

Aalto University  
School of Science  
Degree Programme in Computer Science and Engineering

Yang Zhao

# Anomaly Detection from Patient Visit Data

Master's Thesis  
Espoo, September 28, 2016

**DRAFT! — September 28, 2016 — DRAFT!**

Supervisor: Professor Juho Rousu  
Instructor: Jaakko Hollmén, D.Sc.(Tech.)

Aalto University  
 School of Science  
 Degree Programme in Computer Science and Engineering

ABSTRACT OF  
 MASTER'S THESIS

<b>Author:</b>	Yang Zhao		
<b>Title:</b>	Anomaly Detection from Patient Visit Data		
<b>Date:</b>	September 28, 2016	<b>Pages:</b>	55
<b>Major:</b>	Machine Learning and Data Mining	<b>Code:</b>	T-110
<b>Supervisor:</b>	Professor Juho Rousu		
<b>Instructor:</b>	Jaakko Hollmén, D.Sc.(Tech.)		
<p>Hospital operation cost rises due to the growing demand for outpatient services by increasing elderly population. To reduce the operation cost and serve the patients better, improvements on the efficiency in healthcare service institutes are required. Among several potential aspects of efficiency improvements, smoother patient visits are highly desired. Thanks to the digital era, patient visits to the hospital can be recorded with all details. The Oulu Hospital in Finland starts to gather patient visits data since 2011, using queue system provided by X-Akseli Oy. Utilizing these collected data, this thesis aims at designing a practical way of detecting anomalies from patient visits. With the help this system, the hospital administrative staff could analyse the performance of the queue procedure in the hospital and optimize the procedure. Even better, the system can identify anomalies in real-time so that the patient can get immediate help when it is needed.</p> <p>The thesis explored two categories of methods: clustering method and generative methods. Four candidate algorithms, K-Means, DBSCAN, Markov Chain, and Hidden Markov Model, are discussed. The discussion suggests that DBSCAN and Hidden Markov Model are more practical. Then we proposed a new data representation and used negative binomial distribution in Hidden Markov Model to model patient states durations. The experiment result was visualized using t-SNE and evaluated by user interpretation. The analyses show that both DBSCAN and Hidden Markov Model can effectively detect anomalies from patient visits data. But in terms of time and space complexity, and real-time detection, Hidden Markov Model is a better choice.</p>			
<b>Keywords:</b>	sequence data, clustering, generative Markov models, duration modelling, Poisson distribution, negative binomial distribution		
<b>Language:</b>	English		

# Acknowledgements

I would sincerely thank Aalto University for providing me such a great study opportunity. I have gained a lot in the past three years, not only knowledge, but also beneficial lessons for my life. These memorable moments make me think more clearly about my personal long-term life objects.

Next I would sincerely thank X-Akseli Oy for providing me such an opportunity to finish my master thesis. My colleagues friendly help and trust kept me working on and finally finished this challenge task. Also, I would thank Juho Rousu and Jaakko Hollmén for agreeing to supervise my thesis. Without help from these mentioned people, it is impossible for me to complete my thesis.

Additionally, I would like to thank my colleagues while I was working in the group leaded by Professor Teemu Roos, at University of Helsinki. This invaluable experience is my first step into research world. It was in this group, that I learned how to do a research.

Finally, I am really grateful to the support of my family and friends. Without their support, I am not even able to begin my master study.

Espoo, September 28, 2016

Yang Zhao

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	Why Anomalies Matter? . . . . .	6
1.2	Background: State Flow of the X-Akseli System . . . . .	6
1.3	Problem Description . . . . .	10
1.4	Thesis Structure . . . . .	10
<b>2</b>	<b>Clustering Patient Visits</b>	<b>11</b>
2.1	K-Means . . . . .	11
2.1.1	K-Means as a Special Case of EM Algorithm . . . . .	12
2.1.2	Issues in K-Means . . . . .	13
2.2	Density Based Clustering Methods . . . . .	16
2.2.1	DBSCAN and DFS . . . . .	16
2.2.2	Notions of Density in Clusters . . . . .	18
2.2.3	Analysis of DBSCAN . . . . .	20
<b>3</b>	<b>Scoring Patient Visits by Markov Models</b>	<b>23</b>
3.1	Discrete Markov Process . . . . .	23
3.2	Hidden Markov Model . . . . .	26
3.2.1	Definition of Hidden Markov Model . . . . .	27
3.2.2	Learning and Inference . . . . .	28
<b>4</b>	<b>Experiments</b>	<b>33</b>
4.1	Data Details and Representation . . . . .	33
4.2	Methods . . . . .	35
4.2.1	Choice of Clustering Method . . . . .	36
4.2.2	Choice of Generative Method . . . . .	37
<b>5</b>	<b>Results and Discussion</b>	<b>41</b>
5.1	Clustering Method Results . . . . .	41
5.2	Generative Method Results . . . . .	44
5.3	Discussion . . . . .	45

<b>6 Summary and Conclusion</b>	<b>48</b>
<b>A First appendix</b>	<b>52</b>

# Chapter 1

## Introduction

### 1.1 Why Anomalies Matter?

As the elderly population increases, demand for outpatient services rises, which increases the operation cost in hospital [10][14]. This phenomenon requires improvements on the efficiency in healthcare service institutes. Among several potential aspects of improving efficiency, one expectation is to make the whole visit smoother.

A patient visit to the hospital consists of several phases. From enrolling to doctor treatment, each phase is affected by many factors which could result in an unpleasant experience. Examples are long waiting time, disordered treatment procedures etc. Studying unexpected care-flows is an important way to help provide better visit experience. From studying these abnormal cases, administrative staff can understand the reasons for causing the problems. Thus, the staff can balance resources allocated in the hospital for smoother service in the future. What is better, if real-time anomaly detection is available, then the hospital can provide necessary help to the patient in time. This thesis aims at developing practical anomaly detection methods on patient visits data.

### 1.2 Background: State Flow of the X-Akseli System

Studies in this thesis are based on data extracted from Oulu University Hospital, generated by X-Akseli queue system which aims at making patient reception fluent and effortless. This section describes how the X-Akseli reservation system works, in order to give the reader a general idea about how

Table 1.1: One example visit to the hospital.

reservationid	eventname	time
21332189	ENROLLING	2014-03-26 08:02:42.353
21332189	ARRIVED_IN_HOSPITAL	2014-03-26 08:02:42.517
21332189	WAITING	2014-03-26 08:03:29.007
21332189	CALLED	2014-03-26 08:07:15.061
21332189	IN_TREATMENT_ROOM	2014-03-26 08:07:15.072
21332189	CLOSED	2014-03-26 08:13:11.002

the data was produced. Patient privacy has been protected by using anonymous ids. The reservation system has gone through several updates while recording these data. Considering this situation, the discussion in this thesis adheres to the latest system, version 1.18.3.

Let's assume a patient has already made a reservation online. A typical visit scenario to the hospital is as follow. The patient arrives at the hospital lobby. Then the patient takes a queue number at the self-service kiosk, with information showing in which area the patient should wait. Next, the patient goes to the correct waiting area, shows the queue number to another kiosk in the area to check in. After this, the patient can sit down and wait to be called by the doctor. When the doctor is available, the doctor will call the patient, treat the patient and close the whole visit. If assisted diagnoses are needed, the doctor may pause the treatment. After other diagnoses are finished, the doctor then continues the treatment. In small departments, there may be no separation of waiting area and lobby. In this case, the first two steps will be integrated. The patient will not have to show the queue number to another kiosk. However, there will still be two events recorded in the back end system, but with zero transition time. The whole procedure is shown in Figure 1.1 and the state flow of the back end system is shown in Figure 1.2

One example visit is listed in Table 1.1. Notice that, there are 6 events in this visit. However, the first two events **ENROLLING** and **ARRIVED\_IN\_HOSPITAL** happened in less than 1 second. This also happens to the 4th and 5th events, **CALLED** and **IN\_TREATMENT\_ROOM**. It can be considered that the two events in these two pairs happened simultaneously. To reduce redundant information, in this thesis, only following 7 events are considered. They are: **ENROLLING**, **WAITING**, **IN\_TREATMENT\_ROOM**, **PAUSED**, **IN\_TREATMENT\_ROOM\_FROM\_PAUSED**, **CLOSED**, and **CANCELLED**.

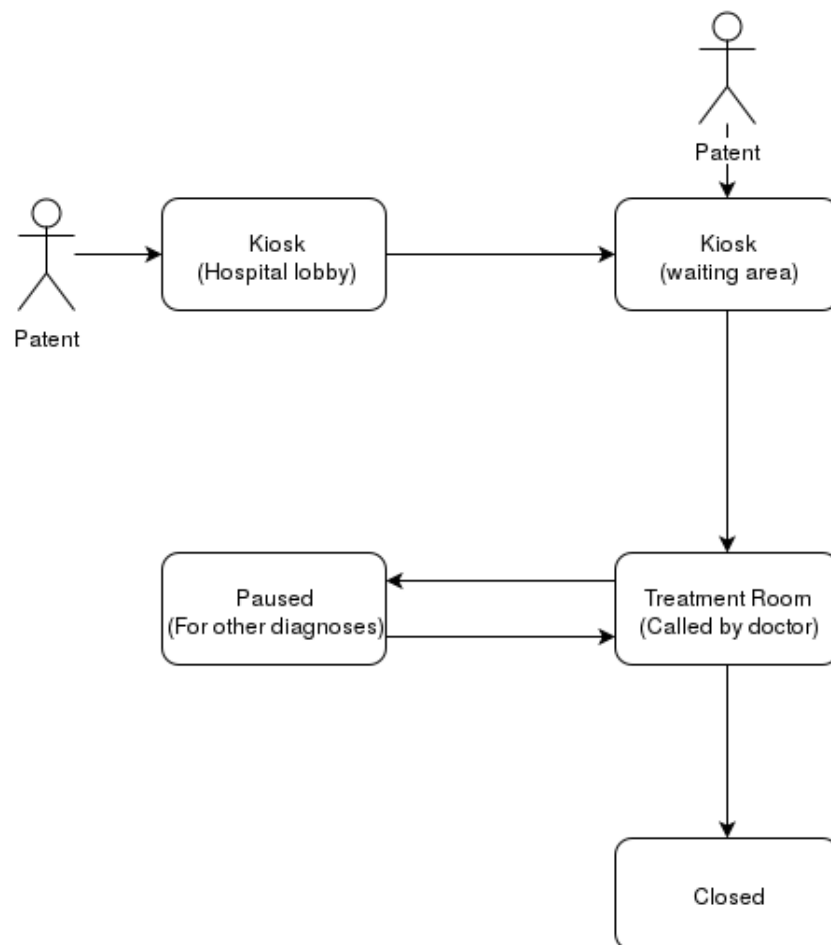


Figure 1.1: A typical visit scenario in hospital.



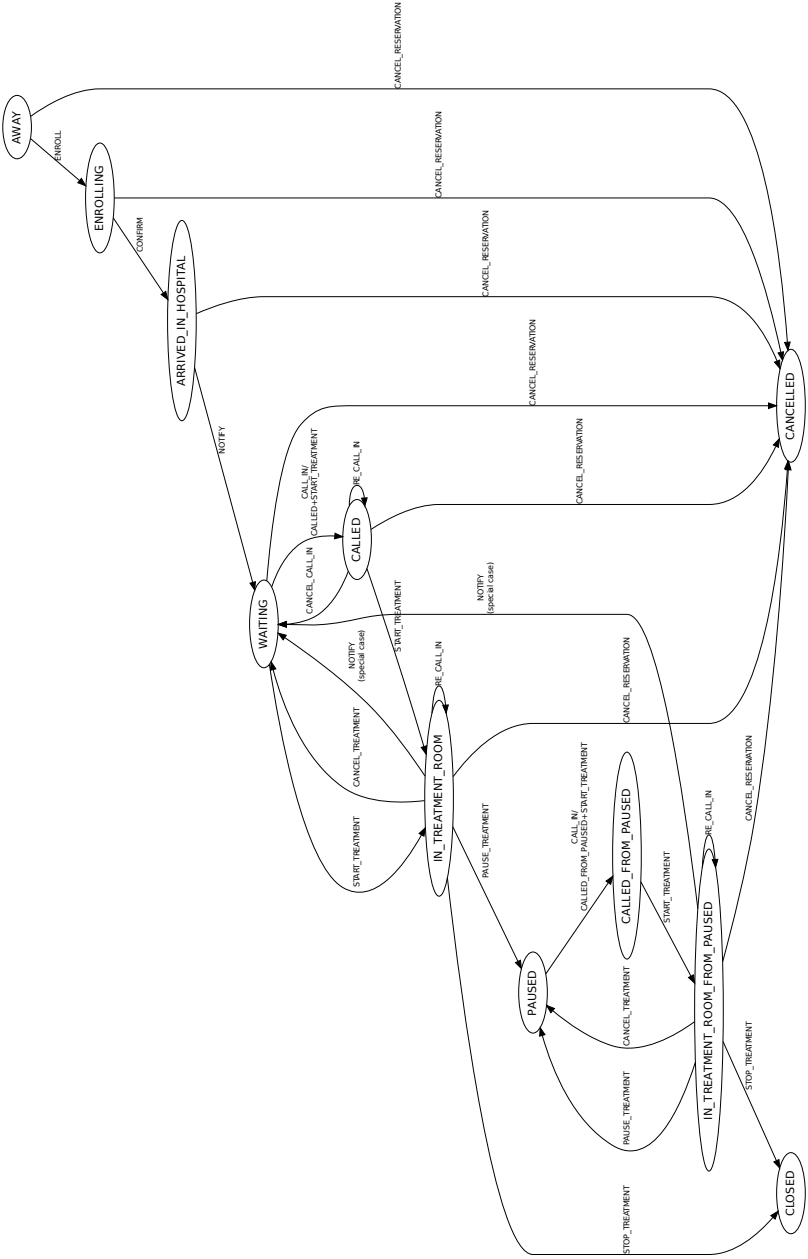


Figure 1.2: State flow in X-Akseli system, version 1.18.3.

### 1.3 Problem Description

Detecting an anomaly from hospital visits can be approached as finding outliers in a database containing time series entries. The assumption is that, the number of anomalies in the database is quite few compared to the number of normal entries. One naive approach to solve this problem is extreme value detection [1]. The idea is to set a threshold and classify any entry which has a value exceeds the threshold to be an anomaly. In the patient visits case, the threshold can be set to be a specific waiting time. Any visits have a longer time than the threshold is an anomaly. However, this approach have some obvious drawbacks. Due to the hard assignment, performance of this method heavily depends on the selection of the threshold value. However, the boundary between anomaly and normal data is not clear. There could be a “gray area” around this threshold value. Additionally, definition of anomalies in patient visit is not only limited to durations. Anomalies may also consists of strange consequences of the patient states, for example, being in `WAITING` states for several times. This may be caused by going to wrong waiting areas.

An alternative approach is to scoring each entry by showing how likely it could be an anomaly [11]. Then, anomalies are detected based on these scores. Typical methods include clustering methods and Markovian models. This approach considers anomalies from overview aspect. The concept is that, a single minor misbehaviour doesn’t necessarily lead to an anomaly. It is a sequence of uncommon behaviours that results in anomalies. This thesis decides to explore the problem using this approach.

### 1.4 Thesis Structure

The thesis aims at suggesting a practical method for detecting anomalies from hospital visits. The structure of this thesis is organized as follow. Chapter 2 introduces potential clustering method. Chapter 3 introduces potential generative model methods. Chapter 4 first describes more details and pre-process executed on the data used in the experiment. Then this chapter compares strengths and drawbacks of these methods, discusses their applicability, and describes the experiment setup. Chapter 5 presents and analyses the obtained results. Finally, Chapter 6 sums up the thesis and discusses potential future work.

## Chapter 2

# Clustering Patient Visits

As described in Chapter 1, all patients data is retrieved from the database directly. These data only records information about patient visits, but without any manual labels indicating which visits are abnormal. Due to lack of labels, unsupervised learning algorithms should be adopted. Clustering is a collection of unsupervised methods, which identifies groups of data points according to a defined similarity metric, such that objects in the same group possess higher similarities compared to objects in other groups. The clustering process does not rely on labels but the choice of similarity metrics. Variations in similarity metrics lead to different clustering methods.

Applying clustering methods in anomaly detection tasks has been studied numerously [12]. This chapter introduces two typical methods, K-Means [18] and DBSCAN [8]. Problem formulation, solutions, and potential issues are formally described using elaborated notations in following sections. However, analysis on the performance and constraints of these two methods are postponed to Chapter 4, which reveals their practicality and infeasibility in the previously described anomaly detection problem.

### 2.1 K-Means

K-Means [18] is one of the simplest unsupervised algorithm which solves clustering problem. Despite its simplicity, K-Means has gained success in various situations, including anomaly detection [4][12], image segmentation and compression [9]. K-Means is also frequently used as pre-processing for more complicated algorithms. The method can be formally defined as follow: Given a data set  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  consisting of  $N$  observations in  $D$ -dimensional space, the object is to partition the data into  $K$  groups, by defining a set of  $K$  centres  $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\}$  in the same space, and assigning each observation to

exactly one center point. Each center point represents a prototype associated with the  $k^{th}$  cluster.

The assignments can be represented using 1-of- $K$  schema. Then, for each data point  $x_n$ , a corresponding  $K$ -dimensional variable consisting of  $K$  binary elements  $r_{nk} \in \{0, 1\}$  is introduced. Among these  $r_{nk}$ , exactly one of them equals 1, which means  $\mathbf{x}_n$  belongs to the  $k^{th}$  cluster. Using this notation, evaluation of the clustering quality can be defined using the object function as follow:

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} D(\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (2.1)$$

where  $D$  is a dissimilarity metric. Common choice of the metric is  $l_1$ -norm or  $l_2$ -norm. Mahalanobis distance is also adopted while considering the covariances between the  $K$ -dimensions [6]. In the following context,  $l_2$ -norm is adopted for discussion. Intuitively, this function can be considered as the distance summation of each point to its corresponding cluster prototype  $\boldsymbol{\mu}_k$ . The K-Means aims at finding a set of  $\boldsymbol{\mu}_k$  which minimizes the object function. Finding the optimal solution for the above object function proves to be NP-Hard [2]. However, employing heuristic algorithms enables finding converged local optimal solutions. Section 2.1.1 describes one iterative algorithm, EM. Section 2.1.2 explores common issues related to K-Means and remedies.

### 2.1.1 K-Means as a Special Case of EM Algorithm

EM algorithm [7] is an iterative algorithm to find local maximum. Each iteration consists of two phases, Expectation and Maximization, which corresponds to minimizing the objective function  $J$  with respect to  $r_{nk}$  and  $\boldsymbol{\mu}_k$  respectively. In the E(expectation) step, the algorithm minimize  $J$  with respect to  $r_{nk}$ , while keeping the  $\boldsymbol{\mu}_k$  fixed. Then in the following M(maximization) step, the algorithm minimizes  $J$  with respect to  $\boldsymbol{\mu}_k$ , while keeping  $r_{nk}$  fixed.

Considering the optimization in E step, a critical observation of (2.1) is that  $\mathbf{r}_n$ 's are independent of each other. Thus, optimization on  $\mathbf{r}_n$ 's can be done separately for each  $n$ . The solution is simply setting the  $r_{nk}$  corresponding to the minimum  $\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$  to 1. Intuitively, the algorithm assigns  $\mathbf{x}_n$  to the closest cluster center. Formally, the solution can be written as

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases} \quad (2.2)$$

In the following M step, the above determined  $r_{nk}$  is clamped. Then,  $\boldsymbol{\mu}_k$  appears only in a quadratic term in  $J$ . Setting derivatives of  $J$  with respect

to  $\mu_k$  to zero, solution for  $\mu_k$  can be expressed in following closed form

$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}} \quad (2.3)$$

This step can be considered as recomputing the cluster prototype by setting  $\mu_k$  to the mean of all points assigned to that cluster.

EM keeps executing these two steps alternatively, until a convergence happens or until number of iterations exceeded a predefined value. Since in each phase, one variable is fixed and updating another variable minimizes the cost function  $J$ , convergence is guaranteed. Formal proof on convergence has been studied by MacQueen [20].

The algorithm is illustrated using dataset generated independently from three Gaussian distributions in Figure 2.1. In this demonstration, the algorithm takes  $K = 3$ , which is the correct number of components. Before running the first iteration, initialized  $\mu_k$  is required. This initialization is done by choosing three objects from the data set randomly.

### 2.1.2 Issues in K-Means

Despite the simplicity of K-Means, several underlying issues exists. The first potential is that, how to choose the value for  $K$ . In the above illustration,  $K$  was set to 3 which is the correct number of components. However, if  $K$  wasn't set to the correct value, unsatisfied clustering may be generated. Example of this issue is shown in Figure 2.2(a)-(c). To solve this problem, one practical way is drawing graph of the cost function versus value of  $K$ , as shown in Figure 2.2(d). Intuitively, when  $K$  is smaller than the true number of clusters, increasing  $K$  will lead to a huge drop of cost function value. However, when  $K$  has reached or exceeded the correct value, increasing  $K$  leads only small cost function value drop. Thus, number of clusters corresponding to the 'elbow' point can be considered as the real number of clusters.

Another issues of K-Means relates to the initialization of  $\mu_k$ . Since EM finds only local optimal solutions, a poor initialization could lead to worse clustering result. To avoid this problem, it is practical to run K-Means for several times and choose the best result according to the value of the cost function.

One more critical issue of K-Means lies in the choice of similarity metric. As mentioned at the beginning of this chapter, variations in similarity metric leads to different clustering methods. Choosing  $l_2$  norm is convenient in terms of computation, but this choice limits the type of data variables to certain types. Using  $l_2$  norm on categorical data is inappropriate since no ordering

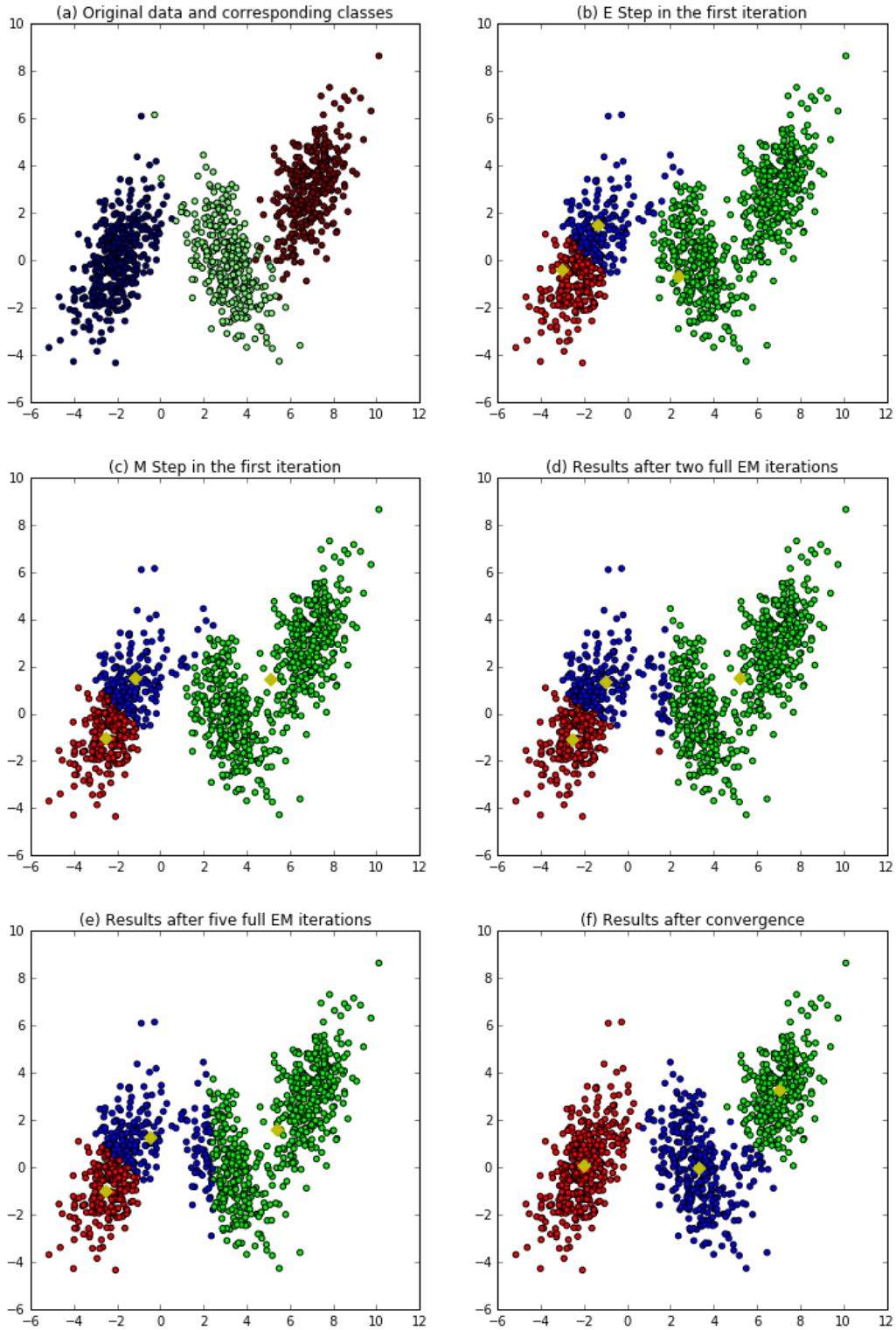


Figure 2.1: Illustration of EM algorithm on K-Means using data generated independently from three Gaussian distribution. (a) Original data and corresponding classes. Classes are denoted in different colors. (b) Assignments of each data after the first E step. The yellow diamonds represent the initialized  $\mu_k$ . (c) Updated  $\mu_k$  after the M step in the first iteration. (d)-(f) Clustering results after several successive full EM iterations until convergence is met.

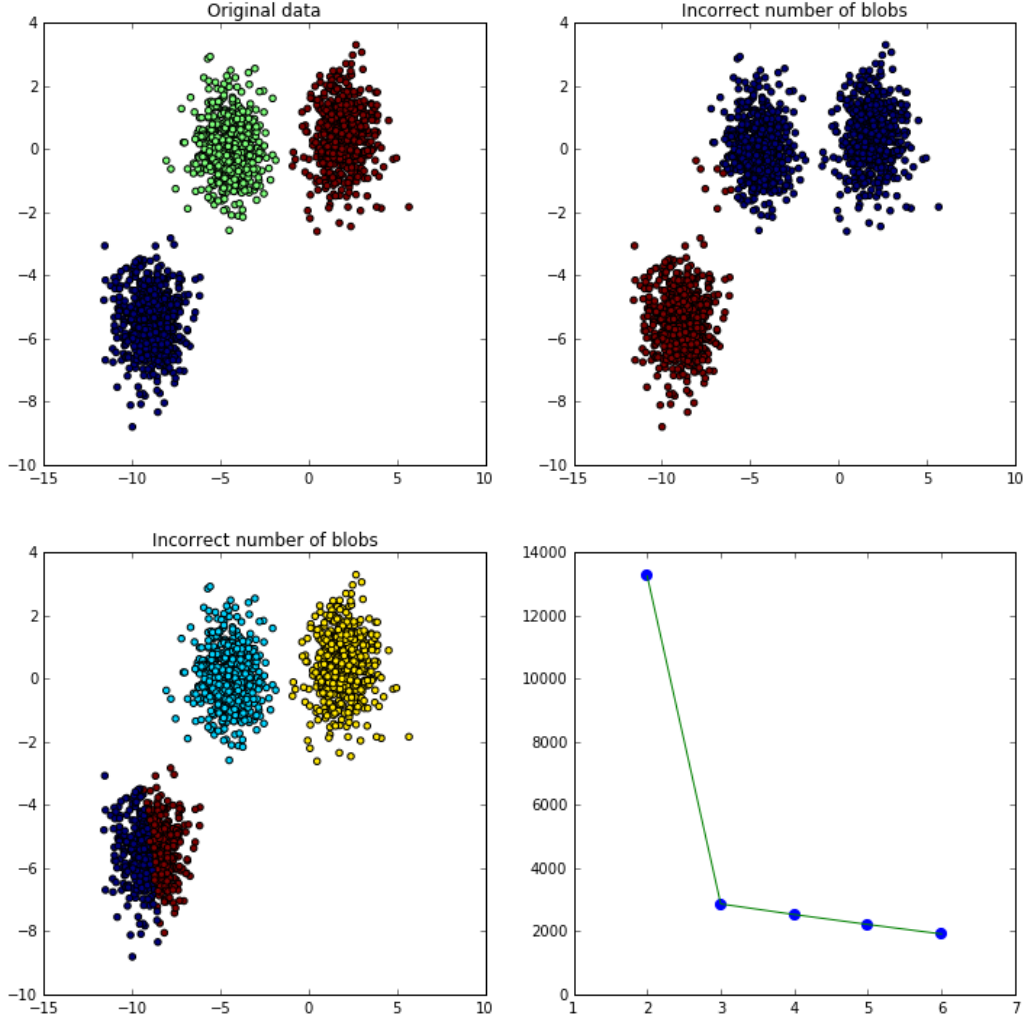


Figure 2.2: Issues of K-Means incurred by choosing inappropriate value for  $K$ . (a) Original data which contains 3 clusters. (b) Clustering result using  $K = 2$ , which is smaller than the real number of clusters. As the result shows, the upper two clusters are merged together. Also, a few points belonging to the third cluster (blue cluster in (a)), are assigned to the first cluster (green cluster in (a)). (c) Clustering result using  $K = 4$ , which is larger than the real number of clusters. As the result shows, an additional cluster was generated by dividing the third cluster into two. (d) The cost function value versus number of clusters. The elbow point appears when  $K = 3$ .

of categorical values exists. Also, this makes computing the mean value a hard problem. To use K-Means on other data types, the similarity metric should be elaborately designed.

Besides, K-Means tends to form clusters into a convex space. As shown in Figure 2.1(b)-(e), the boundary between two different clusters forms a line lying at the midway and is perpendicular to the line connecting two cluster prototypes. However, a cluster is not necessary to be convex. This also limits the application of K-Means.

Finally, K-Means is also sensitive to noise. As shown in Equation 2.3, updating  $\mu_k$  involves computing the mean of all data points assigned to that cluster. When noise objects with large deviation from other points in this cluster exist, update of  $\mu_k$  will be strongly affected.

## 2.2 Density Based Clustering Methods

This section explores another type of clustering method which forms cluster from the view of density aspect. **Density-Based Spatial Clustering of Applications with Noise (DBSCAN)** [8] considers to be one of the most successful method from this category [17]. Compared to K-Means, DBSCAN possesses following advantages: 1). DBSCAN can detect clusters of arbitrary shape. 2). DBSCAN can determine the number of clusters automatically. 3). DBSCAN is robust to noise.

The following of this chapter is organized as follows. Section 2.2.1 compares and reveals the relation of DBSCAN to a basic graph traversal algorithm **Depth First Search(DFS)**. This section will introduce the critical concept used by DBSCAN which differentiate it with DFS. Section 2.2.2 describe this concept in more details. Section 2.2.3 compares DBSCAN and K-Means, and discusses practical issues of DBSCAN.

### 2.2.1 DBSCAN and DFS

Before introducing DBSCAN, it would be helpful to review a basic graph traversal algorithm, DFS, which is an algorithm which traverses or visits a graph. In the discussion of this context, the whole process is divided into two procedures DFS and DFS-VISIT. A typical recursive implementation is described as follow.



DFS( $G$ )

```

1  for each vertex  $u \in G$ 
2       $u.visited = \text{FALSE}$ 
3  for each vertex  $u \in G$ 
4      if  $u.visited == \text{FALSE}$ 
5          DFS-VISIT( $G, u$ )

```

DFS-VISIT( $G, u$ )

```

1   $u.visited = \text{TRUE}$ 
2  for each vertex  $v \in G.Adj[u]$ 
3      if  $v.visited == \text{FALSE}$ 
4          DFS-VISIT( $G, v$ )

```

The first two lines in DFS initializes the algorithm by setting all nodes to be unvisited. Then, the algorithm traverse vertices in the graph. Once an unvisited node is found, DFS-VISIT is called. In the procedure DFS-VISIT, the given node  $u$  is labelled as visited. If an unvisited child  $v$  of  $u$  is found, this procedure goes one layer deeper by calling another DFS-VISIT on  $v$ . After finishing DFS-VISIT( $G, u$ ), the algorithm backtrack to visit other unvisited children of  $u$ , which are siblings of  $v$ .

Assume DFS executes on an undirected graph. Each call of DFS-VISIT on an unvisited node  $u$  explores  $u$  and all nodes reachable to it. These nodes form a component which disconnects with other components by other runs of DFS-VISIT. Thus, each component can be considered as a cluster. The criterion to form a cluster is that each pair of nodes can reach each other through a path consisting of nodes only in this cluster.

DBSCAN can be seen as an application of DFS. However, DBSCAN differs from standard DFS from its usage of DFS-VISIT. In standard DFS, the algorithm goes deeper by calling DFS-VISIT( $G, v$ ) on all unvisited children of node  $u$ . In DBSCAN, however, the algorithm goes deeper if and only if node  $u$  satisfies extra conditions, which are specified by two user given value  $Eps$  and  $MinPts$ . These extra conditions make the component generated from DBSCAN a meaningful cluster viewing from the density side. The pseudo code is shown on next page.

DBSCAN( $S, Eps, MinPts$ )

```

1   $ClusterId = 1$ 
2  for each point  $u \in S$ 
3      if  $u.label \neq \text{NIL}$ 
4          continue
5       $u.neighbor = \text{REGION-QUERY}(u, Eps)$ 
6      if  $\|u.neighbor\| < MinPts$ 
7           $u.label = \text{NOISE}$ 
8      else
9           $\text{EXPAND-CLUSTER}(u, ClusterId, Eps, MinPts)$ 
10      $ClusterId = ClusterId + 1$ 

```

EXPAND-CLUSTER( $u, ClusterId, Eps, MinPts$ )

```

1   $u.label = ClusterId$ 
2  for each point  $v \in u.neighbor$ 
3      if  $v.label == \text{NIL}$ 
4           $v.label = ClusterId$ 
5           $v.neighbor = \text{REGION-QUERY}(v, Eps)$ 
6          if  $\|v.neighbor\| < MinPts$ 
7               $\text{EXPAND-CLUSTER}(v, ClusterId, Eps, MinPts)$ 
8      elseif  $v.label == \text{NOISE}$ 
9           $v.label = ClusterId$ 

```

Similar with DFS traverse, DBSCAN clustering is also divided into two procedures, DBSCAN and EXPAND-CLUSTER. DBSCAN functions as a wrapper function in the same way as DFS. This procedure goes through each point in the data set  $S$ . If the current point  $u$  has been visited, the procedure skips it. Otherwise, further process continues. However, unlike calling DFS-VISIT on every unvisited  $u$  unconditionally as DFS does, in DBSCAN, another procedure EXPAND-CLUSTER is called if and only if  $u$  has sufficient number of neighbours in a given range  $Eps$ .  $\text{REGION-QUERY}(u, Eps)$  returns all the neighbours of  $u$  with a distance no further than  $Eps$ . This is the extra condition mentioned earlier. Similarly, this condition is also checked in EXPAND-CLUSTER at line 6, as opposed to DFS-VISIT. Besides these two places, the rest of the two algorithms are the same. Details of extra conditions are fully explained in next section.

### 2.2.2 Notions of Density in Clusters

From the view of density, points in a space can be classified into three categories: core points, border points, and outliers/noise. The classification

criterion on a point  $u$  bases on the size of its *Eps-neighborhood*, denoted as  $N(u; Eps)$ . The  $N(u; Eps)$  represents the collection of all points whose distances to  $u$  are within  $Eps$ , which is specified by user. More formally,  $N(u; Eps) = \{v \mid dist(u, v) < Eps\}$ .

Based on above definition, a point  $u$  is a core point if and only if  $MinPts \leq \|N(u; Eps)\|$ , where both  $Eps$  and  $MinPts$  are user specified. After defining core points, it is relatively easy to define the other two categories. A border point  $v$  is a neighbour point of a core point  $u$ , but  $v$  is not a core point itself. Formally,  $\|N(u; Eps)\| < MinPts$ , and  $v \in N(u; Eps)$ , where  $\|N(u; Eps)\| \geq MinPts$ . For a outlier, it's a point  $v$  which is neither a core point itself nor a neighbour point of a core point. An example graph is showing in Fig 2.3.

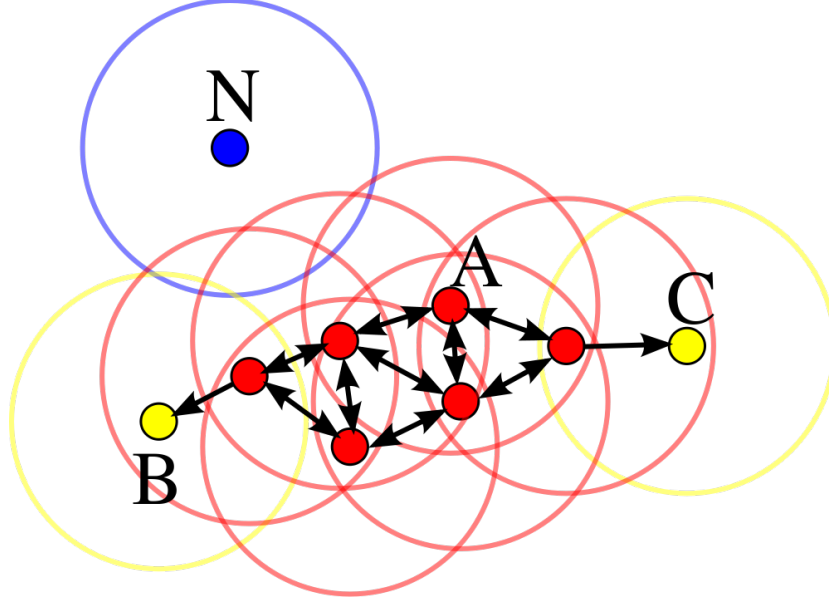


Figure 2.3: Illustration of core points(red), border points(yellow), and outlier(blue) [24]. In the illustration,  $MinPts = 4$ ,  $Eps$  is the radius of those circles.

As show in the figure, red points are core points since in each red circle, there are at least 4 points(including the center point). While for the yellow point, they don't have enough points in their circles, but they embed in one of the red circles. Thus, they are border points. As for the blue point, it has neither enough point in the blue circle nor is in any red circle. Thus, the blue point an outlier/noise.

These definitions reveal the assumption and theory of DBSCAN. DB-

SCAN assumes that, each cluster consists of two types of points, core points and border points. Core points form inner denser areas compared to spaces outside of clusters. Border points form a slightly sparser areas surrounding the inner core areas, but still the density is higher than spaces between different clusters, possessed by noise. In DBSCAN, the density of an area is measured by the number of points centred at  $u$  within an area specified by the given radius  $Eps$ . If there are at least  $MinPts$  points, this area is considered to be a dense area. Thus, this area is part of a cluster. In the procedure EXPAND-CLUSTER, the algorithm will explore the whole dense area and less dense border area belonging to the same cluster. Each call of EXPAND-CLUSTER visits a different cluster. Thus, when the whole procedure DBSCAN finishes, different clusters are formed.

### 2.2.3 Analysis of DBSCAN

DBSCAN applies to a broader range of problems compared to K-Means. In K-Means, one of the obstacle is to compute the mean value. DBSCAN doesn't have this problem as long as the distance between points is computable. This enables DBSCAN to use more complicated distance measurement, such as edit distance, which is very useful in dealing with strings. Another advantage of DBSCAN is its robustness to noises. As mentioned in Section 2.1.2, the mean value will be heavily interfere by noises/outliers. In contrast, DBSCAN has a notion noises and it is able to spot these outliers, which is a core problem of this thesis. Finally, the most important advantage of DBSCAN is that it can find arbitrarily shaped clusters. It can even separate a cluster completely surrounded by a different cluster. An example is shown in Figure 2.4.

Similar to K-Means, DBSCAN also requires user specified parameters. To get the optimal result, a careful choice of these parameters is needed. Finding the appropriate parameter can be achieved in the similar way by finding the elbow point as mentioned in Section 2.1.2. The user can pick a number for  $MinPts$  first. Then, for each point, the distance from its  $k$ th nearest neighbour is computed. After sorting these distances in descending order and plot them, a graph called *sortedk-distgraph* can be obtained. This graph reveals insights about the density distribution of the whole dataset reflected in how the  $k$ -dist varies. Then the  $Eps$  can be set to the value corresponding to the elbow point. This heuristic works well as the graph won't differ too much for  $k > 4$ . An illustration of this approach is shown in Figure 2.5. [8]

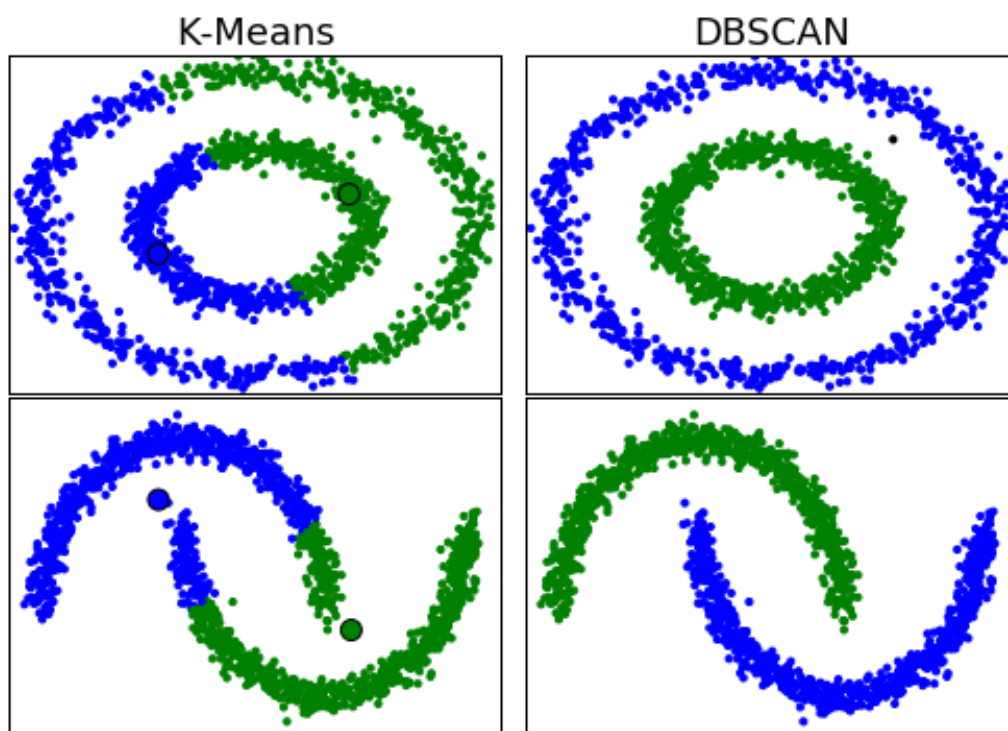


Figure 2.4: A Comparison of K-Means and DBSCAN. In the diagram, two datasets “circles” and “moons” are used. Compared to K-Means, DBSCAN gives more reasonable clustering results on these two irregular shaped datasets.

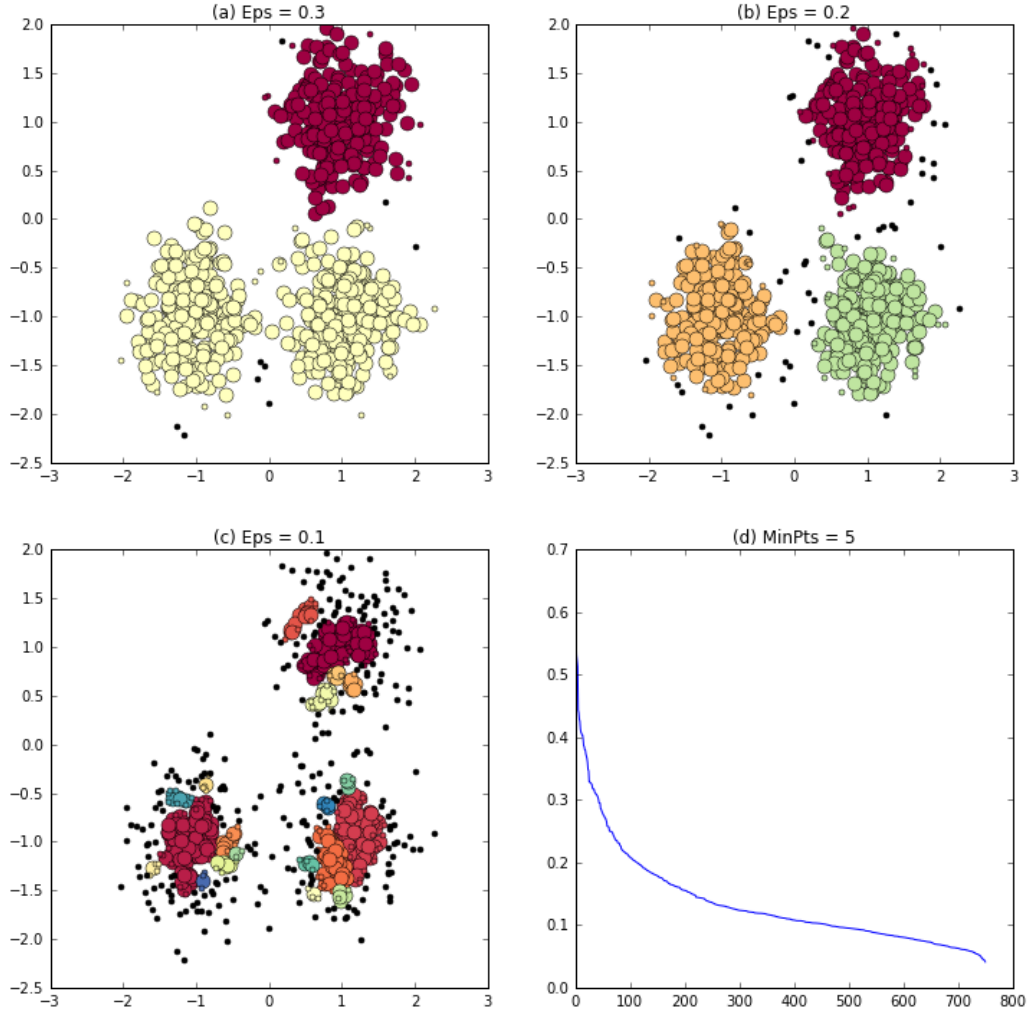


Figure 2.5: Illustration of the heuristic parameter choosing approach for DBSCAN. (a)-(c) illustrates clustering results obtained by using different  $Eps$ . Each cluster is painted with a different color. Core points are represented using larger circles while border points are represented using smaller circles. Black small circles represent outliers. In (a), very few points are labelled as noise and the bottom two clusters are not distinguished. In (b), the result is more reasonable and can be considered as optimal. In (c), too many points are labelled as noise and more than 3 clusters are reported. According to the  $k$ -distgraph shown in (d), 0.2 should be the best value for  $Eps$ , which corresponds to the result in (b).

## Chapter 3

# Scoring Patient Visits by Markov Models

This chapter is going to explore generative methods for anomaly detection. Generative methods are a collection of algorithms which try to build a model that explains the process of how the data generates. Then, the model gives a score indicating how likely one entry is an anomaly. A family of algorithms belonging to this category is the Markov models [15].

The patients visits can be seen as time sequential data consisting of a series of events. The events in one visit are not generated independently and randomly. Instead, past events have an effect on the type of the next possible event. To handle sequential data, Markov models are the correct choice since they consider the relation between consecutive observations. In the following context, Section 3.1 introduces the basic Markov chain model. Later, Section 3.2 expands the Markov chain to a more complicated Hidden Markov Model by introducing hidden variables.

### 3.1 Discrete Markov Process

Consider a system having  $K$  distinct states  $\{S_1, S_2, \dots, S_K\}$ . At any time, the system will be in one of these states. After a given time period  $N$ , a series observation  $\{x_1, x_2, \dots, x_N\}$  can be obtained. (Without loss of generality, the following discussion assumes the variables are all scalar. The assumption holds in the rest of the context unless explicitly stated otherwise) According to the product rule of probability, the joint probability distribution for this

sequence of observations is

$$p(x_1, x_2, \dots, x_T) = \prod_{n=2}^N p(x_n \mid x_1, \dots, x_{n-1}) \quad (3.1)$$

The conditional probability distribution of each observation  $x_n$  depends on all observations having a smaller index than it. The above relations between the observations can be represented graphically in Figure 3.1(a). The graph is fully connected, and no independence property can be obtained from it. Now assume that each observation  $x_n$  only depends on one immediate previous observation  $x_{n-1}$ . Then the joint distribution becomes

$$p(x_1, x_2, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n \mid x_{n-1}) \quad (3.2)$$

This newly obtained model is depicted in Figure 3.1(b), and is referred as *first-order Markov chain*. The term *first-order* indicates the dependence on only one previous observation. Suppose the system has only 3 states, as shown in Figure 3.1(c). Then, to fully represent the system, the only required information is the transition probabilities between different states. The transition probabilities are usually referred as *transition matrix*, denoted as  $\mathbf{A}$ . Each element  $A_{ij}$  represents the probability of transferring from state  $s_i$  to state  $s_j$ . Learning the parameters of this model is very simple. Since the states are exactly the observations,  $A_{ij}$  can be simply obtained by compute the frequency of transferring to  $s_j$  starting from  $s_i$ . The number of free parameters in this model is  $K(K - 1)$ , where  $K$  represents the number of states in the system.

Sometimes, the observations can depend on more than one observations in the past. One simple way to achieve this is creating a higher order Markov chain. By allowing each observation to depend on previous two values, a second-order Markov chain is obtained, as shown in Figure 3.2. Then the joint distribution becomes

$$p(x_1, x_2, \dots, x_N) = p(x_1)p(x_2 \mid x_1) \prod_{n=3}^N p(x_n \mid x_{n-1}, x_{n-2}) \quad (3.3)$$

Using the same state space representation, the *second-order Markov chain* has better capability of modelling complex relations between variables, compared to *first-order Markov chain*. In fact, the higher the order is, the more flexible the model is. However, the number of parameters grows as well, which makes the model difficult to train. For a  $M^{th}$ -order Markov Chain, there will be  $K^M(K - 1)$  parameters. Because the exponential growth in number of parameters, the model gradually becomes impractical as  $M$  grows.



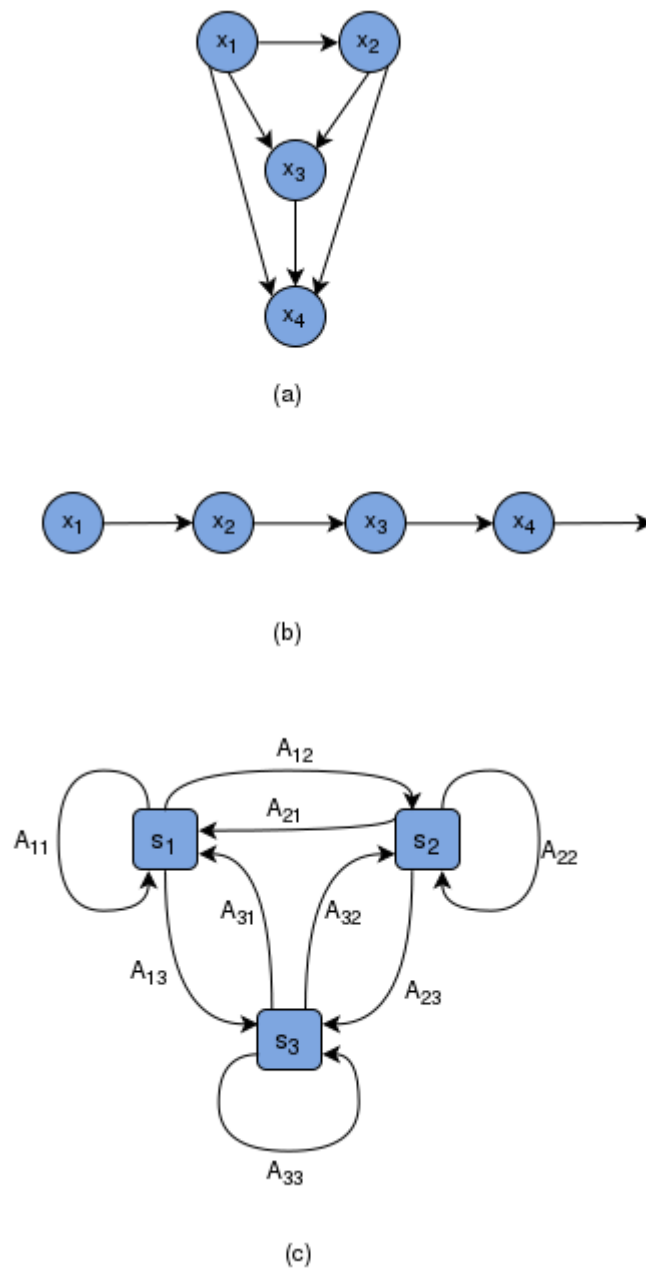


Figure 3.1: Illustration of a Markov Chain of 4 observations possessing 3 states. Variables are represented using filled circles, while states are represented using filled squares.

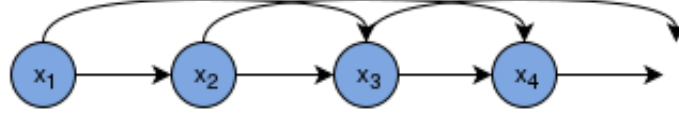


Figure 3.2: Illustration of a second-order Markov model.

### 3.2 Hidden Markov Model

The simple Markov Chain model is not enough for modelling the patient visit sequence. The variables  $\{x_1, x_2, \dots, x_N\}$  can be considered as the patient states, namely **ENROLLING**, etc. However, the visit sequence also contains time part. To integrate the time part into the model, the Markov Chain model can be expanded in another way, by associating an emission distribution  $\mathbf{E}_k$ ,  $k = 1, \dots, K$ , to each state in the system. Thus, two observations  $x_n, y_n$  exist at any time, where  $y_n$  is generated depending on  $x_n$ . If the relation between  $\{x_1, x_2, \dots, x_N\}$  is modelled as a first-order Markov chain, the joint distribution becomes

$$p(x_1, y_1, \dots, x_N, y_N) = p(x_1) \prod_{n=2}^N p(x_n | x_{n-1}) \prod_{n=1}^N p(y_n | x_n) \quad (3.4)$$

where  $x_t$  represents the patient state and  $y_t$  represents the associated duration. For each visit, the patient goes through a series of events, which is the patient state. Each event will then last for a certain period of time. The duration can be seen as generated from a distribution, and the parameters of this distribution depend on the event. One example is shown in Figure 3.3. This model is in fact a special case of Hidden Markov Model. This section explores Hidden Markov Model in details.

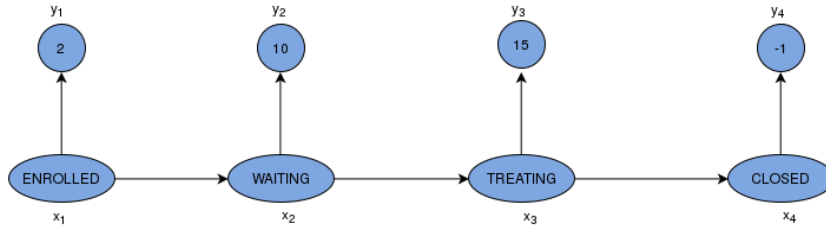


Figure 3.3: Modelling a patient visit as a special case of Hidden Markov Model.

### 3.2.1 Definition of Hidden Markov Model

As mentioned in the last section, a trade-off between flexibility and practicality exists in deciding the order number for a Markov chain model. It would be ideal if a model is not limited to any specific given order, and still only limited number of parameters are required to specify the model. Luckily, these requirements can be satisfied by constructing a Hidden Markov Model using additional latent variables [3].

Suppose a sequence of observations  $\mathbf{X} = \{x_1, \dots, x_N\}$  is obtained. Instead of assuming each observation depends directly on a specific number of previous observations, the new assumption is that, there is a latent variable  $z_t$  corresponding to each observation, and the latent variables form a Markov chain. The latent variables don't have to possess any physical meanings. They can even be of different type to the observations, in terms of distribution and dimensionality. A graphical representation of this model is shown in Figure 3.4. It's easy to get confused by comparing Figure 3.3 and Figure 3.4 since they share the same graphical structure. The difference is that, in Figure 3.4, the  $z_t$ 's are unobserved latent variables, which is depicted using unfilled circles. While in Figure 3.3, both events and duration are observed values. All observed variables are represented using filled circles. Despite the fact there are no unobserved variables in Figure 3.3, it still belongs to HMM. In this model, we are just lucky to have the luxury to obtain all variables. It is possible to add additional latent variables into this model. One potential structure could be the one shown in Figure 3.5. Intuitively, in this model, the value of the newly added latent variable determines which event will generate, then the event determines how long the duration will be. Notice that the latent variables don't have any associated physical meaning or specific distribution form. One can explain them as indication of the functioning status of the system by selecting them to be binary variables. When  $z_t = 1$ , it indicates the queue system in the hospital is working in normal mode. When  $z_t = 0$ , it means the system is working in a problematic way.

In the framework of HMM, the joint distribution over both observed and latent variables is given below

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = p(z_1 | \pi) \prod_{n=2}^N p(z_n | z_{n-1}, \mathbf{A}) \prod_{m=1}^N p(x_m | z_m, \phi) \quad (3.5)$$

where  $\mathbf{X} = \{x_1, \dots, x_N\}$  represents all the observed variables,  $\mathbf{Z} = \{z_1, \dots, z_N\}$  represents latent variables, and  $\boldsymbol{\theta} = \{\pi, \mathbf{A}, \phi\}$  represents the parameters in this model. The  $\pi$  is a prior distribution for deciding the value of the first

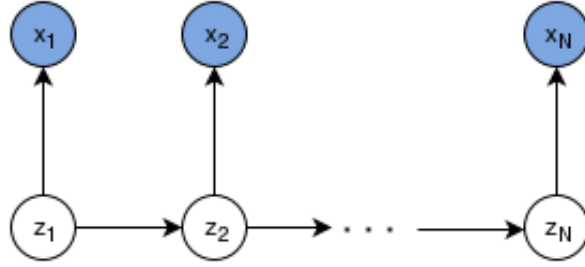


Figure 3.4: Graphical representation of a Hidden Markov Model. Observations are represented using filled circles, while latent variables are depicted using unfilled circles. The latent variables form a first-order Markov chain.

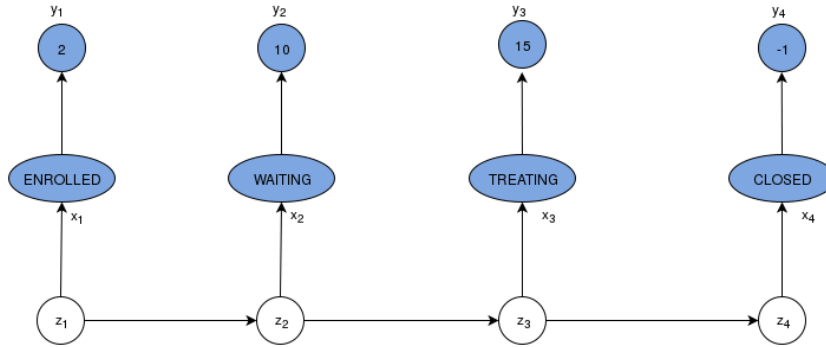


Figure 3.5: Modelling patient visit as Hidden Markov Model. Event types and event duration are represented using filled eclipses and circles respectively. Additional hidden variables are represented using unfilled circles. No specific physical meaning is associated with these latent variables.

variable  $z_1$ . The matrix  $\mathbf{A}$  is the transition matrix among the latent variables. The  $\phi$  are the parameters of the emission distribution associated with  $z_t$  and  $x_t$ .

### 3.2.2 Learning and Inference

There are three basic problems in HMM. [23] These problems are described below using above notations:

- Problem 1: Given a sequence of observations  $\mathbf{X} = \{x_1, \dots, x_N\}$ , what is the probability  $p(\mathbf{X}|\boldsymbol{\theta})$  over the observations, under specific parameters  $\boldsymbol{\theta} = \{\pi, \mathbf{A}, \phi\}$ ?

- Problem 2: What's the value of the parameters which maximizes  $p(\mathbf{X}|\boldsymbol{\theta})$ ?
- Problem 3: Given a sequence of observations  $\mathbf{X} = \{x_1, \dots, x_N\}$ , what is the value of the corresponding latent variables?

If luxury of observing the latent variable is available, these three problems becomes trivial. But if we do not have this luxury, these three problems become complicated. The rest of the context focus on the first two questions, assuming the latent variable are unobservable. The reason is that, once the value of  $p(\mathbf{X}|\boldsymbol{\theta})$  is computed, the decision on whether a given sequence is anomaly can be made by comparing  $p(\mathbf{X}|\boldsymbol{\theta})$  to a threshold value.

Though it seems more intuitive that finding a way to evaluate  $p(\mathbf{X}|\boldsymbol{\theta})$  should come before maximizing it with respect to the parameters, it would be more convenient to start at solving problem 2. After solving problem 2, the solution to the first problem will appear naturally. The following discussion begins by introducing some new concepts and notations.

The distribution over only observed variables  $p(\mathbf{X}|\boldsymbol{\theta})$  is usually referred as *incomplete likelihood*, while distribution over both observed and unobserved variables  $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$  is referred as *complete likelihood*. Using Equation 3.5, the logarithm of incomplete likelihood can be represented as

$$\begin{aligned} \ln p(\mathbf{X}|\boldsymbol{\theta}) &= \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \\ &= \ln p(z_1|\pi) + \ln \sum_{\mathbf{Z}} \left( \prod_{n=2}^N p(z_n | z_{n-1}, \mathbf{A}) \prod_{m=1}^N p(x_m | z_m, \phi) \right) \end{aligned} \quad (3.6)$$

The above equation is a generalization of the *mixture distribution* [3]. Maximizing  $\ln p(\mathbf{X}|\boldsymbol{\theta})$  with respect to the parameters is very difficult since the derivatives don't have a closed form. An alternative practical working algorithm is the *expectation-maximization(EM)* algorithm [7][21]. The EM algorithm is very similar to the K-Means algorithm mentioned in Chapter 2. The algorithm consists of two steps, E-step and M-step. In the E-step, the algorithm fixes the value of parameters and find the posterior distribution of the latent variables  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$ . Here the notation is adopted from Bishop's [3]. The superscription *old* in  $\boldsymbol{\theta}^{old}$  means the parameter is fixed. Then the algorithm computes the expectation of the logarithm of the complete likelihood, with respect to the derived posterior distribution. The newly derived term becomes a function of  $\boldsymbol{\theta}$ , which is shown below

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \quad (3.7)$$

Then in the M-step, the new value of  $\theta$  is updated by maximizing  $Q(\theta, \theta^{old})$ . Compared to K-Means, the E-step corresponds to assign each point to a cluster prototype, and the M-step corresponds to update the value of the prototypes. These two steps are executed alternatively until convergence or maximum number of iteration is reached. In the following text,  $\gamma(\mathbf{z}_n)$  and  $\gamma(\mathbf{z}_{n-1}, \mathbf{z}_n)$  are introduced which stand for the posterior distribution of a single latent variable and the joint posterior distribution over two consecutive latent variables, separately. Instead of assuming the latent variables are scalar, here they are represented using *1-of-K* coding. Namely, each latent variable is a length  $K$  vector, where one and only one of these  $K$  elements equals 1. When  $z_{nk} = 1$ , it means the  $n$ th latent variable is in the  $k$ th state. Using this representation schema, following equations are obtained

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \theta^{old}) \quad (3.8)$$

$$\gamma(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \theta^{old}) \quad (3.9)$$

$$\begin{aligned} Q(\theta, \theta^{old}) = & \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \gamma(\mathbf{z}_{n-1}, \mathbf{z}_n) \ln A_{jk} \\ & + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(x_n | \phi_k) \end{aligned} \quad (3.10)$$

Computation in the M-step is relatively easy. Assume the E-step has been done, so that  $\gamma(\mathbf{z}_n)$  and  $\gamma(\mathbf{z}_{n-1}, \mathbf{z}_n)$  are like constants now. Then following update equation can be obtained

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} \quad (3.11)$$

$$A_{jk} = \frac{\sum_{n=2}^N \gamma(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \gamma(z_{n-1,j}, z_{nl})} \quad (3.12)$$

Update of  $\phi_k$  is more tricky, since it depends on the specific choice of the emission distribution. One good observation is that, only the final term depends on  $\phi_k$ , and different  $\phi_k$  doesn't couple with each other. Thus, each  $\phi_k$  can be updated separately. The term  $\gamma(z_{nk})$  functions as a soft assignment, representing the probability of assigning a point  $x_n$  to each state.

Computation in E-step is more difficult which requires efficient algorithm. The most widely used algorithm is known as *alpha-beta* algorithm. This algorithm can be seen as an application of dynamic programming technique which takes advantage of the tree structure in HMM thus leading to efficiency. To start with the alpha-beta algorithm, following conditional independence

properties should be obtained first [16]

$$p(\mathbf{X}|\mathbf{z}_n) = p(x_1, \dots, x_n|\mathbf{z}_n) \quad (3.13)$$

$$p(x_1, \dots, x_{n-1}|x_n, \mathbf{z}_n) = p(x_1, \dots, x_{n-1}|\mathbf{z}_n) \quad (3.14)$$

$$p(x_1, \dots, x_{n-1}|z_{n-1}, \mathbf{z}_n) = p(x_1, \dots, x_{n-1}|\mathbf{z}_{n-1}) \quad (3.15)$$

These equations can be obtained by using *d-seperation* technique [22], or proved formally using sum and product rules of probability. Using the first independence property and Bayes' theorem, following equations are obtained

$$\begin{aligned} \gamma(\mathbf{z}_n) &= p(\mathbf{z}_n|\mathbf{X}) = \frac{p(\mathbf{X}|\mathbf{z}_n)p(\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{p(x_1, \dots, x_n, \mathbf{z}_n)p(x_{n+1}, \dots, x_N|\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})} \end{aligned} \quad (3.16)$$

where

$$\alpha(\mathbf{z}_n) = p(x_1, \dots, x_n, \mathbf{z}_n) \quad (3.17)$$

$$\beta(\mathbf{z}_n) = p(x_{n+1}, \dots, x_N|\mathbf{z}_n) \quad (3.18)$$

Using the other two conditional independence properties,  $\alpha(\mathbf{z}_n)$  can be expressed recursively in terms of  $\alpha(\mathbf{z}_{n-1})$

$$\begin{aligned} \alpha(\mathbf{z}_n) &= p(x_n|\mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(x_1, \dots, x_{n-1}, \mathbf{z}_{n-1})p(\mathbf{z}_n|\mathbf{z}_{n-1}) \\ &= p(x_n|\mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1})p(\mathbf{z}_n|\mathbf{z}_{n-1}) \end{aligned} \quad (3.19)$$

Similarly,  $\beta(\mathbf{z}_n)$  can also be expressed recursively as

$$\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1})p(x_{n+1}|\mathbf{z}_{n+1})p(\mathbf{z}_{n+1}|\mathbf{z}_n) \quad (3.20)$$

The term  $\alpha(\mathbf{z}_n)$  can be seen as messages propagated from the beginning to the end. Each  $\alpha(\mathbf{z}_n)$  receives messages passed from its predecessor, combines these information with its own information and then pass them to its successor. The logical also applies to the term  $\beta(\mathbf{z}_n)$ , but the messages are from the end to the beginning. Due to the tree structure in HMM, computing each term only depends on one adjacent term, instead of all terms before/after

it. Thus, the computation reduces dramatically which makes the algorithm efficient. To start the whole computation, initial conditions  $\alpha(\mathbf{z}_1)$  and  $\beta(\mathbf{z}_n)$  are required. The initial conditions are given below

$$\alpha(\mathbf{z}_1) = \prod_{k=1}^K \{\pi_k p(x_1 | \phi_k)\}^{z_{1k}} \quad (3.21)$$

$$\beta(\mathbf{z}_N) = 1 \quad (3.22)$$

Having obtained  $\alpha(\mathbf{z}_n)$  and  $\beta(\mathbf{z}_n)$ , the posterior distribution  $\gamma(\mathbf{z}_n)$  can be computed as in equation (3.16). As for  $\gamma(\mathbf{z}_{n-1}, \mathbf{z}_n)$ , it can be computed as following

$$\begin{aligned} \gamma(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}) \\ &= \frac{p(\mathbf{X} | \mathbf{z}_{n-1}, \mathbf{z}_n) p(\mathbf{z}_{n-1}, \mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{p(x_1, \dots, x_{n-1} | \mathbf{z}_{n-1}) p(x_n | \mathbf{z}_n) p(x_{n+1}, \dots, x_N | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1})}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_{n-1}) p(x_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \beta(\mathbf{z}_n)}{p(\mathbf{X})} \end{aligned} \quad (3.23)$$

Up till now, both steps in EM algorithm are introduced, and the problem 2 can be solved efficiently. The left question is how to solve problem 1, computing the likelihood over the incomplete data. The solution comes from Equation 3.16. Notice that  $\gamma(\mathbf{z}_n)$  is a posterior distribution. Integrating both sides of Equation 3.16 over  $\mathbf{z}_n$  gives

$$p(\mathbf{X}) = \sum_{\mathbf{z}_n} \alpha(\mathbf{z}_n) \beta(\mathbf{z}_n) \quad (3.24)$$

where  $\mathbf{z}_n$  is an arbitrary latent variable. If  $n = N$ , then  $\beta(\mathbf{z}_n) = 1$ , which makes the above equation simpler

$$p(\mathbf{X}) = \sum_{\mathbf{z}_N} \alpha(\mathbf{z}_N) \quad (3.25)$$

Then, both problem 1 and problem 2 are solved.



## Chapter 4

# Experiments

This chapter describes the detailed experimental design and implementation. Following topics will be discussed:

- Which specific data are used? How to represent the data? What pre-process has been done? What related functions are defined?
- Which methods are used? What’s the reason to use such methods?

The three topics will be discussed in three sections, respectively. Analysis of the result will be postponed to Chapter 5.

### 4.1 Data Details and Representation

Since the X-Akseli system has gone through several update after its first release, many features have changed, including how the visit log is recorded. Considering this aspect, only data generated after year 2014 is used in the experiment. All the data are from Oulu Hospital, with patient privacy information eliminated. In total, 243K unique visits with 1.93 million entries are retrieved, which spans from January to July. The retrieved data consists of three columns: visit id, event type, and recorded time of the event. Other information, such as which resource generated the data, is not retrieved. Among the retrieved data, 7 event types are adopted. The used event types are: `ENROLLING`, `WAITING`, `IN_TREATMENT_ROOM`, `PAUSED`, `IN_TREATMENT_ROOM_FROM_PAUSED`, `CLOSED`, and `CANCELLED`.

Huang et al [13] proposed to represent the patient visits as a sequence of pairs. Each pair contains the information of (a). what the event is and (b). when this event happened. For example, given a sequence showing below

$$\langle (a, 1), (b, 2), (c, 5), (d, 7) \rangle$$

it means a patient comes to the hospital, and at time 1, the patient encounters event  $a$ . Then at time 2, the patient encounters event  $b$ , and so on. The time unit can be selected arbitrarily, such as minutes, hours, or days. Huang et al. [13] used days as the time unit. In their work, they tried to cluster the patient traces. However, the patient trace has different length. Thus, typical distance metrics such as euclidean distance is not applicable. To address this problem, they also proposed a new distance metric, based on edit distance.

Edit distance [5] is commonly used in comparing strings and biological sequences, such as proteins. Edit distance is defined as the minimum number of allowed operations used, to transform a string  $s$  to another string  $t$ . For example, if the allowed operations are *delete*, *insert*, and two strings  $S = \text{"array"}$ ,  $T = \text{"xray"}$  are given. Then the edit distance between  $s$  and  $t$  is 3, by taking 3 operations. One potential transformation is:

1. Delete the second letter  $r$  in  $S$  by  $x$ . Then  $S$  becomes *"aray"*.
2. Delete the first letter  $a$  in  $S$ . Then  $S$  becomes *"ray"*.
3. Insert letter  $x$  at the beginning of  $S$ . Then  $S$  becomes *"xray"*, which is the same with string  $T$ .

The edit distance problem can be solved effectively by using the dynamic programming technique [5]. Using terminology from dynamic programming, the optimal solution to the edit distance problem can be represented recursively

$$D(i, j) = \begin{cases} D(i-1, j-1) & \text{if } S[i] = T[j] \\ \min\{D(i-1, j), D(i, j-1)\} + 1 & \text{if } S[i] \neq T[j] \end{cases}$$

The edit distance only considers the difference between types of events, when applied to the patience visit data. However, the time associated with each event should also makes an effect when comparing two traces. Huang et al. [13] addressed this problem by providing a modified edit distance. In the old edit distance, events from two patient trace will either increase the distance by 1 if they belong to different type or 0 if they belong to the same type. In the modified distance, however, the increment caused by two events range from  $[0, 1]$  as shown below

$$\delta(\sigma_i, \rho_j) = \begin{cases} 1 & \text{if } \sigma_i(e) \neq \rho_j(e) \\ \frac{|\sigma_i(t) - \rho_j(t)|}{\max\{\sigma_i(t), \rho_j(t)\}} & \text{if } \sigma_i(e) = \rho_j(e) \end{cases}$$

where  $\sigma$  and  $\rho$  are two patient traces.  $\sigma_i$  is the  $i$ th pair of the trace.  $\sigma_i(e)$  and  $\sigma_i(t)$  represent the event type and timestamps of the  $i$ th pair in that trace.

The intuition of the above equation is that, if the event types of two pairs in two traces are different, then they contribute 1 to the edit distance. If the event types are the same, then the distance is determined by the timestamps associated with the two events. The closer the timestamps are, the smaller the distance is.

The modified edit distance seems reasonable. However, some subtle issues exist when the modified edit distance applies to Huang’s [13] representation. Consider two patient traces

$$\begin{aligned} S &= \langle (a, 1), (b, 1000), (c, 1001), (d, 1002) \rangle \\ T &= \langle (a, 1), (b, 2), (c, 3), (d, 4) \rangle \end{aligned}$$

The two traces are very similar, except that the second pair differs greatly. But this difference propagates further to the third and fourth pairs, incurring more penalty. The modified edit distance will equal almost 3. It would be more reasonable if the distance accounts only the huge difference generated in the second pairs, and considers the third and fourth pair the same. To address this problem, in the experiment, we proposed an alternative representation form. Rather than record the absolute timestamps associated with each event, the duration of each event is recorded. Thus, the above two patient traces becomes

$$\begin{aligned} S &= \langle (a, 1), (b, 999), (c, 1), (d, 1) \rangle \\ T &= \langle (a, 1), (b, 1), (c, 1), (d, 1) \rangle \end{aligned}$$

Applying the modified edit distance to the new representation, the answer equals roughly 1, which is more intuitive. Thus, in all experiments, the second representation form is adopted. Minute is used as the time unit.

Another critical point is about preprocessing. After using the second representation form, numerous noise points are observed. It’s believed that the noise points are generated by the system itself for logging reasons. The feature of the noise is that all events have a 0 duration time. The noise points consist of approximately 20% of all data, which incurs great effect in training models. Thus, in the pre-process step, all noise data are manually removed.

## 4.2 Methods

Chapter 2 and Chapter 3 introduced 3 potential methods, K-Means, DBSCAN, and Hidden Markov Model. However, only DBSCAN and Hidden Markov Model are used in the experiment. This section explains the reasons for choosing only these two methods and describe related details.

### 4.2.1 Choice of Clustering Method

As stated in Chapter 2, both DBSCAN and K-Means need an elaborate distance metric. This requirement can be fulfilled by using the modified edit distance. Besides this, K-Means also requires efficient computation of the mean value. However, based on the current data representation, it is not clear how to compute the mean. Also, the different sequence length also makes K-Means not applicable in this setup.

One critical phase in applying DBSCAN is how to efficiently compute the *Eps-neighborhood*. The DBSCAN algorithm will go through several iterations. Each iteration involves computing *Eps-neighborhood* for all points. The efficiency of computing *Eps-neighborhood* directly determines the practicability of DBSCAN. A naive implementation is to compute a  $N$  by  $N$  pair-wise distance table. Then sort the rows in the matrix. After this, finding the *Eps-neighborhood* takes only constant time. This process takes  $O(N^2)$  space and time. The space complexity can be further reduced to  $O(N)$ .

A more efficient method, vantage-point tree [25] can be used in the experiment. Vantage-point tree is a recursively built balanced binary tree. Each non-leaf node consists of two fields, a point functions as a center and a radius, and has two children. The center point is randomly selected from a set of available points. Initially, the set contains all data points. Then, distances from the center to all the rest points in the set is computed. Next, the radius is set to the median of the distances. After this, the point set is divided into two subsets, one consists of points with a distance shorter than the radius, the other consists of points with a distance longer than the median value. The first set is stored in the left child of current node, and the second set is stored in the right child of the current node. Intuitively, this works as drawing a circle centred on the selected center point. The circle partitions the other points into two parts, with half inside the circle, half outside the circle. This process continues until a leaf node is encountered. This building process takes  $O(N \log N)$  time and  $O(N \log N)$  space.

After building the vantage-point tree, finding the *Eps-neighborhood* for a single point takes  $O(\log N)$  time. Suppose the query point is  $p$ . Now the query enters a node centered on point  $q$  and the radius of the node is  $\tau$ . Suppose the distance from  $p$  to  $q$  is  $\delta < \tau$ . Then the query explores only the left child of  $q$  if  $\delta + \epsilon \leq \tau$ . Otherwise, the right child is also explored. The intuition is that, if one *Eps-neighborhood* of query point  $p$  is exactly  $\epsilon$  far to  $p$ , then the distance between this neighbour point and the centred point  $q$  is at most  $\delta + \epsilon$ . If this distance is no larger than the radius associated to the centred node, then there is no need to explore points outside the circle, which reduces the search space by half. As a result, find *Eps-neighborhood*

for all points takes  $O(N \log N)$  time.

### 4.2.2 Choice of Generative Method

In the experiment, the simple Hidden Markov Model with observed latent variable is selected for several reasons. The reason is that, the type of next potential event closely depends on the previous one, rather than depending on an implicit status. For example, after the event **ENROLLING**, it is very likely the next event is **WAITING**. Sometimes, due to special situations such as cancellation of the reservation, the event **CANCELLED** follow. But it is impossible a **IN\_TREATMENT\_ROOM** event comes, which skips the **WAITING** phase. Another argument is that, anomaly happened in one phase does not necessarily affect the coming phase. For example, a patient may encounter problems while being in the **WAITING** phase. But this does not mean the patient may also encounter problems while in **IN\_TREATMENT\_ROOM** phase. In other words, being in abnormal status in each phase is independent, which is quite contrary to the assumption of HMM that the value of current hidden variable depends on value of the previous hidden variables. As a result, simple Hidden Markov Model with observed latent variable is selected in the experiment.

In the HMM model, the events are represented using *1-of-K* coding schema. Thus, the events are multinomial variables. An important part of the model is the choice for the emission distribution. Usually, Gaussian distribution is used for modelling continuous variables. However, waiting time is not symmetrically distributed with respect to a mean. The shortest waiting time can only be 0 and the longest waiting time can be very long. The distribution has a skew with a long tail. According to studies in queuing theory, Poisson distribution is more appropriate. Several methods exist to test if Poisson distribution applies, for example, measuring the dispersion. Before stepping into numerical computations, visualizing the distribution is a good start. Related histogram of duration distribution associated with each event is shown in Figure 4.1. Intuitively, the distributions of **ENROLLING**, **IN\_TREATMENT\_ROOM**, and **IN\_TREATMENT\_ROOM\_FROM\_PAUSED** look like Poisson distribution. The distribution of **WAITING** seems rather different, which resembles mixture of two Poisson distributions. One plausible explanation is that, these data is obtained from all departments in the hospital. It is possible different departments have different typical waiting time. Luckily, information of department can be also obtained. After extracting data from the largest department, 8915 entries remains. The histogram of these data is shown in Figure 4.2. The newly obtained data seems to be more likely from Poisson distribution. Notice that, in these data, only three events remains.

Next is to compute the dispersion to see if the data distribution really

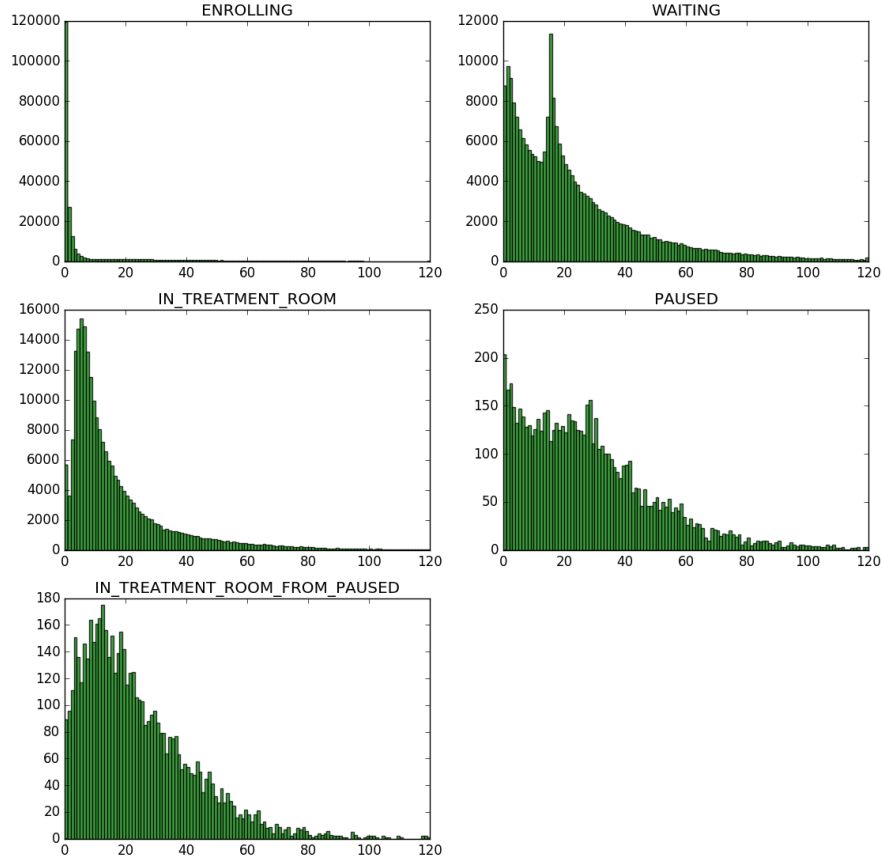


Figure 4.1: Histogram of duration associated with all events

matches Poisson distribution. The formula to compute dispersion is as follow

$$D = \frac{\sigma^2}{\mu}$$

where  $\sigma^2$  and  $\mu$  is the variance and mean value of the data, respectively. One caveat is very large duration time in each distribution. For example, the longest duration in **WAITING** can be over 1000. Though this situation is very rare, it incurs very large difference when computing  $\mu$  and  $\sigma$ . In fact, such rare values can be considered as the anomalies we are trying to detect. Thus, it makes sense to ignore these cases. It's hard to decide a threshold

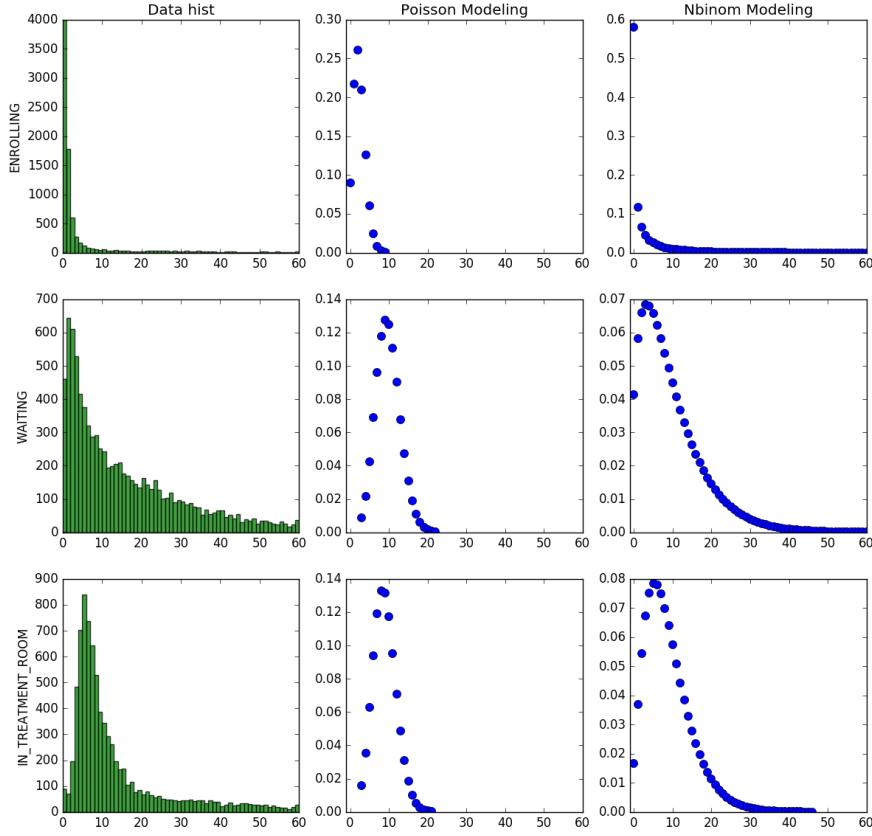


Figure 4.2: Histogram of duration associated with all events using data generated by the largest department, together with Poisson distribution fitting and negative binomial distribution fitting.

that determines which part of data should be discarded. In the experiment, the threshold is set to 30. All duration longer than this threshold is not used while computing  $\mu$  and  $\sigma$ . The reason of selecting 30 is that a typical event seldom lasts more than 30 minutes. In fact, 84% events in this data have less than 30 minutes duration. The computed values are listed in Table 4.1 When dispersion equals 1.0, it means the distribution follows Poisson distribution. Typically, data generated from Poisson distribution can have slightly larger dispersion than 1.0. However, as shown in Table 4.1, the dispersions are much larger than 1.0. This suggests that the distributions are more likely to come

Table 4.1: Mean, variance, and dispersion of data generated from largest department.

events	mean	variance	dispersion
ENROLLING	2.4	28.9	12.0
WAITING	9.8	68.0	6.93
IN_TREATMENT_ROOM	8.9	36.0	4.0

from negative binomial distribution rather than Poisson distribution. Compared to Poisson distribution, negative binomial distribution has “heavier” tails and larger variance. Considering this aspect, negative binomial distribution modelling is also implemented in the experiment. The fitting result is shown in Figure 4.2. As shown in the figure, negative binomial distribution has a better fitting of the data. Thus, in the final, Hidden Markov Model with observed latent variable with negative binomial distribution is chosen as the generative method.



## Chapter 5

# Results and Discussion

Since the data doesn't have labels indicating which entries are anomalies, precise quantitative evaluation is not possible. To examine how the methods perform, user interpretation of the data is the only standard. To help strengthen intuitive understanding, a visualization method t-SNE[19] was applied. This method only requires a distance metric between pairs of data entries. Then it projects the entries into a 2D space showing potential structures underlying the data. t-SNE typically places "important" points in the center of the drawing. The concept of "importance" can be understood as a special kind of density. Intuitively, points lying in dense area are less likely to be outliers.

Section 5.1 and 5.2 describes results obtained using clustering method and generative method respectively. Then, Section 5.3 compares these two methods.

### 5.1 Clustering Method Results

The first step of applying DBSCAN is to choose parameters. As stated in Section 2.2.3, DBSCAN is relatively robust with  $k > 4$ , where  $k$  stands for the  $k$ th nearest neighbour. In this experiment,  $k$  was set to 7. The distance of the 7th nearest neighbour of all points is drawn in Figure 5.1. The figure shows that the distance begins to increase dramatically beyond 0.5. Thus, in the experiment,  $Eps$  was finally set to 0.5, with  $MinPts = 7$ . The clustering result is shown in Table 5.1 and visualized in Figure 5.2.

As shown in the table, 12 clusters are generated by DBSCAN. Among the 12 clusters, cluster 0 is labelled as noise/anomalies which consists of 352 visits. Cluster 2 is the largest cluster which consists of 7547 visits. This cluster is believed to consist of the most typical visits. For the rest small

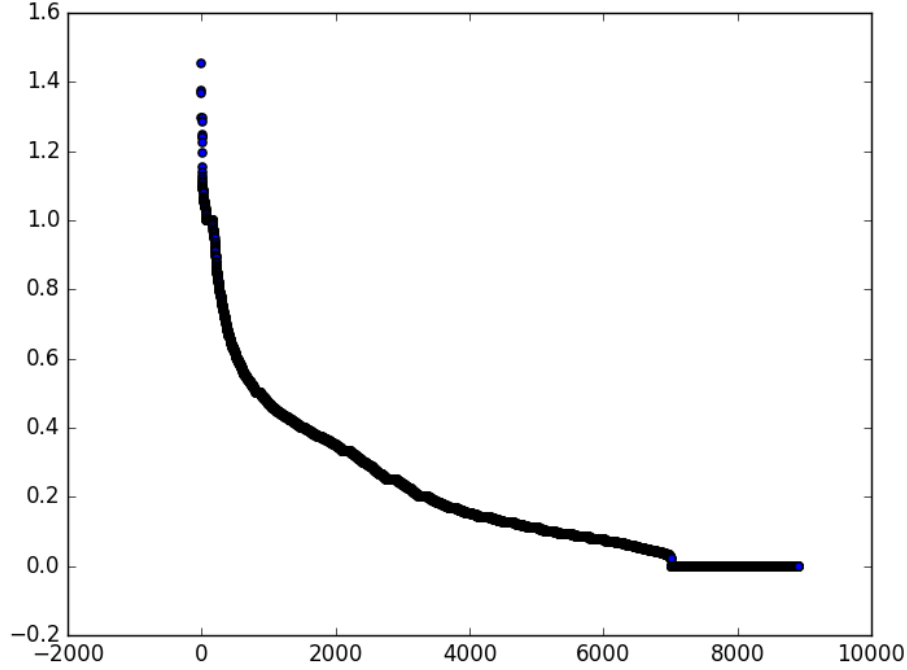


Figure 5.1: 7th nearest neighbour distance for selecting  $Eps$  of DBSCAN

Table 5.1: Mean, variance, and dispersion of data generated from largest department.

Cluster No.	0	1	2	3	4	5	6	7	8	9	10	11
Num Points	352	10	7547	134	215	283	64	104	44	47	21	94

clusters, they can be considered as representing some non-typical but normal visits. One potential reason of generating such sub-clusters is that, there are many different resources/machines for diagnosing. Data in these sub-clusters are generated from these less frequently used resources/machines. The result is visualized in Figure 5.2.

In Figure 5.3, cluster 0 is plotted using red while the rest clusters are all painted in blue. As the figure shows, many red points located on the border of the figure, which indicates they come from a sparser area in their original

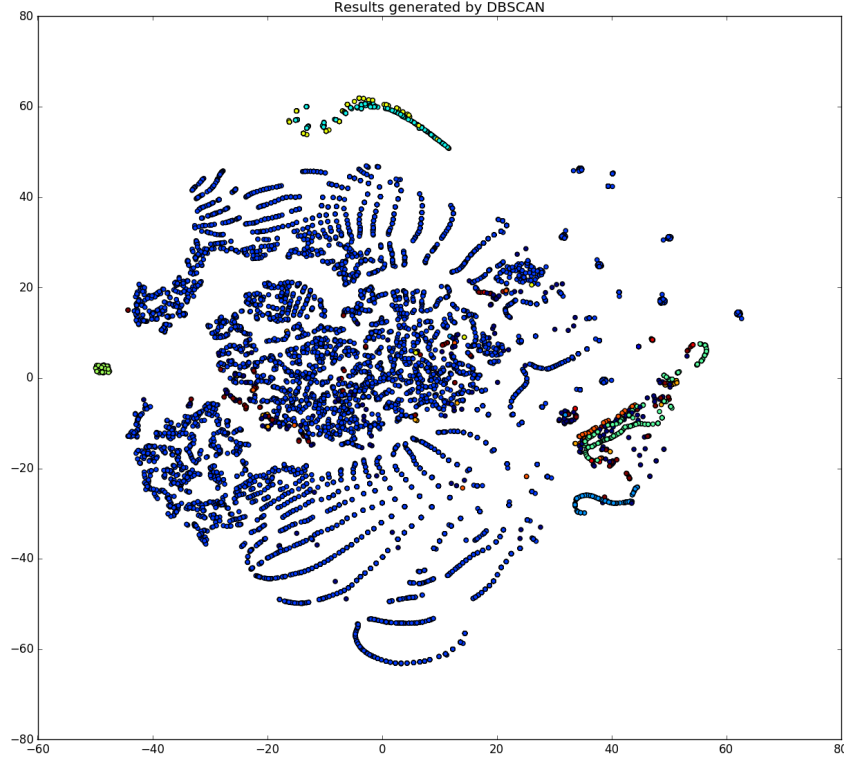


Figure 5.2: Results generated by DBSCAN, visualized using t-SNE. All clusters are shown. Each cluster is plotted in a different color.

space and have less importance. This corresponds the intuition interpretation.

To further verify the suspicion, 10 samples from each cluster are listed in Table A.1. Visits in Cluster 0 seem very uncommon. Some visits just ended without neither closed by the doctor nor cancelled by the patient. Cluster 1 consists of similar visits. Cluster 2 seems to have many reasonable visits. Visits from this cluster typically consists of four events and each event has duration no longer than 30 minutes. Thus, this cluster can be interpreted as the collection of normal visits as assumed. The rest clusters also exhibit some intuitive patterns. Some small clusters can be also considered as anomalies in addition to cluster 0, for example, cluster 10. The reason there are many clusters is that the distance between border points in two clusters are too

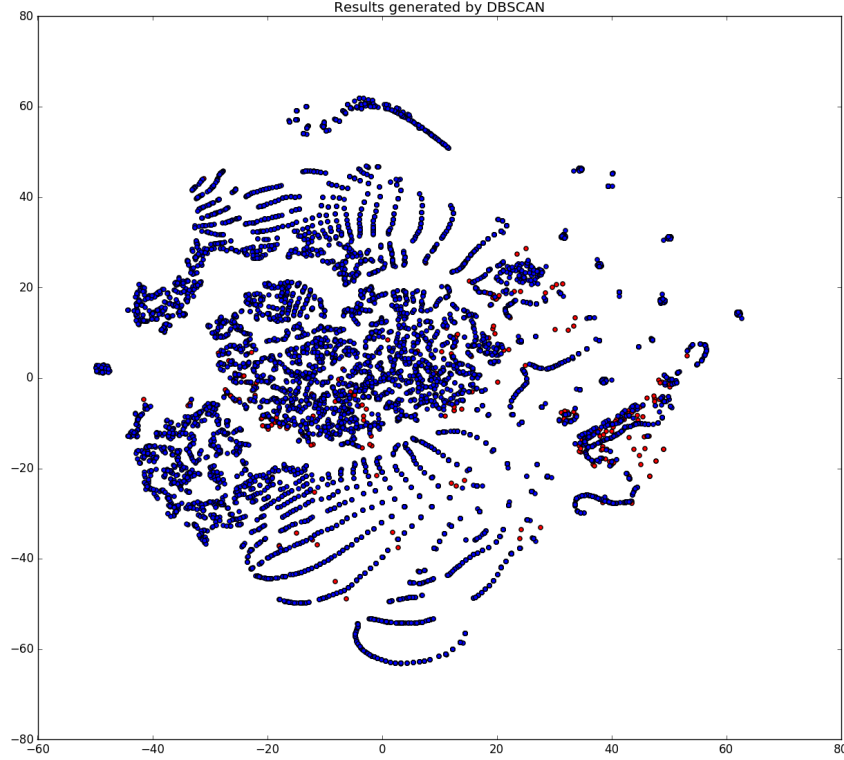


Figure 5.3: Results generated by DBSCAN, visualized using t-SNE. Cluster 0 is plotted in red, while all the rest clusters are plotted in blue.

large so that the two clusters did not merge.

## 5.2 Generative Method Results

Compared to DBSCAN, Markov Chain method is easier to implement. However, some special rules need to be set before running the method. The first rule is how to handle the -1 appeared in each event, which means the visit terminates. Since negative binomial only defines probability for non-negative inputs, the probability of duration equals -1 is undefined. In the experiment, we set this probability to be the frequency of being -1 happened in this event. So the summation of the probability of all possible values is slightly larger

than 1. But this does not incur a large effect.

Another caveat is the numerical issues while computing likelihood. Since the likelihood of a probability will typically be so small that precision problem may occur. To avoid this, the log-likelihood is computed instead. Visits having longer sequence of events tend to have smaller likelihood, but this does not mean the visit is less likely to happen. Considering this problem, the final log-likelihood is normalized by dividing the length of the sequence. The log-likelihood of visits sorted in decreasing order is shown in Figure 5.4. As shown by the figure, most visits have a log-likelihood larger than -10, the rest few visits with log-likelihood much smaller than -10 are very like to be anomalies. After selecting a threshold to be -10, the detection result is shown in Figure 5.5. Anomalies are painted in red. Again, these suspected anomalies locates on border areas in the figure.

Similarly, for every 1000 visits, 5 samples are selected with their log-likelihood listed in Table A.2 for further exploration. As listed in the table, visits in the first several blocks are very similar and seem to be very normal. It was only from the last but one block, visits start to behave in different ways. And in the last block, which consists of visits have very large negative values of log-likelihood, these visits are very bizarre and are exactly the anomalies we tried to find.

### 5.3 Discussion

Above result suggests both DBSCAN and Markov Chain can spot out anomalies. However, compared to DBSCAN, we think Markov Chain is a better method for several reasons.

Firstly, clusters formed by DBSCAN are slightly contaminated. For example, the 4th visit in cluster 7 seems very abnormal and should appear in other clusters. Other clusters also have entries does not resemble other visits in this log. A reason for such behaviour is the hard assignment to clusters in DBSCAN. In Markov Chain, however, each visit is assigned by a score which indicates how “normally” this visit is. This “soft-assignment” is a better description of the entries. Besides, the user has to interpret the meaning of each cluster by themselves, which is typically unexpected by the user.

Secondly, Markov Chain has better time and space complexity for detecting future anomalies. When determine if a new visit log is anomaly, DBSCAN will compare this new visit to all past visits and then assign this new visit to the cluster fitting it best. Thus, DBSCAN requires to maintain all past visits, and new detection takes  $O(n)$  time. This requirement will gradually becomes impractical. In contrast, Markov Chain only needs to maintain the

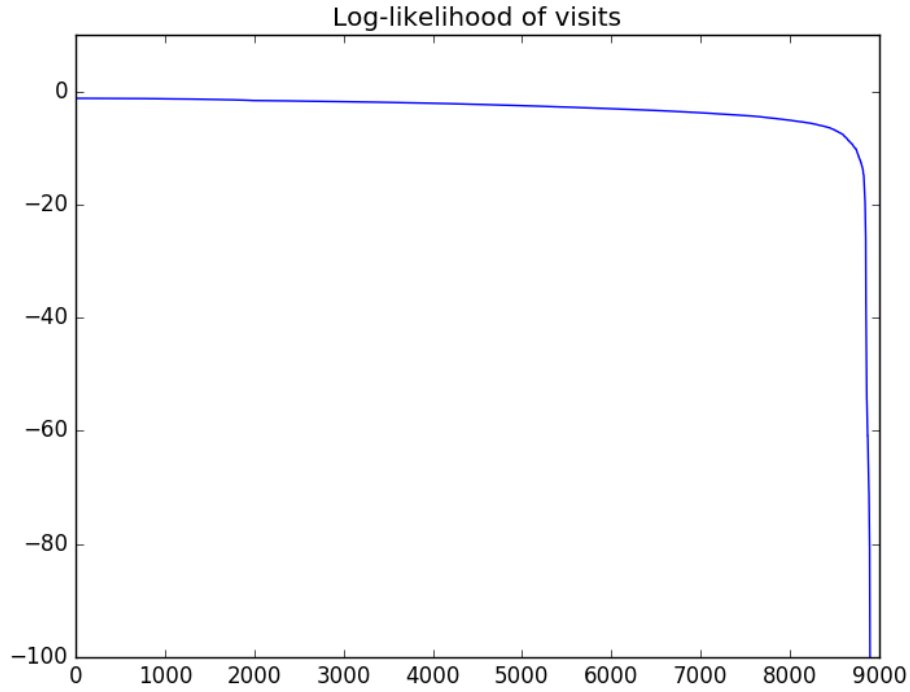


Figure 5.4: Log-likelihood of all visits sorted in decreasing order.

computed parameters of transition matrix and emission functions. Computing the log-likelihood takes  $O(1)$  time for each new visit log. Thus, speaking from this aspect, Markov Chain is a much better method than DBSCAN.

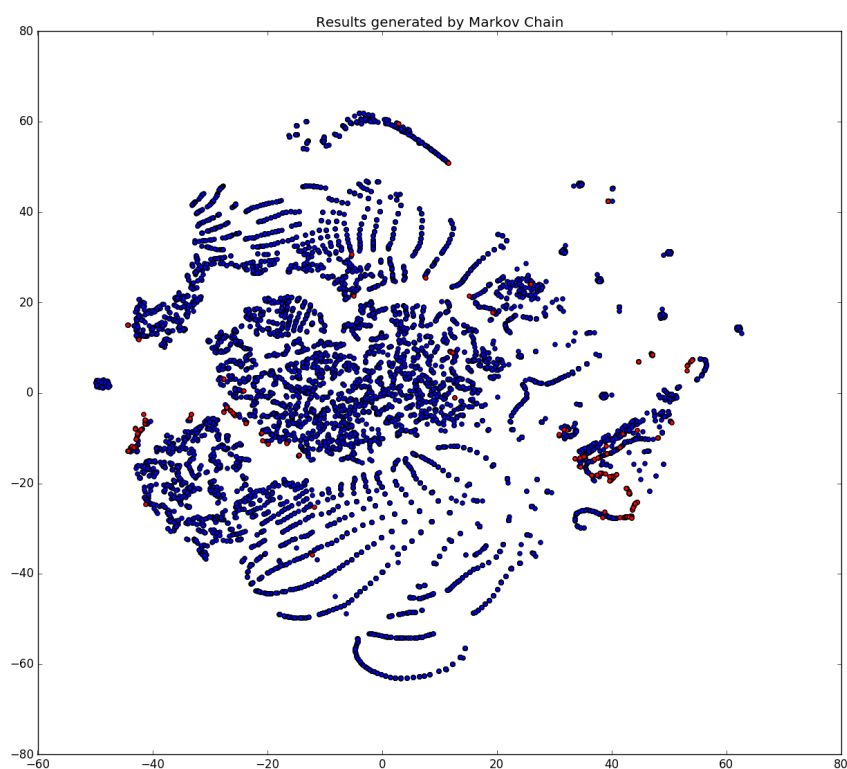


Figure 5.5: Results generated by first-order Markov Chain using negative binomial distribution as emission function, visualized using t-SNE. Anomalies are painted in red.

## Chapter 6

# Summary and Conclusion

In this thesis, we proposed a new representation for patient visit data, and explored four anomaly detection methods, K-Means, DBSCAN, Markov Chain, and Hidden Markov Model. Based on our experiments results, we suggest using Hidden Markov Model for anomaly detection from patient visit data. The model takes patient states as latent variable while takes event duration as observed variables. We further improved this method by using negative binomial to model event duration. Thanks to the luxury of observing patient states, training the model becomes easy. Once training completes, the Hidden Markov Model has the ability to score each patient in real-time and suggest potential anomalies. With information provided by the method, the hospital will be able to provide smoother visit procedure.

The thesis has considered the fact that patient visit may vary in different departments. The experiment demonstrated the variation indeed exist. It is likely the patient visit may also change in different seasons, due to reasons such as holidays. To further improve the model, factors of seasons, days can also be taken into consideration in future work.



# Bibliography

- [1] AGGARWAL, C. C. An introduction to outlier analysis. In *Outlier Analysis*. Springer, 2013, pp. 1–40.
- [2] ALOISE, D., DESHPANDE, A., HANSEN, P., AND POPAT, P. Np-hardness of euclidean sum-of-squares clustering. *Machine learning* 75, 2 (2009), 245–248.
- [3] BISHOP, C. M. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- [4] CAMPELLO, R. J., MOULAVI, D., ZIMEK, A., AND SANDER, J. Hierarchical density estimates for data clustering, visualization, and outlier detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 10, 1 (2015), 5.
- [5] CORMEN, T. H. *Introduction to algorithms*. MIT press, 2009.
- [6] DAVIS, J. C., AND SAMPSON, R. J. *Statistics and data analysis in geology*, vol. 646. Wiley New York et al., 1986.
- [7] DEMPSTER, A. P., LAIRD, N. M., AND RUBIN, D. B. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)* (1977), 1–38.
- [8] ESTER, M., KRIEGEL, H.-P., SANDER, J., XU, X., ET AL. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Kdd* (1996), vol. 96, pp. 226–231.
- [9] FORSYTH, D. A., AND PONCE, J. *Computer vision: a modern approach*. Prentice Hall Professional Technical Reference, 2002.
- [10] GUPTA, D., AND DENTON, B. Appointment scheduling in health care: Challenges and opportunities. *IIE transactions* 40, 9 (2008), 800–819.

- [11] GUPTA, M., GAO, J., AGGARWAL, C., AND HAN, J. Outlier detection for temporal data. *Synthesis Lectures on Data Mining and Knowledge Discovery* 5, 1 (2014), 1–129.
- [12] HE, Z., XU, X., AND DENG, S. Discovering cluster-based local outliers. *Pattern Recognition Letters* 24, 9 (2003), 1641–1650.
- [13] HUANG, Z., LU, X., AND DUAN, H. Anomaly detection in clinical processes. In *AMIA Annual Symposium Proceedings* (2012), vol. 2012, American Medical Informatics Association, p. 370.
- [14] HULSHOF, P. J., KORTBEEK, N., BOUCHERIE, R. J., HANS, E. W., AND BAKKER, P. J. Taxonomic classification of planning decisions in health care: a structured review of the state of the art in or/ms. *Health systems* 1, 2 (2012), 129–175.
- [15] ISAACSON, D. L., AND MADSEN, R. W. *Markov chains, theory and applications*, vol. 4. Wiley New York, 1976.
- [16] JORDAN, M. I. An introduction to probabilistic graphical models, 2003.
- [17] KDD ORG. 2014 sigkdd test of time award. webpage, 2014. <http://www.kdd.org/News/view/2014-sigkdd-test-of-time-award>. Accessed 11.8.2016.
- [18] LLOYD, S. Least squares quantization in pcm. *IEEE transactions on information theory* 28, 2 (1982), 129–137.
- [19] MAATEN, L. V. D., AND HINTON, G. Visualizing data using t-sne. *Journal of Machine Learning Research* 9, Nov (2008), 2579–2605.
- [20] MACQUEEN, J., ET AL. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability* (1967), vol. 1, Oakland, CA, USA., pp. 281–297.
- [21] MCLACHLAN, G., AND KRISHNAN, T. *The EM algorithm and extensions*, vol. 382. John Wiley & Sons, 2007.
- [22] PEARL, J. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 2014.
- [23] RABINER, L. R. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE* 77, 2 (1989), 257–286.

- [24] WIKIPEDIA. DBSCAN — Wikipedia, the free encyclopedia. webpage, 2016. <https://en.wikipedia.org/wiki/DBSCAN>, Accessed 9.8.2016.
- [25] YIANNILOS, P. N. Data structures and algorithms for nearest neighbor search in general metric spaces. In *SODA* (1993), vol. 93, pp. 311–21.

# Appendix A

## First appendix

Table A.1: 10 samples from each cluster formed by DBSCAN

Cluster No.	Samples
Cluster 0	945438 ENROLLING 0 WAITING 77 WAITING -1 1230357 ENROLLING 45 WAITING 13 WAITING -1 1960464 ENROLLING 3 WAITING 24 WAITING -1 14813553 ENROLLING 8 WAITING 33 WAITING -1 15253762 ENROLLING 1 WAITING 70 WAITING -1 15254301 ENROLLING 110 WAITING -1 16146920 ENROLLING 128 WAITING 4 IN_TREATMENT_ROOM 52 CLOSED -1 16335930 ENROLLING 1 WAITING 16 WAITING 5 WAITING 1 CANCELLED -1 16759935 ENROLLING 33 WAITING 55 WAITING -1 16760057 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 53 CANCELLED -1
Cluster 1	8490996 ENROLLING 0 WAITING 71 WAITING 0 WAITING -1 13585666 ENROLLING 0 WAITING 78 WAITING -1 15636771 ENROLLING 2 WAITING 23 WAITING -1 16760229 ENROLLING 0 WAITING 76 WAITING -1 17383819 ENROLLING 106 WAITING 0 WAITING -1 17383955 ENROLLING 0 WAITING 18 WAITING 0 WAITING -1 18441848 ENROLLING 222 WAITING 0 WAITING -1 20080659 ENROLLING 0 WAITING 30 WAITING 0 WAITING -1 20330934 ENROLLING 97 WAITING 0 WAITING -1 20454780 WAITING 1 WAITING 239 WAITING 0 WAITING 0 WAITING -1
Cluster 2	15253761 ENROLLING 1 WAITING 27 IN_TREATMENT_ROOM 43 CLOSED -1 15833795 ENROLLING 28 WAITING 48 IN_TREATMENT_ROOM 6 CLOSED -1 15845553 ENROLLING 0 WAITING 16 IN_TREATMENT_ROOM 9 CLOSED -1 15856901 ENROLLING 1 WAITING 4 IN_TREATMENT_ROOM 4 CLOSED -1 15982252 ENROLLING 3 WAITING 2 IN_TREATMENT_ROOM 12 CLOSED -1 16080179 ENROLLING 1 WAITING 3 IN_TREATMENT_ROOM 5 CLOSED -1 16239380 ENROLLING 1 WAITING 1 IN_TREATMENT_ROOM 10 CLOSED -1 16240164 ENROLLING 0 WAITING 41 IN_TREATMENT_ROOM 13 CLOSED -1 16261761 ENROLLING 0 WAITING 1 IN_TREATMENT_ROOM 21 CLOSED -1 16285157 ENROLLING 22 WAITING 16 IN_TREATMENT_ROOM 4 CLOSED -1
Cluster 3	16270356 ENROLLING 1 WAITING 124 WAITING -1 16718646 ENROLLING 1533 CANCELLED -1 16802740 ENROLLING 480 CANCELLED -1 16803188 ENROLLING 35 CANCELLED -1 16808873 ENROLLING 1533 CANCELLED -1 16820174 ENROLLING 173 CANCELLED -1 16847486 ENROLLING 3 CANCELLED -1 16848514 ENROLLING 91 CANCELLED -1 16896527 ENROLLING 84 CANCELLED -1

	16911990 ENROLLING 100 CANCELLED -1
Cluster 4	16759931 ENROLLING 1 WAITING 0 IN_TREATMENT_ROOM 89 CLOSED -1 16759933 ENROLLING 12 WAITING 0 IN_TREATMENT_ROOM 35 CLOSED -1 16759945 ENROLLING 93 WAITING 0 IN_TREATMENT_ROOM 5 CLOSED -1 16760072 ENROLLING 33 WAITING 0 IN_TREATMENT_ROOM 5 CLOSED -1 16760633 ENROLLING 8 WAITING 0 IN_TREATMENT_ROOM 45 CLOSED -1 16760641 ENROLLING 24 WAITING 0 IN_TREATMENT_ROOM 2 CLOSED -1 16802066 ENROLLING 1 WAITING 0 IN_TREATMENT_ROOM 19 CLOSED -1 16802074 ENROLLING 3 WAITING 0 IN_TREATMENT_ROOM 5 CLOSED -1 16802391 ENROLLING 1 WAITING 0 IN_TREATMENT_ROOM 22 CLOSED -1 16802516 ENROLLING 2 WAITING 0 IN_TREATMENT_ROOM 23 CLOSED -1
Cluster 5	16760211 ENROLLING 0 WAITING 36 CANCELLED -1 16802397 ENROLLING 0 WAITING 58 CANCELLED -1 16802538 WAITING 49 CANCELLED -1 16802716 ENROLLING 0 WAITING 2 CANCELLED -1 16802724 ENROLLING 0 WAITING 9 CANCELLED -1 16803196 ENROLLING 0 WAITING 2 CANCELLED -1 16808893 ENROLLING 0 WAITING 0 WAITING 23 CANCELLED -1 16824757 WAITING 3 CANCELLED -1 16825742 WAITING 1 CANCELLED -1 16841048 ENROLLING 0 WAITING 0 WAITING 2 CANCELLED -1
Cluster 6	16760216 ENROLLING 27 WAITING 92 IN_TREATMENT_ROOM 0 CLOSED -1 16802089 ENROLLING 0 WAITING 92 IN_TREATMENT_ROOM 0 CLOSED -1 16803170 ENROLLING 0 WAITING 17 IN_TREATMENT_ROOM 0 CLOSED -1 16883234 WAITING 1 IN_TREATMENT_ROOM 0 CLOSED -1 16889884 ENROLLING 0 WAITING 6 IN_TREATMENT_ROOM 0 CLOSED -1 17137323 ENROLLING 2 WAITING 4 IN_TREATMENT_ROOM 0 CLOSED -1 17144229 ENROLLING 2 WAITING 66 IN_TREATMENT_ROOM 0 CLOSED -1 17266404 ENROLLING 34 WAITING 48 IN_TREATMENT_ROOM 0 CLOSED -1 17266727 ENROLLING 44 WAITING 21 IN_TREATMENT_ROOM 0 CLOSED -1 17383049 ENROLLING 0 WAITING 33 IN_TREATMENT_ROOM 0 CLOSED -1
Cluster 7	16760390 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 13 CLOSED -1 16801915 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 7 CLOSED -1 16802224 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 8 CLOSED -1 16802368 ENROLLING 195 WAITING 37 IN_TREATMENT_ROOM 2 WAITING 2929 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 21 CLOSED -1 16802390 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 4 CLOSED -1 16802750 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 7 CLOSED -1 16803047 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 3 CLOSED -1 16803314 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 8 CLOSED -1 16824756 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 41 CLOSED -1 16824867 ENROLLING 0 WAITING 0 IN_TREATMENT_ROOM 10 CLOSED -1
Cluster 8	16825520 ENROLLING 72 WAITING 7 IN_TREATMENT_ROOM 39 CLOSED -1 16825742 ENROLLING 2 WAITING 0 CANCELLED -1 16847801 ENROLLING 53 WAITING 1 CANCELLED -1 16896665 ENROLLING 3 WAITING 8 WAITING 0 CANCELLED -1 16898355 ENROLLING 50 WAITING 1 CANCELLED -1 17149264 ENROLLING 11 WAITING 0 CANCELLED -1 17266883 ENROLLING 1 WAITING 55 WAITING 0 CANCELLED -1 17383018 ENROLLING 3 WAITING 113 IN_TREATMENT_ROOM 7 CLOSED -1 17383043 ENROLLING 9 WAITING 42 IN_TREATMENT_ROOM 73 CLOSED -1 17383842 ENROLLING 11 WAITING 1 CANCELLED -1
Cluster 9	16896538 ENROLLING 1 WAITING 32 CANCELLED -1 16898635 ENROLLING 1 WAITING 6 CANCELLED -1 17073845 ENROLLING 1 WAITING 5 IN_TREATMENT_ROOM 0 WAITING 6 IN_TREATMENT_ROOM 6 CLOSED -1 17137406 ENROLLING 44 WAITING 3 IN_TREATMENT_ROOM 56 CLOSED -1 17137437 ENROLLING 1 WAITING 56 CANCELLED -1 17137811 ENROLLING 1 WAITING 18 CANCELLED -1 17147824 ENROLLING 0 WAITING 1323 CANCELLED -1 17250342 ENROLLING 1 WAITING 68 CANCELLED -1

	17383392 ENROLLING 42 WAITING 3 IN_TREATMENT_ROOM 50 CLOSED -1 17383812 ENROLLING 1 WAITING 7 CANCELLED -1
Cluster 10	16896542 WAITING 0 WAITING 4 WAITING -1 16898754 ENROLLING 10 WAITING 55 IN_TREATMENT_ROOM 2 CLOSED -1 16911970 WAITING 239 CANCELLED -1 18042974 ENROLLING 8 WAITING 1 IN_TREATMENT_ROOM 71 CLOSED -1 18206521 ENROLLING 0 WAITING 0 WAITING 136 CANCELLED -1 18206522 ENROLLING 0 WAITING 0 WAITING 136 CANCELLED -1 19336998 ENROLLING 0 WAITING 1 IN_TREATMENT_ROOM 21 WAITING 0 WAITING 536 CANCELLED -1 19488319 WAITING 13 WAITING 0 WAITING 0 CANCELLED -1 19517799 WAITING 215 CANCELLED -1 19551936 ENROLLING 3 WAITING 0 WAITING 154 CANCELLED -1
Cluster 11	17141316 ENROLLING 3 WAITING 24 CANCELLED -1 17267028 ENROLLING 59 WAITING 13 IN_TREATMENT_ROOM 69 CLOSED -1 17333621 ENROLLING 1 WAITING 1 CANCELLED -1 17384117 ENROLLING 40 WAITING 11 IN_TREATMENT_ROOM 52 CLOSED -1 17384134 ENROLLING 1 WAITING 105 IN_TREATMENT_ROOM 92 CLOSED -1 17445733 ENROLLING 1104 ENROLLING -1 17472038 ENROLLING 9 WAITING 37 IN_TREATMENT_ROOM 70 CLOSED -1 17472501 ENROLLING 52 WAITING 118 IN_TREATMENT_ROOM 22 CLOSED -1 17484544 ENROLLING 0 WAITING 1457 ENROLLING 0 WAITING -1 17512538 ENROLLING 0 WAITING 16 WAITING 1002 CANCELLED -1

Table A.2: Samples with log-likelihood computed using first-order Markov Chain using negative binomial distribution as emission function.

log-likelihood	sample
-1.22501248649	21332189 ENROLLING 0 WAITING 3 IN_TREATMENT_ROOM 5 CLOSED -1
-1.22501248649	21332152 ENROLLING 0 WAITING 3 IN_TREATMENT_ROOM 5 CLOSED -1
-1.22501248649	21243121 ENROLLING 0 WAITING 3 IN_TREATMENT_ROOM 5 CLOSED -1
-1.22501248649	21242603 ENROLLING 0 WAITING 3 IN_TREATMENT_ROOM 5 CLOSED -1
-1.22501248649	20753590 ENROLLING 0 WAITING 3 IN_TREATMENT_ROOM 5 CLOSED -1
-1.31252682278	20531975 ENROLLING 0 WAITING 9 IN_TREATMENT_ROOM 5 CLOSED -1
-1.31252682278	20335903 ENROLLING 0 WAITING 9 IN_TREATMENT_ROOM 5 CLOSED -1
-1.31252682278	20332176 ENROLLING 0 WAITING 9 IN_TREATMENT_ROOM 5 CLOSED -1
-1.31252682278	20330946 ENROLLING 0 WAITING 9 IN_TREATMENT_ROOM 5 CLOSED -1
-1.31252682278	19524600 ENROLLING 0 WAITING 9 IN_TREATMENT_ROOM 5 CLOSED -1
-1.31252682278	18895611 ENROLLING 0 WAITING 9 IN_TREATMENT_ROOM 5 CLOSED -1
-1.6200398225	20090272 ENROLLING 0 WAITING 20 IN_TREATMENT_ROOM 5 CLOSED -1
-1.6200398225	17137238 ENROLLING 0 WAITING 20 IN_TREATMENT_ROOM 5 CLOSED -1
-1.6200398225	17137111 ENROLLING 0 WAITING 20 IN_TREATMENT_ROOM 5 CLOSED -1
-1.6200398225	16803024 ENROLLING 0 WAITING 20 IN_TREATMENT_ROOM 5 CLOSED -1
-1.62051607252	21419996 ENROLLING 0 WAITING 14 IN_TREATMENT_ROOM 13 CLOSED -1
-1.81955956629	18946307 ENROLLING 1 WAITING 0 IN_TREATMENT_ROOM 10 CLOSED -1
-1.81955956629	17010037 ENROLLING 1 WAITING 0 IN_TREATMENT_ROOM 10 CLOSED -1
-1.81974417979	20538959 ENROLLING 1 WAITING 5 IN_TREATMENT_ROOM 13 CLOSED -1
-1.81974417979	16824641 ENROLLING 1 WAITING 5 IN_TREATMENT_ROOM 13 CLOSED -1
-1.82019772133	19969197 ENROLLING 2 WAITING 7 IN_TREATMENT_ROOM 6 CLOSED -1
-2.09969260165	21332560 ENROLLING 6 WAITING 4 IN_TREATMENT_ROOM 8 CLOSED -1
-2.09974363796	21419663 ENROLLING 1 WAITING 1 IN_TREATMENT_ROOM 19 CLOSED -1
-2.10057254235	19664612 ENROLLING 1 WAITING 19 IN_TREATMENT_ROOM 11 CLOSED -1
-2.10063359561	17267455 ENROLLING 3 WAITING 9 IN_TREATMENT_ROOM 12 CLOSED -1
-2.10077419383	20538310 ENROLLING 1 WAITING 20 IN_TREATMENT_ROOM 10 CLOSED -1
-2.52467325597	18335666 ENROLLING 1 WAITING 33 IN_TREATMENT_ROOM 10 CLOSED -1
-2.52507176943	18243950 ENROLLING 1 WAITING 35 IN_TREATMENT_ROOM 4 CLOSED -1
-2.52507176943	16911955 ENROLLING 1 WAITING 35 IN_TREATMENT_ROOM 4 CLOSED -1
-2.52523586276	18160697 ENROLLING 2 WAITING 31 IN_TREATMENT_ROOM 5 CLOSED -1
-2.52574525145	19841457 ENROLLING 7 WAITING 20 IN_TREATMENT_ROOM 7 CLOSED -1
-3.06107298199	19954450 ENROLLING 0 WAITING 61 IN_TREATMENT_ROOM 3 CLOSED -1
-3.06260151935	21127113 ENROLLING 11 WAITING 30 IN_TREATMENT_ROOM 9 CLOSED -1
-3.06387434043	19230307 ENROLLING 28 WAITING 10 IN_TREATMENT_ROOM 11 CLOSED -1
-3.06458227144	18247745 ENROLLING 2 WAITING 10 IN_TREATMENT_ROOM 33 CLOSED -1
-3.06629701355	17141536 ENROLLING 26 WAITING 13 IN_TREATMENT_ROOM 2 CLOSED -1
-3.77652544852	20624690 ENROLLING 7 WAITING 54 IN_TREATMENT_ROOM 11 CLOSED -1
-3.77716200844	19230163 ENROLLING 0 WAITING 32 IN_TREATMENT_ROOM 43 CLOSED -1
-3.77716858350	17266875 ENROLLING 3 WAITING 25 IN_TREATMENT_ROOM 36 CLOSED -1
-3.77756155251	16848287 ENROLLING 58 WAITING 10 IN_TREATMENT_ROOM 8 CLOSED -1
-3.77963738529	17384149 ENROLLING 1 WAITING 42 CANCELLED -1
-5.10571464997	20032830 ENROLLING 28 CANCELLED -1
-5.10571464997	18004184 ENROLLING 28 CANCELLED -1
-5.10571464997	18004183 ENROLLING 28 CANCELLED -1
-5.10766253514	18945732 ENROLLING 30 WAITING 11 IN_TREATMENT_ROOM 48 CLOSED -1
-5.10842170221	18043295 ENROLLING 0 WAITING 19 IN_TREATMENT_ROOM 71 CLOSED -1
-5.11050052549	20538023 ENROLLING 60 WAITING 8 CANCELLED -1
-227.515556821	18883776 ENROLLING 0 WAITING 4435 ENROLLING -1
-294.647987389	19593547 ENROLLING 0 WAITING 5754 ENROLLING -1
-653.392679598	19229957 ENROLLING 37265 ENROLLING 40 WAITING 1 IN_TREATMENT_ROOM 4 CLOSED -1
-876.235768456	21127118 ENROLLING 50174 ENROLLING 0 WAITING 1 IN_TREATMENT_ROOM 3 CLOSED -1
-1227.10081314	17144893 ENROLLING 0 WAITING 23 WAITING 48064 ENROLLING 0 WAITING 18 CANCELLED -1