

# Similar Patients Retrieval & Events Sequence Anomaly Detection System

Chengming Xie (cx2234), Hongshan Lin (hl3353)

Topic : **Medical Insurance Fraud Detection** (Group 20)

Github: <https://github.com/steph2234/EECS6895---project>

Presentation Slides:

[https://docs.google.com/presentation/d/1LN\\_BDqEhdxfpA-u52lKkBJhX3WLurn3tqK01L4ukGyo/edit#slide=id.g777c31b5b8\\_0\\_162](https://docs.google.com/presentation/d/1LN_BDqEhdxfpA-u52lKkBJhX3WLurn3tqK01L4ukGyo/edit#slide=id.g777c31b5b8_0_162)

Date: 5/14/2020

## 1. Abstract of our project

Aiming to reduce costs for either insurance companies or the government, our project tries to build a clinical retrieval system to detect anomaly cases among electronic medical records. Previous work has shown that machine learning for data-driven diagnosis has been actively studied in the medical industry to provide better healthcare. However, the challenges associated with characteristics of electronic clinical records may include: imbalanced-class massive data, irregularity in time, and sparsity hinder performance. To address this challenge, we deal with these medical events sequences based on natural language processing techniques. At the final stage, we develop a visual analytics system to support comparative studies of patient records. Through its interactive interface, the user can quickly identify patients of interest and conveniently review both the temporal and multivariate aspects of the patient records.

## 2. Introduction

Healthcare is a major industry in the United States and it is important in the lives of many citizens, but unfortunately the continually rising in Medicare care overall experience and high costs of health-related services leave many patients with limited medical care. The impact of Medicare Fraud is estimated to be between 3% to 10 % among the nation's total healthcare spending, and reached approximately \$2 billion including 165 medical professionals in 2018. In order to reduce costs for either insurance companies or the government, our project aims to build a clinical retrieval system to detect anomaly cases among electronic medical records. Our project also helps medical workers to get better decisions by identifying sequences that deviate

from the typically occurring behavior. For example, a doctor may be interested in finding patients whose postoperative response is different from other patients who have had the same surgery, so that the doctors can provide personalized care plans for similar patients in the future.

Previous work has shown that Machine learning for data-driven diagnosis has been actively studied in medicine to provide better healthcare and publicly available medicare claims data can be leveraged to construct models capable of automating fraud detection. However, the challenges associated with characteristics of electronic clinical records include: imbalanced-class big data, irregularity in time, and sparsity hinder performance. Effective use of event sequence data requires identifying sequences that deviate from the typically occurring behavior. Moreover, most prior work in event sequence anomaly detection focuses on detecting only one type of anomaly with a specific application domain. Such approaches do not provide flexible ways of identifying diverse types of anomalies, and provide only limited support for result interpretation.

To address the challenge, we try several different methods. We calculate similarities between medical records by dealing with these medical records as event and sequence embeddings. We first use Skip-gram to learn the vector representations of medical events codes and set the window size to learn the temporal relationship within the patient's record. To make a comparison of sequences with different lengths easier, we incorporate Dynamic time warping sequence based alignment (*DTW*), and distances are calculated using the events' vector representations. In the second approach, we employ Variational AutoEncoders (*VAE*) in order to project events sequences with different lengths to vectors of the same length. The model learns latent representations for patients' event sequences and detects anomalies that deviate from the overall distribution.

These two approaches are then followed by Local Outlier Factor (LOF) and other anomaly detection analysis separately. We evaluated models performances by comparing the outlier results with several different embedding approaches and their following analysis. At the final stage, we develop a visual analytics system to support comparative studies of patient records. Through its interactive interface, the user can quickly identify patients of interest and conveniently review both the temporal and multivariate aspects of the patient records.

### 3. Related Work

In *Visual Progression Analysis of Event Sequence Data* by Shunan Guo, Zhuochen Jin, David Gotz, Fan Du, Hongyuan Zha and Nan Cao, they proposed leveraging the techniques of segmenting a sentence into thematically consistent sub-phrases in text mining area to analyze the medical records of patients. They introduced an event embedding algorithm which converts each event into a vector derived from the context of surrounding events within a given group of sequences. The co-occurrence likelihood of events is also captured by the distance between their vectors. They used a skip-gram based three layer neural network to learn the latent representation of events formulary code.

On top of the embedding matrix learned by the neural network, they employed Dynamic Time Warping to align a group of patients' event sequences with variable lengths and event orders which addressed the challenges of different lengths of electronic health records.

In their next paper, *Visual Anomaly Detection in Event Sequence Data*, they proposed an unsupervised anomaly detection algorithm based on Variational AutoEncoder which learns latent representations for all sequences in the dataset and detects anomalies that deviate from the overall distribution. As they said in their paper, "the model can estimate an underlying normal progression for each given sequence represented as occurrence probabilities of events along the sequence progression. Events in violation of their occurrence probability (i.e., event occurrences with small occurrence probability, and absent events with large occurrence probability) are identified as abnormal."

In the paper, *Learning Representations from Healthcare Time Series Data for Unsupervised Anomaly Detection*, Joao Pereira reduced the dimensionality of their events sequence embedding results using PCA and t-SNE to visually help finding the sequences (ECG heartbeats in their case) which do not conform with the normal pattern. They also used Clustering Analysis and Wasserstein distance-based anomaly metrics to detect the anomaly among their investigated patients' medical records.

Based on these three major papers and other literature we found, we implemented our model on the MIMIC-III dataset which is described below.

## 4. Our Data & Data Preprocessing

In this project, we used detailed and comprehensive patients' data from MIMIC-III dataset which contains 26 tables and medical records for 41000+ critical care patients. MIMIC-III (Medical Information Mart for Intensive Care III) is a large, freely-available database consisting of identified health-related data associated with over forty thousand patients who stayed in critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012. It includes demographics, vital signs, laboratory tests, medications, and more. Details are available on the MIMIC website: <https://mimic.physionet.org/>

We extracted timestamped medical records for all patients which include their diagnoses, procedures and prescriptions. Some data preprocessing steps include dropping the missing data, merging relational tables together, converting the data format, and storing the cleaned data into json files.

```
In [20]: patient_dict = dict()
        for index in admissions.SUBJECT_ID.unique():
            patient = diag_df[diag_df['SUBJECT_ID'] == index][['ADMITTIME', 'ICD_9', 'SHORT_TITLE', 'SEQ_NUM']]
            patient_dict[index] = []
            for i in range(len(patient)):
                stay = patient.iloc[i,:].to_dict()
                stay['event_type'] = 'Diagnose'
                patient_dict[index].append(stay)

In [27]: patient_dict[2]

Out[27]: [{'ADMITTIME': '2138-07-17 19:04:00',
            'ICD_9': 'V3001',
            'SEQ_NUM': 1.0,
            'SHORT_TITLE': 'Single lb in-hosp w cs',
            'event_type': 'Diagnose'},
          {'ADMITTIME': '2138-07-17 19:04:00',
            'ICD_9': 'V053',
            'SEQ_NUM': 2.0,
            'SHORT_TITLE': 'Need prphyl vc vrl hepat',
            'event_type': 'Diagnose'},
          {'ADMITTIME': '2138-07-17 19:04:00',
            'ICD_9': 'V290',
            'SEQ_NUM': 3.0,
            'SHORT_TITLE': 'NB obsrv suspct infect',
            'event_type': 'Diagnose'}]

In [21]: procedure_dict = dict()
        for index in admissions.SUBJECT_ID.unique():
            patient = procedure_df[procedure_df['SUBJECT_ID'] == index][['ADMITTIME', 'ICD_9', 'SHORT_TITLE', 'SEQ_NUM']]
            procedure_dict[index] = []
            for i in range(len(patient)):
                stay = patient.iloc[i,:].to_dict()
                stay['event_type'] = 'Procedure'
                procedure_dict[index].append(stay)
```

Figure I. Cleaned data Overview

## 5. Learn vector representations of medical events formulary code

After data processing, we fetch the medical records of patients whose disease category falls into digestive systems, Circulatory & Respiratory systems related, injury & poisoning categories. We then train our models respectively in these three categories of patients' event sequence data.

digestive system related patients				events embedding	
patient_id	68				
event_type	<div>✓ Diagnose</div> <div>Prescriptions</div> <div>Procedure</div>				
					id
	285178	2173-12-15 16:16:00	042	Diagnose	68
	285179	2173-12-15 16:16:00	486	Diagnose	68
	285180	2173-12-15 16:16:00	4254	Diagnose	68
	285181	2173-12-15 16:16:00	42820	Diagnose	68
	285182	2173-12-15 16:16:00	4280	Diagnose	68
	285183	2173-12-15 16:16:00	5849	Diagnose	68
	285184	2173-12-15 16:16:00	5859	Diagnose	68
	285185	2173-12-15 16:16:00	2639	Diagnose	68
	285186	2173-12-15 16:16:00	2848	Diagnose	68
	285187	2173-12-15 16:16:00	2761	Diagnose	68
	285188	2173-12-15 16:16:00	2859	Diagnose	68
	285189	2173-12-15 16:16:00	V141	Diagnose	68
	285190	2173-12-15 16:16:00	V071	Diagnose	68
	285191	2174-01-04 22:21:00	042	Diagnose	68
	285192	2174-01-04 22:21:00	486	Diagnose	68
	285193	2174-01-04 22:21:00	4254	Diagnose	68
	285194	2174-01-04 22:21:00	2762	Diagnose	68
	285195	2174-01-04 22:21:00	5849	Diagnose	68
	285196	2174-01-04 22:21:00	570	Diagnose	68

Figure II. Medical Records Overview  
(screenshot of our interface)

To begin with, in order to learn the similarities between events formulary code (i.e. ICD-9 and Drug Code), we implement a Skip-gram based Word2Vec model to learn the vector representation of medical events.

The training objective of our Word2Vec model is to learn event vector representations (of size 200) that aim to predict its context in our defined time window

Mathematically, given a sequence of training events,  $w_1, w_2, \dots, w_n$  (e.g. all medical events of patients who have digestive system disease,) the objective is to maximize the average

log-likelihood  $1/T \{ \sum_{t=1}^n \sum_{j=-k}^k \log p(w_{t+j} | w_t) \}$ , where  $k$  is the size of the time window. In our

model implementation, we chose 7 days as time window (i.e.  $w_{t+j}$  is any medical event happening in seven days before or after  $w_t$ .)

The probability of correctly predicting event  $w_t$  given event  $w_j$  is determined by the softmax model, which is  $p(w_{t+j} | w_t) = \exp(u^T_{w_t} v_{w_j}) / \sum_{l=1}^V \exp(u^T_{\hat{f}} v_{w_j})$ , where  $u_w$  and  $v_w$  are vector representations of  $w$  as event and its context event respectively, and  $V$  is the vocabulary containing all events code in the event corpus.

Finally, we are able to generate the vector representations for all events codes in three categories respectively. In the Word2Vec model implementation, we leverag the power of Spark Clusters and used Spark MLlib to enhance the time efficiency of the model training.

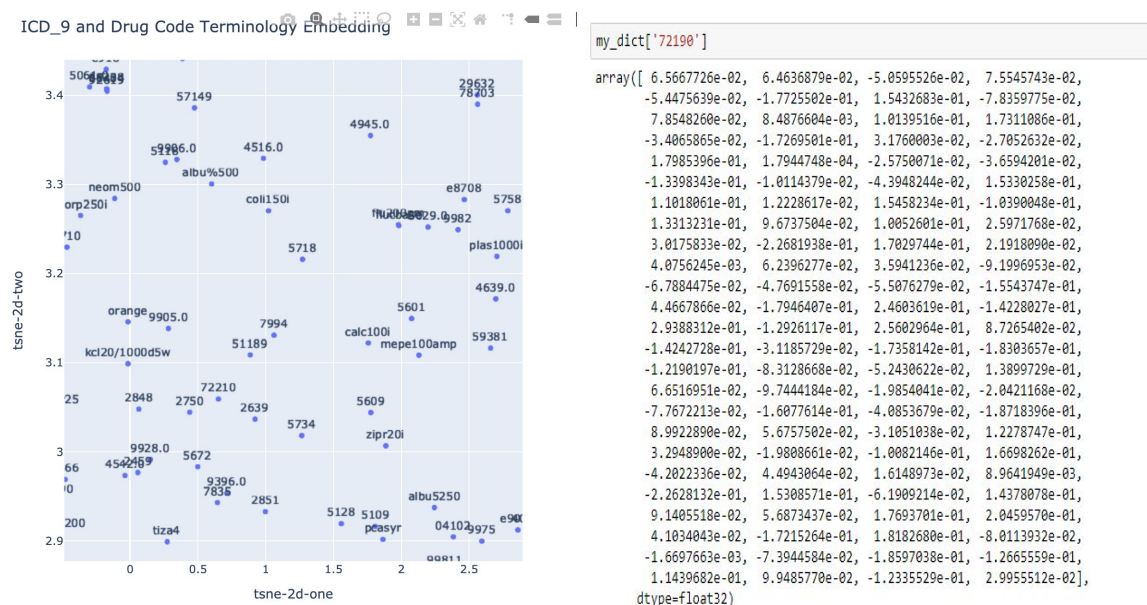


Figure III: Screenshot from our web app (i.e. t-SNE visualization)

Figure IV: Embedding vector for Code ‘72190’

## 6. Events Sequence Anomaly Detection

One biggest challenge of this project is the development of an unsupervised anomaly detection that can efficiently handle the complex temporal structure of event sequences. The methods we use in our second approach (i.e., computing DTW pairwise distances of patients and

using traditional anomaly detection models such as LOF) make more sense intuitively but fail to capture complex temporal structures in the data.

The first approach is to use LSTM-based Variational AutoEncoder to extract sequential patterns from patients' event sequence data and project them into a latent space. Then we implemented Local Outlier Factor analysis to detect the anomalies using the latent representations of each patient's event sequence.

The second approach is to implement a Dynamic programming algorithm (Dynamic Time Warping) to align events sequence data with different lengths for each patient based on the cosine similarity between two medical events. We could then compute the dtw distance (accumulate cost of cosine similarity score along the optimized warping path) between two patients' events sequence data (i.e. achieved by DTW package.)

### **5.1. First Approach: LSTM-based Variational AutoEncoder:**

The VAE encoder and the VAE decoder are using LSTM-based Recurrent neural networks to better extract sequential patterns from event sequence data. Our LSTM-based VAE model is trying to replicate the one in the paper, *Visual Anomaly Detection in Event Sequence Data*, Shunan Guo, 2019

The VAE encoder is trained to project the input sequence data, multi-hot encoding vectors of size (training\_size, time\_steps, vocab\_size), into a latent feature vector that describes a sequential distribution of events occurring in the sequence. Each coordinate represents an event type and is marked 1 if such event type happens in a certain time step and 0 otherwise.

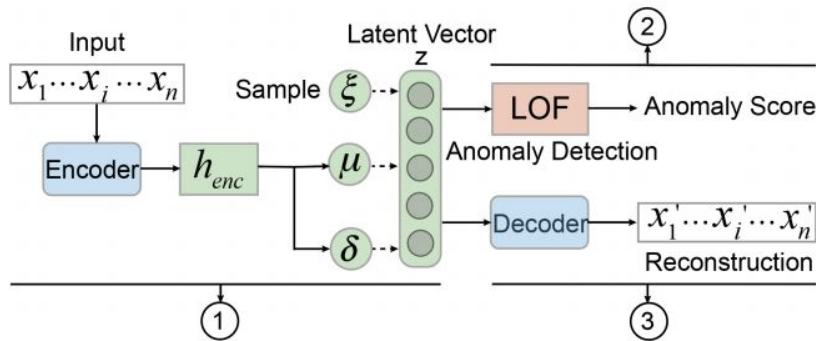
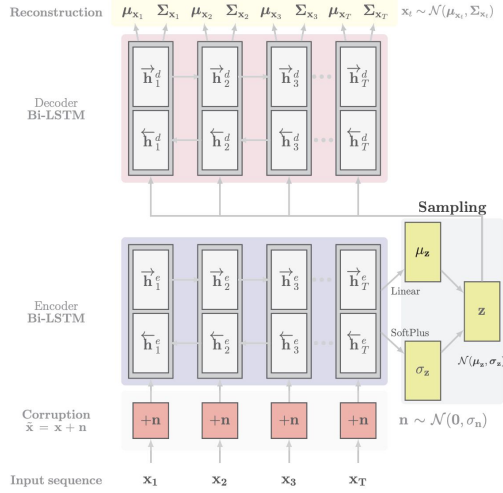


Figure V. Model Architectures



Suppose  $h_{enc} = \text{encoder}(X)$  is the hidden state vector of the last layer of LSTM RNN. The vector is then projected into two vectors  $\mu$  and  $\sigma$  to parameterize a normal distribution. When encoding, we want our latent space representation to be close to the true input (i.e.,  $q_\phi(z|x)$  and  $P(z)$ ) where we can use Kullback- Leibler Divergence to measure the closeness of two distributions.

The  $KL$  loss we define is:

$$L_{KL} = -1/2 \sum_{i=1}^{M_z} (1 + \sigma^2 - \exp(\sigma^2) - \mu^2)$$

Then we can draw a low-dimensional latent vector  $z$  by randomly sampling from the distribution.  $Z$  is sampled using a reparameterization trick:

$$z = \mu_\phi(x) + \epsilon * \Sigma_\phi^{1/2}(x) \text{ where } \epsilon \sim N(0, 1)$$

In the decoder, we can reconstruct the input sequence from the latent feature vector  $z$  (i.e.,  $z$  is fed to each layer of the RNN to estimate the probability distribution of events for each time slot).

Thus we use a softmax activation function which is then followed by a weighted cross entropy loss to calculate reconstruction loss. The reconstruction loss is defined as .

$$L_{reconstruction} = -1/n \sum_{i=1}^n \sum_{j=1}^{|E|} (w_{e_j} X_{ij} \log(x'_{ij}) + (1 - x_{ij})(\log(1 - x'_{ij})))$$

where  $w_{e_j} = 1/\log(n_j)$  is the weight to balance the prediction for unbalanced class of events in the training data and reduce the marginal importance of events with high occurrence.

The total loss is



$$L = L_{reconstruction} + w_{kl} * L_{kl}$$

where  $w_{kl}$  is the KL weight in the loss term which we set to 1 in our first training for now.

The parameter we set in our first training is roughly the same in the paper mentioned in the beginning of this section. We employ bidirectional LSTM with 300 hidden nodes in the encoding layer and LSTM with 300 hidden nodes as a decoding layer. We set the dimensions of the latent representation layer as 16 in the paper and optimize the loss function with the Adam optimizer with training data batch size of 100 for each training step.

Then we can use our encoding layer to project our training event sequence data into latent representations into a same dimensional space. We are then able to use the LOF Scikit-Learn function to calculate the “outlier” score for each patient. If the outputs of `localOutlierFactor.fit_predict()` function equals to -1, we can question the patient as an anomaly.

In our web application, we visualized the latent representation using PCA and t-SNE to allow users to have a better view of the results (visualizations included in System Overview.)

## **5.2. Second approach: DTW Distance based approach**

After fetching all vector representations for all medical events of patients in a specific disease category, sequences are then aligned temporally using a dynamic programming algorithm (DTW) , and distances are calculated using the event vectors. The algorithm measures similarity between sequences by estimating the similarity between each pair of events respectively in these sequences based on the cosine/euclidean similarity of the corresponding event vectors.

```
distance, path = fastdtw(icd_2080,icd_4579,dist=euclidean)
print(distance)
print(path)

index_a,index_b=zip(*path)

plt.plot(index_a,index_b, color='k', linewidth=2)
plt.savefig('demo_dtw', bbox_inches='tight')
```

26.71888577368587  
[(0, 0), (1, 1), (2, 1), (3, 2), (4, 3), (5, 4), (6, 5), (7, 6), (8, 7), (9, 8), (10, 9), (11, 10), (12, 11), (13, 11), (14, 11), (15, 11), (16, 11), (17, 12), (18, 13), (19, 14), (20, 15)]

Figure VI. We then compute the dtw distance (accumulate cost of cosine similarity score along the optimized warping path) between two patients' events sequence data.

Based on the DTW distance, we also create a similar patient retrieval system.

similarPatients

Patient	4579			
Calculate				
Patient historical record				
event_type	SHORT_TITLE	ADMITTIME	ICD_9	
0	Diagnose	Opn skul vlt fx w/o coma	2193-06-14 19:05:00	80061
1	Procedure	Skin closure NEC	2193-06-14 19:05:00	8659
2	Diagnose	NaN	2193-06-14 19:05:00	5185
3	Procedure	Insert endotracheal tube	2193-06-14 19:05:00	9604
4	Diagnose	Pulmonary collapse	2193-06-14 19:05:00	5180
5	Procedure	Entral infus nutrit sub	2193-06-14 19:05:00	966
6	Diagnose	Alcohol abuse-unspec	2193-06-14 19:05:00	30500
7	Procedure	Cont inv mec ven 96+ hrs	2193-06-14 19:05:00	9672
8	Diagnose	Hyperosmolality	2193-06-14 19:05:00	2760
9	Diagnose	Contusion of eyeball	2193-06-14 19:05:00	9213
10	Diagnose	Fx facial bone NEC-close	2193-06-14 19:05:00	8028
11	Diagnose	Oth coll stndng obj-driv	2193-06-14 19:05:00	E8230
12	Diagnose	Cl skul base fx w/o coma	2193-06-14 19:05:00	80111
13	Diagnose	Nasal bone fx-closed	2193-06-14 19:05:00	8020
14	Diagnose	NaN	2193-06-14 19:05:00	3488
15	Diagnose	Open wound of forehead	2193-06-14 19:05:00	87342
16	Diagnose	Mydriasis not d/t mydrtc	2193-06-14 19:05:00	37943
17	Diagnose	Contusion face/scalp/nck	2193-06-14 19:05:00	920
Top5 similar patient [('4579', 0.0), ('29051', 53.90031917837744), ('16133', 53.99492440454526), ('31020', 54.06990283559117), ('28632', 54.52492599230485)]				

Figure VII: Screenshot from our web app (i.e. similar-patient-retrieval)

### 5.2.1 Unsupervised Outlier Detection using Local Outlier Factor (LOF)

The Local Outlier Factor (LOF) algorithm computes the local density deviation of a given data point with respect to its neighbors. It is local in that the anomaly score depends on how isolated the object is with respect to the surrounding neighborhood. More precisely, locality is given by k-nearest neighbors, whose distance is used to estimate the local density. By comparing the local density of a sample to the local densities of its neighbors, one can identify samples that have a substantially lower density than their neighbors. These are considered outliers.

**5.2.1.1** Firstly, We implemented the LOF algorithm from scratch by computing reachability distance from each patient to his or her k-nearest neighbors.

$$reachability\ distance(a, b) = \max \{k\ distance(b),\ dist(a, b)\}$$

where  $k\ distance(b)$  is the distance of b to its k-th closest point

Then we computed local reachability density (lrd) for each patient  $a$  :

$$lrd(a) = 1/(\sum_{n \in N_k(a)} reachability\ distance(a, n) / k)$$

Then the LOF is basically the average ratio of the  $lrd$  of  $a$ 's k nearest neighbours to the  $lrd$  of  $a$ .

**5.2.1.2** Secondly, we use dtw statistics summary as features for each patient. After we set up our feature properly, we are able to implement anomaly detection models (ex, Local Outlier Factor (LOF), Isolation Forest).

We then visualize the results and check if the classification makes sense. We normalize and fit the metrics to a PCA to reduce the number of dimensions and then plot them in 2D highlighting the anomalies. We find both global outliers and local outliers in our graph.

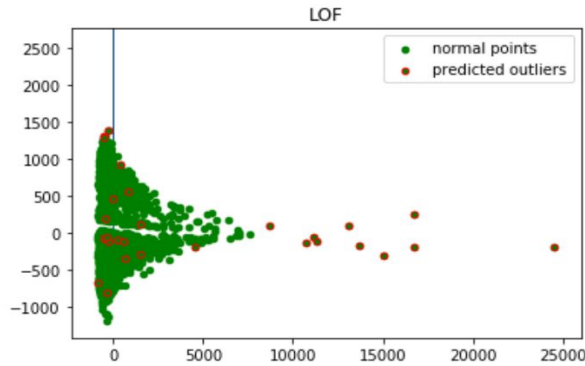


Figure VIII. Result of LOF Model in Injury & Poison related category. The red points in the left among green cluster are considered local outlier

### **5.2.2 Outlier Detection using Isolation Forest**

Another algorithm we used in DTW Distance Approach is Isolation Forest. The algorithm isolates each point in the data and splits them into outliers or inliers. Isolating an outlier means fewer loops than an inlier.

The Isolation Forest algorithm isolates observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature. The logic argument goes: isolating anomaly observations is easier because only a few conditions are needed to separate those cases from the normal observations. On the other hand, isolating normal observations requires more conditions. Therefore, an anomaly score can be calculated as the number of conditions required to separate a given observation.

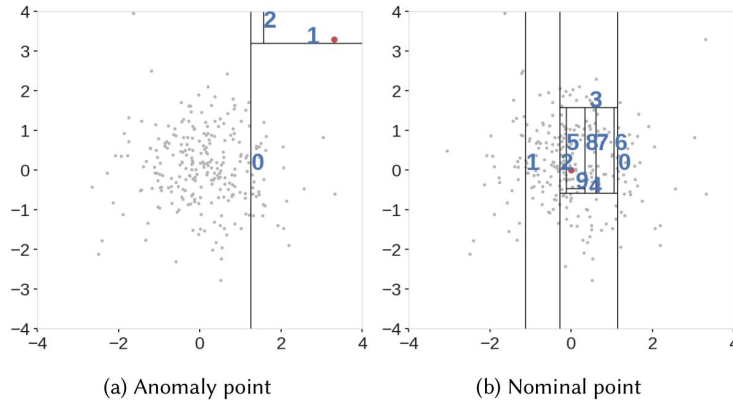


Figure IX. Concept of Isolation Forest (Inliers are harder to isolate than outliers.)

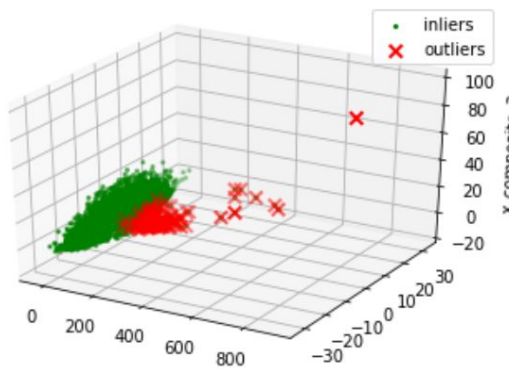


Figure X. Result of Isolation Forest Model in Injury & Poison related category. The red points in the left among green cluster are considered local outlier

## 7. Experiments & Challenges

We met a lot of challenges during the project, and solved them by doing experiments and comparing with previous works.

Our **first challenge** is data cleaning. Due to the characteristics of medical records: high dimensionality, irregularity in time, and sparsity, we spend a lot of time dealing with these raw data. To solve this problem, it took time to understand the ER between tables and figure out a way to transform these 26 tables into a format in which selected features can be used for our models. We set the time window to learn the temporal relationship between medical records.

**Moreover**, although with the help of the GCP virtual machine, we faced a big challenge for model training during the skip-gram implementation. For instance, when creating one-hot encoding vectors for input layer in the skip-gram model, we are suffering from memory leak since we had over **10,000,000+** events pairs which were too many to generate one hot vectors for

each of them; Also, when building neural network model, the model required a large amount of memory in our virtual machine. More specifically, we have tried 100+ GB extended memory VM but it still failed in the model training process and the kernel still consistently broke.

Thus, we leveraged the resilient distributed dataset (*RDD*) and *Spark* programming technique in GCP Dataproc *Clusters*. Using Scala-Spark allowed us to better scale the Word2Vec algorithms while enhancing the time efficiency by dozens of times. We could query vector representations for each medical event code. We can also order the most relevant events for each medical event by the cosine similarity of their vectors representations.

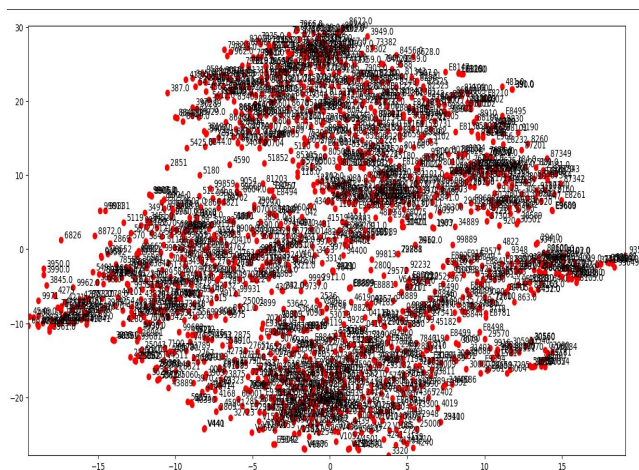
```
{
  "word": "57410",
  "vector": {
    "type": "1",
    "values": [0.0860849610805511, 0.21708840131759644, 0.38294118,
    {"word": "3782.0", "vector": {
      "type": "1",
      "values": [0.440932959318161, 0.24971017241477966, 0.364751666,
      {"word": "meg40l", "vector": {
        "type": "1",
        "values": [-0.13486158847808838, 0.0833863914012909, 0.123991,
        {"word": "macr100", "vector": {
          "type": "1",
          "values": [-0.1039232611656189, -0.1138115525456665, -0.4782,
          {"word": "29644", "vector": {
            "type": "1",
            "values": [-0.12293147295713425, -0.7968072891235352, 0.3628941,
            {"word": "8876.0", "vector": {
              "type": "1",
              "values": [0.445857435464859, 0.10378375961065292, 0.429222077,
              {"word": "35781", "vector": {
                "type": "1",
                "values": [-0.17551539838314056, -0.03375619277358055, 0.049713,
                {"word": "45355", "vector": {
                  "type": "1",
                  "values": [0.027856092802167244, 0.5323464741706848, -0.2306157,
                  {"word": "2760", "vector": {
                    "type": "1",
                    "values": [0.029653968410335015, 0.07552369853258133, 0.00801576,
                    {"word": "63.0", "vector": {
                      "type": "1",
                      "values": [0.7132171748260498, 0.11299943923958195, 0.3288291096,
                      {"word": "44023", "vector": {
                        "type": "1",
                        "values": [0.274781863527298, 0.19476436078548431, 0.3672872291,
                        {"word": "515", "vector": {
                          "type": "1",
                          "values": [0.29041996598243713, 0.05612929165363312, 0.20240484209,
                          {"word": "28803", "vector": {
                            "type": "1",
                            "values": [-0.39649176597595215, 0.18379515409469604, 0.6873161,
                            {"word": "8052", "vector": {
                              "type": "1",
                              "values": [-0.1912718117237091, 0.349385510520893506, 0.115886718,
                              {"word": "nac123.4", "vector": {
                                "type": "1",
                                "values": [0.25673702359199524, 0.2667843997478485, -0.1153,
                                {"word": "28249", "vector": {
                                  "type": "1",
                                  "values": [0.6722646355628967, -0.26024430990219116, 0.25085425,
                                  {"word": "29284", "vector": {
                                    "type": "1",
                                    "values": [0.637754499912262, 0.5026821494102478, 0.14450210332,
                                    {"word": "4233.0", "vector": {
                                      "type": "1",
                                      "values": [0.2743739187717438, -0.1364138126373291, 0.25268724,
                                      {"word": "metosusp10l", "vector": {
                                        "type": "1",
                                        "values": [1.1064503192901611, -1.4282987117767334, -1.36,
                                        {"word": "28984", "vector": {
                                          "type": "1",
                                          "values": [0.06335902959108353, -0.0102286571636796, 0.12249884,
                                          {"word": "9802.0", "vector": {
                                            "type": "1",
                                            "values": [-0.20101068913936615, -0.4324425458908881, -0.18690,
                                            {"word": "7873", "vector": {
                                              "type": "1",
                                              "values": [0.17501193284988403, 0.262361079454422, 0.20222994685,
                                              {"word": "53510", "vector": {
                                                "type": "1",
                                                "values": [0.4570658504962921, -0.014163077808916569, -0.06137,
                                                {"word": "5961", "vector": {
                                                  "type": "1",
                                                  "values": [0.1095343753695488, 0.2535054087638855, -0.1685651987,
                                                  {"word": "53510", "vector": {
                                                    "type": "1",
                                                    "values": [-0.3653593063354492, -0.7899649739265442, 0.18348076,
                                                    {"word": "4408", "vector": {
                                                      "type": "1",
                                                      "values": [0.22112111747264862, -0.049642618745565414, 0.1364563,
                                                      {"word": "99663", "vector": {
                                                        "type": "1",
                                                        "values": [0.1513883428144455, -0.3396414816379547, -0.42690059,
                                                        {"word": "hyd250", "vector": {
                                                          "type": "1",
                                                          "values": [-0.7257155680655433, 0.14191462089552155, -0.60029,
                                                          {"word": "4464", "vector": {
                                                            "type": "1",
                                                            "values": [0.27153319120407104, 0.3893878387840271, 0.358955736,
                                                            {"word": "more160bsyr", "vector": {
                                                              "type": "1",
                                                              "values": [0.17323750257492065, 0.1872567981481552, -0.30,
                                                              {"word": "rbe300", "vector": {
                                                                "type": "1",
                                                                "values": [0.43563398718833923, 1.0950158834457397, 1.0394678,
                                                                {"word": "fent2.550", "vector": {
                                                                  "type": "1",
                                                                  "values": [0.0022087967954576015, 0.23442459106445312, -0.3,
                                                                  {"word": "4919", "vector": {
                                                                    "type": "1",
                                                                    "values": [-0.020358916372060776, -0.3617551922798157, 0.1873705,
                                                                    {"word": "4290", "vector": {
                                                                      "type": "1",
                                                                      "values": [0.13315902650356293, -0.5839620232582092, 0.355841130,
                                                                      {"word": "2911.0", "vector": {
                                                                        "type": "1",
                                                                        "values": [0.9862813520431519, 0.4688606262207031, 0.411473065,
                                                                        {"word": "levo100l", "vector": {
                                                                          "type": "1",
                                                                          "values": [-0.18033763766288757, -0.9843612313270569, -0.112,
                                                                          {"word": "78900", "vector": {
                                                                            "type": "1",
                                                                            "values": [0.22901701927185059, 0.16912700235843658, -0.1003790,
                                                                            {"word": "30300", "vector": {
                                                                              "type": "1",
                                                                              "values": [0.7120336294174194, 0.00423880061134696, -0.00153554,
                                                                              {"word": "vor150", "vector": {
                                                                                "type": "1",
                                                                                "values": [0.19665685296058655, 0.19013066589832306, 0.1601507,
                                                                                {"word": "mic50l", "vector": {
                                                                                  "type": "1",
                                                                                  "values": [0.15486463904380798, -0.515732848023987, 0.951334,
                                                                                  {"word": "5853", "vector": {
                                                                                    "type": "1",
                                                                                    "values": [0.10879230499267578, -0.1274573807947098, -0.42196183,

```

Figure XI: Json format of event embedding results (Vector Representation)

word	similarity
topr50	0.6053675413131714
hctz25	0.5762994289398193
42769	0.5510246157646179
49320	0.5212512016296387
4823.0	0.5142900347709656
4881.0	0.5130040049552917
71590	0.49900346994400024
4589	0.49728673696517944
4439	0.47791099548339844
2720	0.47099921107292175
2765	0.46940693259239197
7048	0.4664931297302246
5180	0.46129298210144043
v4582	0.46122118830680847
cill150	0.4498251974582672
ator20	0.44093233346939087
566	0.4358844459056854
53081	0.4232785105705261
lev500	0.39566564559936523
gaba100	0.38527387380599976

Figure XII. Spark rdd dataframe: Most relevant medical event to Aten25. (left)



```
import org.apache.spark.ml.feature.{Word2Vec, Word2VecModel}

val rdd = sc.textFile("gs://stephixie-6895/listfile4.txt")
val rdd2 = rdd.map(line => line.toLowerCase).
  map(_.split(" ")).
  map(Tuple1.apply).
  toDF("text")
```

Waiting for a Spark session to start...

```
rdd = gs://stephixie-6895/listfile4.txt MapPartitionsRDD[1] at textFile
rdd2 = [text: array<string>]
```

```
[text: array<string>]
```

```
rdd2.take(1)
```

```
[aten25, gaba100]
```



*The third challenge* we meet is the wide distribution of word embeddings, we find our sequence records vary, nonetheless represent similar disease progressions. Patients may have doctor visits at different times, they may have different comorbidities, they may have distinct sets of comorbidities. It is necessary to account for these semantics.

ferences of length in the event sequence data. After  
dramatically in length and exact event orders that  
multiple patients suffering from the same disease may  
symptoms appear or recede at different speeds, and they  
to align event sequences based on their implicit

Figure 1 consists of two plots, (a) and (b), illustrating the relationship between Time Series A and Time Series B.

Plot (a) shows Time Series A (blue line) and Time Series B (green line) plotted against time. The x-axis is labeled "time". The y-axis is labeled "Time Series A". The plot includes a grid overlay, and the data points are represented by red dots.

Plot (b) shows the same Time Series A and Time Series B, but with a shaded region (orange) indicating a specific time interval. The x-axis is labeled "time".

mic time warping algorithms to align each patient's  
s based on their DTW distance (accumulated cost of

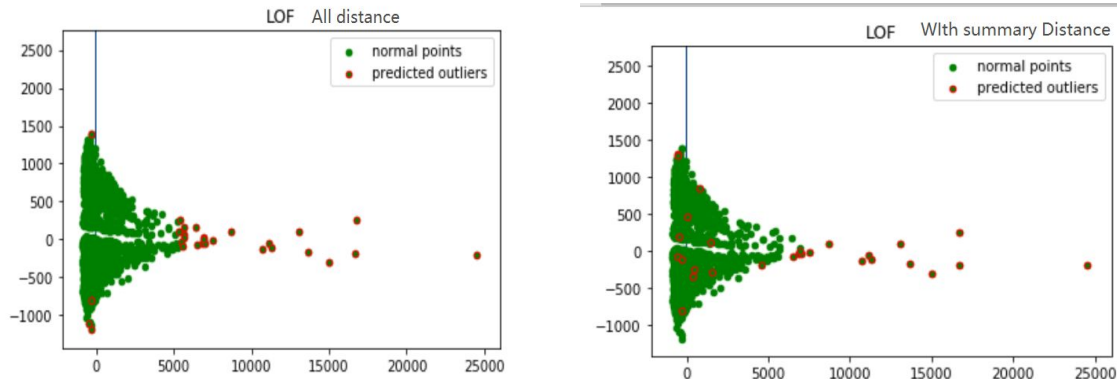
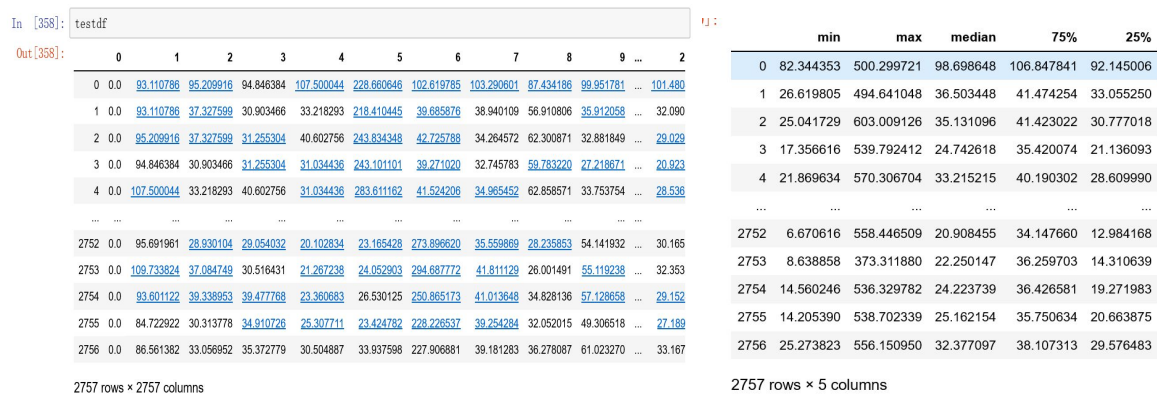
***The fourth challenge*** is the time-consuming DTW pairwise distance matrix computation in our second approach. We tried to employ our second approach on the injury & poisoning disease related patients. Since we have 2757 patients in this sub data set, we need to generate a  $2757 \times 2757$  dimensions distance matrix. In other words, we have to calculate the DTW distance for each pair of patients, a  $2757 \times 2757$  times of DTW algorithm computation. To tackle the time complexity of our algorithm implementation, we used Spark DataFrame UDF function to calculate the pairwise distance matrix between patients while optimizing the running time efficiency using RDD technique.

```

1 df_pairs.withColumn("distance",udf get_distance(df_pairs.vector2).cast(DoubleType)))).first()
2
3 Row(id=61419, vector1=[DOPABAE, *BISAIOR, *BISAS, *NACFLUFLUR, *BENB1RT, *HEPPHRRM, *HFAF1, *INFL0,5F0,
4 *D50D0M, *D50D0M, *A10M105T, *FERTNSL00M, *8856,0, *42741, *9671,0, *3961,0, *1019, *570, *78001, *3481,2,
5 *762, *2761, *5781, *2767, *4275, *4101, *4019, *49390, *4271, *78451, *V1741, *V707, *3612,0, *3613,0,
6 *3748, *1722,0, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271, *4271,
7 *000, *KCLBAE2E, *D10W100, *D50W00, *C140P,0, *C1SA10, *D5W100, *C1SA2001, *LV04V1, *3012,0, *KCL20P,
8 *R, *CFKPI1, *A1T5001, *KCLBAE2E, *B48E, *LEV04B4E, *N5100, *MAGZPM, *KCL20P, *INHVR1, *GL0C11, *A1015
9 *50, *KCL20P, *D5W100, *C140P,0, *C1SA10, *D5W100, *C1SA2001, *LV04V1, *3012,0, *KCL20P, *R, *CFKPI1,
10 *80, *INSULIN, *C1SA2001, *D5W250, *KAT51E, *KAT51E, *C1SA10, *C140G/105M, *CL0P9, *KCLBAE2E, *OXCYS,
11 *50, *KCL20P, *D5W100, *C140P,0, *C1SA10, *D5W100, *C1SA2001, *LV04V1, *3012,0, *KCL20P, *R, *CFKPI1,
12 *5325, *L10D5, *N5100, *KTR3301, *HABEEL120, *MORP21, *KCL20P, *MORP21, *FENK21, *D0C100, *D0C100, *AC
13 *T25, *KCL20P, *KCL20P, *B48I/405, *OXCYS, *MORP21, *IPR42E, *KCLBAE2E, *GUA10, *MORP1, *P02A0, *AT
14 *50, *KCL20P, *D5W100, *C140P,0, *C1SA10, *D5W100, *C1SA2001, *LV04V1, *3012,0, *KCL20P, *R, *CFKPI1,
15 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
16 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
17 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
18 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
19 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
20 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
21 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
22 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
23 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
24 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
25 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
26 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
27 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
28 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
29 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
30 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
31 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
32 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
33 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
34 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
35 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
36 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
37 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
38 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
39 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
40 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
41 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
42 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
43 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
44 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
45 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
46 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
47 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
48 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
49 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
50 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
51 *0, *OXC1C, *KCL20P, *CAP125, *D0C100, *PACT50, *BENB1RT, *MAGZPM, *OXCYS, *OXC1C, *MAGZPM, *BAG50, *M
52 *ETK500, *MAGZPM, *L0R41, *BAG50, *OXCYS, *ACT450R, *D0C100, *M0P0101, *N0M30L, *HPC20, *ACT3E25, *INFL
53 *0, *OXC1C, *KCL20P, *CAP1
```

**The fifth challenge** is high dimensionality of our dtw distance matrix, as you can see above, the shape of this matrix is 2754\*2754 (num of patients \* num of patients). Methods like LOF perform poorly in high dimensional data. To address this problem, we found not all the pairwise distances are significant and necessary in model training, so we use statistical summary as our features. For instance, we use aggregation functions to obtain the quantile distance of each patient and successfully reduce the number dimension from 2754 to 5 while retaining most of the information. By comparing the result, we find the results are almost identical but we significantly reduce the running time. However, by looking at the graph, we find that there are several red points among green clusters and we call them local outliers. We realized that we can not simply use dtw distance as lof feature because when k is larger it possible to miss local outliers (Here K=100).

Figure XVIII Output before (left) and after processing (right)



***One biggest challenge*** of this project is the development of an unsupervised anomaly detection that can efficiently handle the complex temporal structure of event sequences which we have described before.

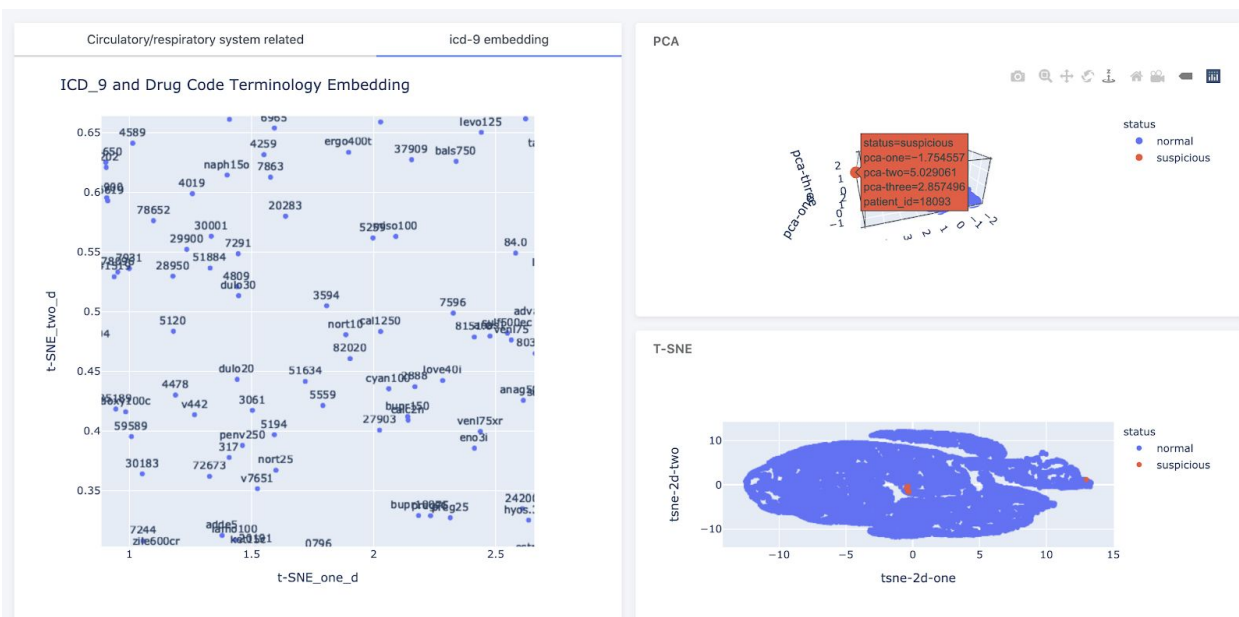
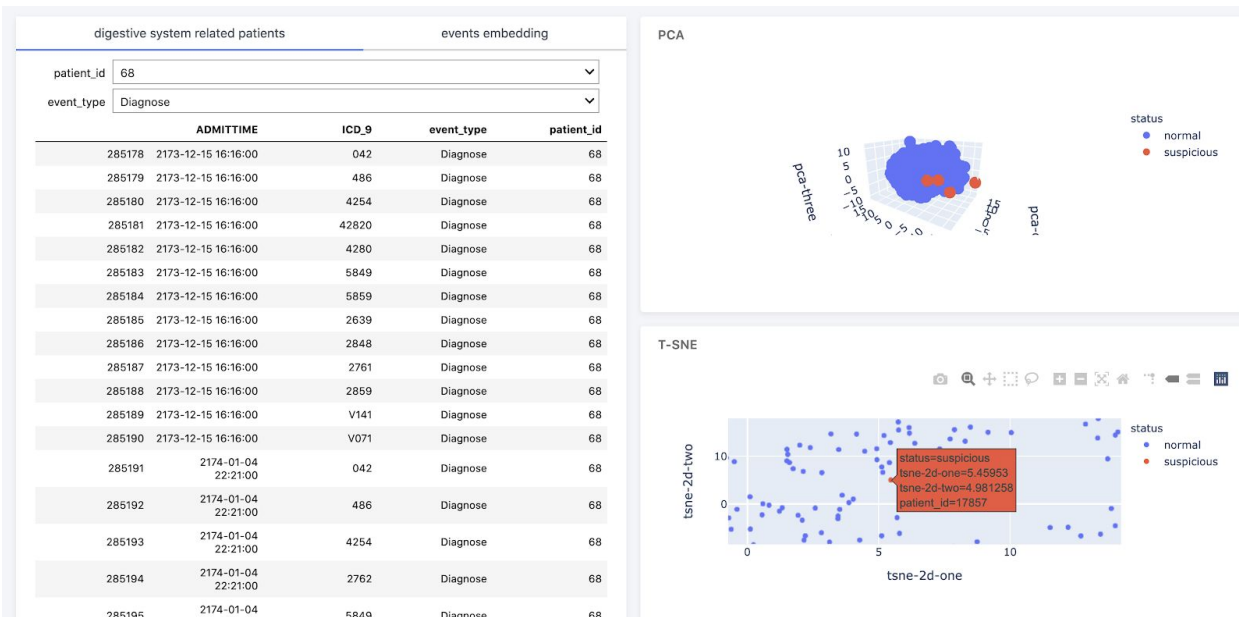
Under such circumstances, there are several recent papers regarding time series anomaly detection using LSTM-based VAE (eg. *Sequential VAE-LSTM for Anomaly Detection on Time Series*, Run-Qing Chen, 2019; *Visual Anomaly Detection in Event Sequence Data*, Shunan Guo, 2019, etc.) As Shunan Guo stated in their paper, “VAE uses a probabilistic encoder for modeling the distribution of the latent variables...such probabilities give more principled criteria for identifying anomalies and do not require model-specific thresholds. As a result, VAE can better facilitate objective judgements for deciding the boundary of anomalous sequences compared to other unsupervised algorithms.” Thus, based on the papers, we are trying to implement LSTM-based VAE on our patients’ event sequence data as our first approach to implement anomaly detection.

## **8. System Overview**

### **7.1 Digestive system & Respiratory system**

In Page 1 and Page 2 of our dashboard, we include detailed timestamped patients’ medical records which can be queried by their unique patient id and event type (i.e. Diagnosis, Prescriptions and Procedures). The first column has two tabs which include patients’ medical records as well as the events code embedding results (visualized using t-SNE) in their respective category. The right column includes the visualization after PCA and t-SNE dimensionality reduction upon LSTM-based VAE latent space representations. The anomalies are labeled as red points and users can fetch the anomaly patients id when hovering on the data point.





## **7.2 Similar patient retrieval**

In the similar patient retrieval page, we include Top 5 most similar patient retrieval and we also plot the DTW alignment for pairs of patients.

The section above shows the top 5 similar patients for a particular patient based on their DTW distance. The section below plots the DTW alignment graph for pairs of patients and shows how their sequence matches together.

## **7.3 Similar ICD-9 Retrieval**

In this page, based on the result of word2vec word embedding, we plot the distance map of ICD-9 code to help users easily find correlated icd-9 code from patients' electronic records. This can also help physicians and users to make predictions and better decisions about future treatments.

## **9. Conclusion and Future Work**

Analyzing similar patients' medical records to detect potential anomaly cases will help reduce costs for either insurance companies or the government. It also helps doctors in clinical decision making. Our work provides two embedding approaches and combines them with different unsupervised learning methods for measuring the similarity of patients and detecting anomaly cases, which can deal with high-dimensionality, irregularity, and sparsity in medical records. The visual analytics system built on top of our methods enables effective analysis of medical records from both temporal and multivariate perspectives. In the future, we will add more features and design a more comprehensive and user friendly system. Moreover, we will focus on the interpretation of these anomaly cases by delving into the case studies .

## Reference:

S. Guo, Z. Jin, Q. Chen, D. Gotz, H. Zha and N. Cao, *Visual Anomaly Detection in Event Sequence Data*, 2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA, USA, 2019, pp. 1125-1130.

Joao Pereira, *Learning Representations from Time Series Data for Unsupervised Anomaly Detection* 2019

S. Guo, et al., *Visual Progression Analysis of Event Sequence Data in IEEE Transactions on Visualization & Computer Graphics*, vol. 25, no. 01, pp. 417-426, 2019. doi: 10.1109/TVCG.2018.2864885.

Graphen AI.

Guo, Rongchen & Fujiwara, Takanori & Li, Yiran & Lima, Kelly & Sen, Soman & Tran, Nam & Ma, Kwan-Liu. (2020). *Comparative visual analytics for assessing medical records with sequence embedding*. Visual Informatics. 10.1016/j.visinf.2020.04.001.

Anne Fischer, *Generating Text From latent Space: creating novel news headlines from latent space using a Variational AutoEncoder*, 2019

Rasim Alguliyev, Ramiz Aliguliyev, Lyudmila Sukhostat, *Anomaly detection in Big data based on clustering*, Stat., Optim. Inf. Comput., Vol. 5, December 2017

MIMIC-III, a freely accessible critical care database. Johnson AEW, Pollard TJ, Shen L, Lehman L, Feng M, Ghassemi M, Moody B, Szolovits P, Celi LA, and Mark RG. Scientific Data (2016). DOI: 10.1038/sdata.2016.35. Available from: <http://www.nature.com/articles/sdata201635>