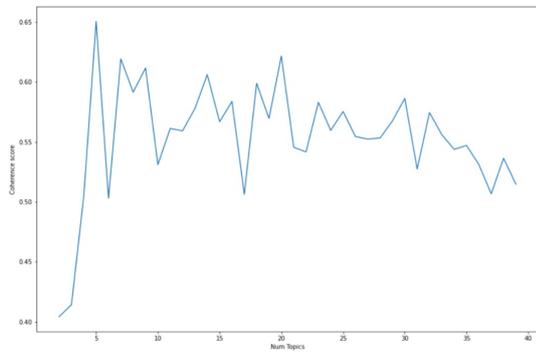


Figure 2: Optimal number of topics by coherence of the topic cluster

4.3.1 Discretize Numerical Feature Values

To detect event-based patterns, we convert the values of numerical features into categorical features by using the Equal Frequency discretization which distributes the values into equal size bins. so that numerical feature values are converted into discrete values referred to as “feature value” (meaning the discrete value for that feature). To be consistent with our existing work in PDD (Wong et al., 2021) we use the term Attribute Value (AV) instead.

4.3.2 Association Disentanglement

In order to measure the association between a pair of AVs (i.e. certain values of one attribute co-occurs with the value of another attribute), we use the statistical measure of adjusted standardized residual, abbreviated by SR, to represent the statistical weights of the AV pair, which is denoted as $SR(AV_1 \leftrightarrow AV_2)$ (shorten as $SR(AV_{12})$) and calculated by Eqn. (1) below.

$$SR(AV_{12}) = \frac{Occ(AV_{12}) - Exp(AV_{12})}{\sqrt{Exp(AV_{12})}} \times \left(1 - \frac{Occ(AV_1)}{T} \frac{Occ(AV_2)}{T}\right) \quad (1)$$

where $Occ(AV_1)$ and $Occ(AV_2)$ are the number of occurrences of AV; $Occ(AV_{12})$ is the total number of co-occurrence for two AVs in a AV pair; and $Exp(AV_{12})$ is the expected frequency and T is the total number of records.

An association matrix, treated as a vector space, is then generated to represent the strength of associations between each pair of AVs. Each row of the matrix, corresponding to a distinct AV, represents an AV-vector with SRs between that AV associated with all other AVs corresponding to the column vectors as its coordinates. We call the matrix the SR Vector Space (SRV). SRV is an $n \times n$ dimensional vector space consisting of n distinct AV-vectors.

We then use PCA to decompose SRV (Wong et al., 2021) (Wong et al., 2018) into principal components to reveal AV associations orthogonal to others AV associations, i.e. $PC = PC_1, PC_2, \dots, PC_k$ which are ranked according to the weights of the associations (eigenvalues). We then reproject the projections of AV-vectors on the principal components onto the SRV again, to obtain a set of reprojected-SRVs (abbreviated by RSRV). We refer to the PC together with its RSRV as a disentangled space.

The above process is called *Pattern Disentanglement* which allows us to take the reprojected components/vectors from PCA and use the reprojected values as new measurements/criteria to represent the strength of associations between AVs in different orthogonal disentangled spaces.

4.3.3 Obtain Attribute Value Groups with Disentangled Associations

In an RSRV, after screening in the statistical residual values (referred to as RSR) greater than 1.96, only the significant pairs of AV associations remain. Statistically, under the null hypothesis that the two AVs are independent, the adjusted residuals will have a standard normal distribution. So, an adjusted residual that is more than 1.96 (2.0 is used by convention) indicates the association is significantly greater than what would be expected (with a significance level of 0.05 or 95% confidence level) if the hypothesis were true. We can also set a threshold as 1.44 with 85% confidence, or 1.28 with 80% confidence level.

As an unsupervised learning approach, on each RSRV, we generate AV groups such that each group contains a set of AVs. We build the set of AVs up iteratively by adding AVs that are associated with AVs in the set. That is to say an AV (e.g., AV_i) that is significantly associated with another AV (e.g. AV_j) in the group will join the group, otherwise, a new AV group is generated for AV_i . Theoretically, in one projected principal component, usually two AV groups on the opposite sides are generated as two opposite groups. When such opposite groups do not exist, we may obtain AV groups only on one side of the PC. The output of this step is one or two AV groups, and each group contains a set of AVs.

Furthermore, to obtain detailed separated groups, several AV subgroups can be generated for each AV group using a similarity measure such that the similarity between two AV subclusters is specified as the percentage of the overlapping records covered by each AV subcluster. We denote each AV subgroup by a three-digit code [#PC, #Group, #SubGroup]. The AV groups or subgroups can reveal the characteristics of the records at specific groups with disentangled patterns to provide statistical evidence for further clustering or prediction. Furthermore, patient record groups are obtained according to their specific characteristics (disentangled patterns) discovered in the AV groups or subgroups.

4.3.4 Pattern Discovery on Attribute Value SubGroups

Traditional pattern clustering algorithm /citezhou2016effective, without PCA, can group patterns based on their "similarity", which is limited and time-consuming. In this case, after disentanglement and generating AV groups/subgroups, only a few AVs remain to be candidate patterns, which can reduce time consumption when high-order patterns are growing. The high-order pattern describes a statistically significant association among more than two AVs.

So far, each AV subgroup contains a set of AVs considered as candidate patterns. We then test the candidates from order > 2 (i.e. consisting of more than 2 AVs) to high order sets to determine their pattern status. Hence, we obtain a compact set of patterns which are statistically significant and interpretable. Hence PDD reduces the computational complexity drastically and produces very small and succinct pattern sets for interpretation and tracking. The disease related record groups of patients can then be explicitly revealed.

4.4 Output

The output of PDD is organized into an all-in-one representational framework known as PDD Knowledge Base. It consists of three parts: a Knowledge Section showing the hierarchical clusters such that each cluster unveil distinct characteristics of a related group of records; a Pattern Section listing the discovered patterns showing detailed associations between AVs; and the Data Section listing the record ID's, the knowledge source and pattern(s) associated with each patient by linking the patient to the Knowledge and Pattern Sections

5 Experimental Result

5.1 Topic Modeling Result

From a clinical perspective, the generated topic models correspond reasonably well with each ICD9 diagnosis. In the 20-topic model, septicemia - a widespread infection of the body, was predicted by topics containing relevant words such as "infection", "bacteria", and "culture". Conversely, topics that contained cardiovascular-related terms such as "ventricular" or "aorta" predicted the heart-related diagnoses. Additionally, the algorithm was able to discern the heart-related diagnoses from one another: dividing acute myocardial infarction

(410) from the more chronic and congenital diseases (414, 424). The algorithm may have discerned that words representing severe prognoses or procedures, such as "angioplasty", "emergency", and "death" were more correlated with acute myocardial infarction. Taken together, topic modeling and PDD provides an interpretable methodology to predict ICD9 diagnosis with reasonable accuracy when given unstructured clinical text as input.

5.2 Comparison of Unsupervised Learning

Although the process of clustering individuals does not require class label information, the entity clustering performance can be evaluated from the clustering results by two statistical measures using the presumed class labels as ground truth. In this study, since the numbers of records belonging to different classes are imbalanced, the correct prediction of the majority classes will overwhelm that of the minority classes. In this case, we followed the same evaluation method in (Van Aken et al., 2021), *balanced accuracy* (Balanced Acc. in Table 1) and *weighted F1-scores* (Weighted F1 in Table 1), to evaluate performance of both supervised and unsupervised results. Balanced accuracy is defined as the average of recall obtained in each class (Broder sen et al., 2010) and the weighted F1-score is calculated by averaging the support-weighted mean per class F1-scores (i.e. weights on class distribution) (Chakravarthi et al., 2020). Both above results are referred to the *sklearn.metrics* package in Python 3.0 (Pedregosa et al., 2011).

We compared the clustering results of PDD with the classical clustering algorithm, K-mean, as the baseline, and also two supervised learning algorithms: Random Forest (Breiman, 2001) and CNN (Brownlee, 2020). The data were split into 70% training and 30% for testing.

As for K-means, we use the *sklearn.clusters* package in Python 3.0 (Pedregosa et al., 2011) with all default parameter settings and assign the number of clusters as four. For Random Forest, we apply the default parameter settings from the package of *sklearn.ensemble.RandomForestClassifier* in Python 3.0 (Pedregosa et al., 2011).

For CNN (Brownlee, 2020), we trained a CNN model with the input layer as a reshaped cleaned dataset with probabilities of topics or extracted words and ICD9 labels. The architecture is as follows: a 1D CNN layer, followed by batch normalization, then a dropout layer for regularization (Li

et al., 2019), and finally a 1D max-pooling layer. After the CNN and pooling, the learned features are flattened to one long vector and passed through a fully connected layer before the output layer for prediction. We used Adam optimizer with a learning rate of 0.001 trained on 25 epochs with a batch size of 32.

As the baseline comparison for features, we also applied all supervised and unsupervised learning algorithms on the dataset with words extracted using TFIDF (Jones, 1972). In a corpus, frequent words in one document tend to be frequent in all other documents. TFIDF (term-frequency-inverse document frequency) is an algorithm that scores words that are distinctively frequent in a particular document but not necessarily within the general corpus. TFIDF can be computed as:

$$\text{tf-idf}(t,d) = \text{tf}(t,d) \times \text{idf}(t)$$

where tf refers to the term frequency (proportion of a particular term t over all terms); and

$$\text{idf}(t) = \log \frac{1+n}{1+\text{df}(t)} + 1$$

where n is the total number of documents in the set and df is the number of documents containing the term t.

To discover associations among features and class labels and to make the interpretation meaningful, we did not keep all words in TFIDF, but selected the top 40 words with a feature selection algorithm by Random Forest.

The comparison results are shown in Table 1. It is interesting to observe that PDD outperformed other models but underperformed when applied on the TFIDF results, which consist of the results of K-means. Both supervised learning algorithms, Random Forest and CNN perform better on the TFIDF dataset. The reason should be that the top 40 words (feature) are selected based on classification results.

When topic modeling results are used as a dataset, PDD outperforms K-means and even the two other supervised learning algorithms, with balanced acc.=0.78 and weighted F1-score=0.78, when only 5 topics are used. As for Random Forest, it performs better when applied to the topic modelling results with 20 topics than another the two experiments running on 5 topics and 30 topics. While as for CNN, the results of experiments on 30 topics are slightly better than the results on 20 topics.

Comparison	Balanced Acc.	Weighted F1
K-means		
<i>TFIDF</i> ₄₀	0.48	0.42
<i>TM</i> ₅	0.62	0.57
<i>TM</i> ₂₀	0.50	0.54
<i>TM</i> ₃₀	0.51	0.42
Random Forest		
<i>TFIDF</i> ₄₀	0.81	0.81
<i>TM</i> ₅	0.63	0.66
<i>TM</i> ₂₀	0.72	0.74
<i>TM</i> ₃₀	0.71	0.74
CNN		
<i>TFIDF</i> ₄₀	0.86	0.85
<i>TM</i> ₅	0.63	0.67
<i>TM</i> ₂₀	0.71	0.73
<i>TM</i> ₃₀	0.70	0.73
PDD		
<i>TFIDF</i> ₄₀	0.45	0.41
<i>TM</i> ₅	0.78	0.78
<i>TM</i> ₂₀	0.74	0.72
<i>TM</i> ₃₀	0.73	0.71

Table 1: Experimental Result Comparison.

One important notion we would like to bring forth is that, even if the accuracy score reflects the algorithm performance to some extent, class labels may not always be reliable in supervised classification algorithms. On the contrary, clustering merely recognizes patterns in the data and holds no such risk.

5.3 Interpretability

From the perspective of interpretability, when the topic modeling dataset with top 5 and top 20 topics were compared, the clustering performance of PDD is superior to all the other methods. As an example, we present the PDD Knowledge Base on 5 topics and 20 topics as shown in Figure 3.

The first three columns show the knowledge space, which are clustering results of PDD and statistical measurement of each pattern. The clusters are identified by a three-digital code [#PC, #Group, #Subgroup] (PC: Principal Component, Group: pattern groups in the same principal component, Subgroup: pattern Sub-group in the same pattern group). We observe that, in the first principal component, two opposite groups are discovered: one where ICD9=4XX, and the other where ICD9 = 038. All ICD9=4XX are diseases related to heart disease, while ICD9=038 is related to Septicemia,

so these are two opposite groups. Then in the second principal components, ICD9=424(diseases of the endocardium) was separated, still showing opposite patterns with ICD9=38. Finally, in the third principal component, ICD9=424 was separated from ICD9=410(acute myocardial infarction).

Then, the pattern space shows the discovered significant associations between ICD9 code and the extracted topics. To be more specific, the unveiled knowledge can be summarized as below.

- ICD9=424,410,414 (heart diseases) show similar patterns with Topic 0 (Medication) showing low probabilities.
- ICD9=424 (endocardium disease) and 414 (chronic ischemic heart disease) show more closed patterns compared to 410 (acute myocardial infarction), topic 4 (Intensive Care/Infection) showing low probability. And the unique characteristic of ICD9=424 (endocardium disease) is that Topic 1 (Cardiovascular 1) showing high probability.
- ICD9=38(septicemia) shows opposite characteristics compared to ICD9=4XX, with Topic 0 (Medication) showing high probability, Topic 2 (Cardiovascular 2) showing low probability, and Topic 4 (Intensive Care/Infection) showing high probability.

The data space shows the records IDs that are covered by the patterns. For example, the first association pattern listed in the first row of the knowledge base can be covered by the records with ID = 2,11,44,53,63 and so on. And all above records belong to the group labeled as ICD9=410, which is same with the discovered pattern.

In addition, Figure 4 shows the partial knowledge base on 20 topics dataset. As same with the above results, in the first principal component, two opposite groups are discovered: one where ICD9=4XX (heart diseases), and the other where ICD9 = 038 (septicemia). But the difference is that three subgroups (i.e. 424, 414, 410) are detected related to three different ICD9 codes in the first group in the first principal component.

Similar to the above results using 5 topics, the discovered significant patterns can be summarized for 20 topics as below. Since the most of topics are not clear, we highlighted the meaning for partial topics.

PDD Knowledge Base										
Knowledge Space				Pattern Space						Data Space
				Attributes (i.e. Topics in this study)						
PC	Group	SubGroup	Residual	ICD9	Topic 0	Topic 1	Topic 2	Topic 3	Topic 4	Records ID
1	1	1	24.99	410	[0.00 0.01]		[0.03 0.17]	[0.13 0.95]	[0.07 0.36]	#2, #11, #44, #53, #63, ...
1	1	1	11.71	414	[0.00 0.01]		[0.17 0.94]	[0.13 0.95]	[0.00 0.07]	#62, #88, #93, ...
1	1	1	13.64	424	[0.00 0.01]	[0.42 0.97]	[0.17 0.94]		[0.00 0.07]	#1, #63, #184, ...
1	2	1	51.07	38		[0.18 0.42]	[0.00 0.03]	[0.03 0.13]	[0.36 0.97]	#35,#53,#77,#80,...
1	2	1	86.06	38	[0.01 0.84]	[0.00 0.18]	[0.00 0.03]		[0.36 0.97]	#84, #96, #99,...
1	2	1	56.5	38	[0.01 0.84]	[0.00 0.18]		[0.03 0.13]	[0.36 0.97]	#84,#126,#130,...
2	1	1	10.55	424		[0.42 0.97]	[0.17 0.94]		[0.00 0.07]	#1, #63, #176,...
2	2	1	85.89	38		[0.00 0.18]	[0.00 0.03]	[0.03 0.13]	[0.36 0.97]	#12, #83, #84, ...
3	1	1	18.99	424		[0.42 0.97]	[0.00 0.03]	[0.03 0.13]	[0.00 0.07]	#206, #225, ...
3	2	1	19.1	410	[0.00 0.01]	[0.18 0.42]	[0.17 0.94]		[0.07 0.36]	#8, #64, #75,...
3	2	1	31.56	410	[0.00 0.01]	[0.00 0.18]		[0.13 0.95]	[0.07 0.36]	#2, #42, #53, ...

Note: PC=Principal Component; Group=Attribute Value Group; SubGroup = Attribute Value Sub-Group;

Figure 3: The PDD Knowledge Base for 5 topics are used as input.

PDD Knowledge Base													
Knowledge Space				Pattern Space									Data Space
				Attributes (i.e. Topics in this study)									
PC	Group	SubGroup	Residual	ICD9	Topic 0	Topic 1	Topic 2	...	Topic 16	Topic 17	Topic 18	Topic 19	Records ID
1	1	1	19.76	424	[0.01 0.42]	[0.03 0.54]	[0.03 0.44]	...					#1, #9, #13,...
1	1	2	9.39	410	[0.01 0.42]		[0.03 0.44]	...		[0.07 0.45]			#2, #4, #5, #7,...
1	1	3	26.59	414	[0.01 0.42]		[0.03 0.44]	...					#3, #6, #16,...
1	2	1	50.27	38	[0.00 0.01]	[0.00 0.01]	[0.00 0.03]	...	[0.00 0.02]		[0.00 0.01]		#9, #12, #16,...
2	1	1	24.46	424	[0.01 0.42]		[0.00 0.03]	...	[0.02 0.05]			[0.02 0.04]	#1, #9, #13,...
2	1	2	33.81	414	[0.01 0.42]	[0.03 0.54]	[0.00 0.03]	...	[0.02 0.05]		[0.01 0.03]	[0.02 0.04]	#3, #6, #16,...
2	2	1	15.28	410		[0.00 0.01]	[0.03 0.44]	...					#2, #4, #5, #7,...

Note: PC=Principal Component; Group=Attribute Value Group; SubGroup = Attribute Value Sub-Group;

Figure 4: The PDD Knowledge Base when Top 20 topics are used as input.

- ICD9=424 (diseases of the endocardium) and 414 (chronic ischemic heart disease) shows similar patterns, for example:
 - high** probabilities appear in the topics 1,2(Cardiovascular/Surgery),5,16;
 - and topics with **low** probabilities are topics 6, 7 (Status/Consciousness), 8 (Lung disease), 9
- while 038 (septicemia) shows opposite patterns, for example:
 - topics with **high** probabilities are topics 3, 4 (Intensive care/Infection), 7 (Status/Consciousness), 8 (Lung disease)
 - and **low** probabilities appear in the topics 0(Heart anatomy) 1, 2 (Cardiovascular/Surgery), 5, 12 (Cardiovascular), 16, 18;

6 Conclusion

In this work, we propose a novel two-step algorithm, combining NLP techniques with pattern discovery to solve the interpretability and unsupervised learning tasks for clinical data analysis. Experiments show results from both clustering accuracy and interpretability.

As for the clustering results, PDD performs better than K-means, especially when applied to the dataset extracted by topic modeling. Clustering results of PDD based on the discovered patterns may reflect the functional sources of the original dataset instead of class labels.

References

Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew

628	McDermott. 2019. Publicly available clinical bert embeddings. <i>arXiv preprint arXiv:1904.03323</i> .	
629		
630	David M Blei, Andrew Y Ng, and Michael I Jordan.	
631	2003. Latent dirichlet allocation. <i>the Journal of machine Learning research</i> , 3:993–1022.	
632		
633	Leo Breiman. 2001. Random forests. <i>Machine Learning</i> , 45:5–32.	
634		
635	Roselie A Bright, Summer K Rankin, Katherine Dowdy,	
636	Sergey V Blok, Susan J Bright, and Lee Anne M	
637	Palmer. 2021. Finding potential adverse events in the	
638	unstructured text of electronic health care records:	
639	Development of the shakespeare method. <i>JMIRx</i>	
640	<i>Med</i> , 2(3):e27017.	
641	Kay Henning Brodersen, Cheng Soon Ong, Klaas Enno	
642	Stephan, and Joachim M Buhmann. 2010. The balanced	
643	accuracy and its posterior distribution. In <i>2010</i>	
644	<i>20th international conference on pattern recognition</i> ,	
645	pages 3121–3124. IEEE.	
646	Jason Brownlee. 2020. 1d convolutional neural network	
647	models for human activity recognition .	
648	Bharathi Raja Chakravarthi, Ruba Priyadharshini,	
649	Vigneshwaran Muralidaran, Shardul Suryawanshi,	
650	Navya Jose, Elizabeth Sherly, and John P McCrae.	
651	2020. Overview of the track on sentiment analysis	
652	for dravidian languages in code-mixed text. In <i>Forum</i>	
653	<i>for Information Retrieval Evaluation</i> , pages 21–24.	
654	Jinying Chen, John Lalor, Weisong Liu, Emily Druhl,	
655	Edgard Granillo, Varsha G Vimalananda, and Hong	
656	Yu. 2019. Detecting hypoglycemia incidents reported	
657	in patients’ secure messages: using cost-sensitive	
658	learning and oversampling to reduce data imbalance.	
659	<i>Journal of medical Internet research</i> , 21(3):e11990.	
660	Phil Culliton, Michael Levinson, Alice Ehresman,	
661	Joshua Wherry, Jay S Steingrub, and Stephen I	
662	Gallant. 2017. Predicting severe sepsis using text	
663	from the electronic health record. <i>arXiv preprint</i>	
664	<i>arXiv:1711.11536</i> .	
665	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	
666	Kristina Toutanova. 2018. Bert: Pre-training of deep	
667	bidirectional transformers for language understand-	
668	ing. <i>arXiv preprint arXiv:1810.04805</i> .	
669	Jinyue Feng, Chantal Shaib, and Frank Rudzicz. 2020.	
670	Explainable clinical decision support from text. In	
671	<i>Proceedings of the 2020 Conference on Empirical</i>	
672	<i>Methods in Natural Language Processing (EMNLP)</i> ,	
673	pages 1478–1489.	
674	Tushaar Gangavarapu, Aditya Jayasimha, Gokul S Kr-	
675	ishnan, and Sowmya Kamath. 2020. Predicting	
676	icd-9 code groups with fuzzy similarity based su-	
677	pervised multi-label classification of unstructured	
678	clinical nursing notes. <i>Knowledge-Based Systems</i> ,	
679	190:105321.	
	Ryan B Ghannam and Stephen M Techtmann. 2021.	680
	Machine learning applications in microbial ecology,	681
	human microbiome studies, and environmental moni-	682
	toring. <i>Computational and Structural Biotechnology</i>	683
	<i>Journal</i> .	684
	Marzyeh Ghassemi, Tristan Naumann, Finale Doshi-	685
	Velez, Nicole Brimmer, Rohit Joshi, Anna	686
	Rumshisky, and Peter Szolovits. 2014. Unfolding	687
	physiological state: Mortality modelling in intensive	688
	care units. In <i>Proceedings of the 20th ACM SIGKDD</i>	689
	<i>international conference on Knowledge discovery</i>	690
	<i>and data mining</i> , pages 75–84.	691
	Chi Yuan Xia Feng Xiahui Jiang Yanchao Li	692
	Hamed Jelodar, Yongli Wang and Liang Zhao. 2018.	693
	Latent dirichlet allocation (lda) and topic modeling:	694
	models, applications, a survey . <i>Multimedia Tools</i>	695
	<i>and Applications</i> , 78.	696
	Jiawei Han, Hong Cheng, Dong Xin, and Xifeng Yan.	697
	2007. Frequent pattern mining: current status and	698
	future directions. <i>Data mining and knowledge dis-</i>	699
	<i>covery</i> , 15(1):55–86.	700
	Steven Horng, David A Sontag, Yoni Halpern, Yacine	701
	Jernite, Nathan I Shapiro, and Larry A Nathanson.	702
	2017. Creating an automated trigger for sepsis clinical	703
	decision support at emergency department triage	704
	using machine learning. <i>PloS one</i> , 12(4):e0174708.	705
	Zhengxing Huang, Wei Dong, and Huilong Duan. 2015.	706
	topic model for clinical risk stratification from elec-	707
	tronic health records. <i>Journal of Biomedical Inform-</i>	708
	<i>atics</i> , 58:28–36.	709
	Alistair EW Johnson, Tom J Pollard, Lu Shen,	710
	H Lehman Li-Wei, Mengling Feng, Moham-	711
	mad Ghassemi, Benjamin Moody, Peter Szolovits,	712
	Leo Anthony Celi, and Roger G Mark. 2016. Mimic-	713
	iii, a freely accessible critical care database. <i>Scien-</i>	714
	<i>tific data</i> , 3(1):1–9.	715
	Karen Sparck Jones. 1972. A statistical interpretation	716
	of term specificity and its application in retrieval.	717
	<i>Journal of Documentation</i> , 28:16.	718
	Efsun Sarioglu Kayi, Kabir Yadav, and Hyeong-Ah	719
	Choi. 2013. Topic modeling based classification of	720
	clinical reports. In <i>51st Annual Meeting of the Asso-</i>	721
	<i>ciation for Computational Linguistics Proceedings</i>	722
	<i>of the Student Research Workshop</i> , pages 67–73.	723
	Been Kim. 2021. Interpretability .	724
	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon	725
	Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang.	726
	2020. Biobert: a pre-trained biomedical language	727
	representation model for biomedical text mining.	728
	<i>Bioinformatics</i> , 36(4):1234–1240.	729
	Xiang Li, Shuo Chen, Xiaolin Hu, and Jian Yang. 2019.	730
	Understanding the disharmony between dropout and	731
	batch normalization by variance shift. In <i>Proceed-</i>	732
	<i>ings of the IEEE/CVF Conference on Computer Vi-</i>	733
	<i>sion and Pattern Recognition</i> , pages 2682–2690.	734

735 Stefan Naulaerts, Pieter Meysman, Wout Bittremieux,
736 Trung Nghia Vu, Wim Vanden Berghe, Bart Goethals,
737 and Kris Laukens. 2015. A primer to frequent itemset
738 mining for bioinformatics. *Briefings in bioinformat-*
739 *ics*, 16(2):216–231.

740 Miha Pavlinek and Vili Podgorelec. 2017. Text classi-
741 fication method based on self-training and lda topic
742 models. *Expert Systems with Applications*, 80:83–93.

743 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel,
744 B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,
745 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,
746 D. Cournapeau, M. Brucher, M. Perrot, and E. Duch-
747 esnay. 2011. Scikit-learn: Machine learning in
748 Python. *Journal of Machine Learning Research*,
749 12:2825–2830.

750 Michael Röder, Andreas Both, and Alexander Hinneb-
751 urg. 2015. [Exploring the space of topic coherence](#)
752 [measures](#). In *Proceedings of the Eighth ACM Interna-*
753 *tional Conference on Web Search and Data Mining,*
754 *WSDM ’15*, page 399–408, New York, NY, USA.
755 Association for Computing Machinery.

756 Betty Van Aken, Jens-Michalis Papaioannou, Manuel
757 Mayrdorfer, Klemens Budde, Felix A Gers, and
758 Alexander Löser. 2021. Clinical outcome prediction
759 from admission notes using self-supervised knowl-
760 edge integration. *arXiv preprint arXiv:2102.04110*.

761 Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh,
762 and Matt Gardner. 2019. Do nlp models know num-
763 bers? probing numeracy in embeddings. *arXiv*
764 *preprint arXiv:1909.07940*.

765 Yanshan Wang, Yiqing Zhao, Terry M Therneau, Eliz-
766 abeth J Atkinson, Ahmad P Tafti, Nan Zhang,
767 Shreyasee Amin, Andrew H Limper, Sundeep Khosla,
768 and Hongfang Liu. 2020. Unsupervised machine
769 learning for the discovery of latent disease clusters
770 and patient subgroups using electronic health records.
771 *Journal of biomedical informatics*, 102:103364.

772 Andrew KC Wong and Gary CL Li. 2008. Simulta-
773 neous pattern and data clustering for pattern cluster
774 analysis. *IEEE Transactions on Knowledge and Data*
775 *Engineering*, 20(7):911–923.

776 Andrew KC Wong, Ho Yin Sze-To, and Gary L Johan-
777 ning. 2018. Pattern to knowledge: Deep knowledge-
778 directed machine learning for residue-residue inter-
779 action prediction. *Scientific reports*, 8(1):1–14.

780 Andrew KC Wong, Pei-Yuan Zhou, and Zahid A Butt.
781 2021. Pattern discovery and disentanglement on rela-
782 tional datasets. *Scientific reports*, 11(1):1–11.