

Developing a Low Dimensional Patient Class Profile in Accordance to Their Respiration-Induced Tumor Motion

Rittika Shamsuddin
University of Texas, Dallas,
Richardson, TX 75080, USA
Rittika.Shamsuddin
@utdallas.edu

Amit Sawant
Radiation Oncology, University
of Maryland, School of
Medicine, MD, USA
asawant
@som.umaryland.edu

Balakrishnan
Prabhakaran
University of Texas, Dallas,
Richardson, TX 75080, USA
bprabhakaran
@utdallas.edu

ABSTRACT

Tumor location displacement caused by respiration-induced motion reduces the efficacy of radiation therapy. Three medically relevant patterns are often observed in the respiration-induced motion signal: *baseline shift*, *ES-Range shift*, and *D-Range shift*.

In this paper, for patients with lower body cancer, we develop class profiles (a low dimensional pattern frequency structure) that characterize them in terms of these three medically relevant patterns. We propose an adaptive segmentation technique that turns each respiration-induced motion signal into a multi-set of segments based on *persistent* variations within the signal. These multi-sets of segments is then probed for *base behaviors*. These base behaviors are then used to develop the group/class profiles using a modified version of the clustering technique described in [1]. Finally, via quantitative analysis, we provide a medical characterization for the class profiles, which can be used to explore breathing intervention technique.

We show that, with i) carefully designed feature sets, ii) the proposed adaptive segmentation technique, iii) the reasonable modifications to an existing clustering algorithm for multi-sets, and iv) the proposed medical characterization methodology, it is possible to reduce the time series respiration-induced motion signals into a compact class profile. One of our co-authors is a medical physician and we used his expert opinion to verify the results.

1. INTRODUCTION

A patient profile, in medical terms, is a list of data on an individual patient collected during their treatment and might include different kinds of measurements, such as age, heart rate, blood pressure, etc. Often, full or partial data from the profile can be used for disease analysis, where different kinds of measurements can be treated as different attributes or variables. However, sometimes patient data

consists of one or more continuous time series, with unannotated patterns. Moreover, the data may include signals only from patients with no supporting control dataset. Another level of difficulty gets added if there is high variations between patients and also, between signals from the same patient collected during different times. An example of such a dataset in the medical field, would be the respiration induced tumor motion arising from patient treatment during radiation therapy.

Radiation therapy (RT) involves delivering a tumoricidal dose to cancerous tissue, while ensuring minimal damage to healthy tissues and organs at risk. Respiration induced motion causes significant geometric and dosimetric uncertainties in radiotherapy for thoracic and abdominal tumors. Such uncertainties have greater impact in case of hypofractionated regimens such as lung stereotactic body radiotherapy (SBRT) [13], where very high, biologically potent doses are delivered in relatively few fractions; 3-5 compared to 30 for conventionally-fractionated RT.

1.1 Dataset and Patterns of Interest

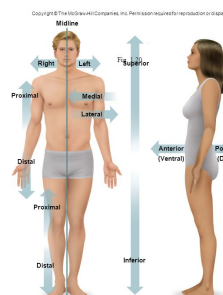


Figure 1: Superior-Inferior (up-down), Anterior-Posterior (in-out) and Left-Right axes of human physiology.

The tumor location displacement (Figure 2) is an inevitable consequence of respiration, which involves the coordinated movement of several abdominal muscle groups, and results in dynamic variations in the motion characteristics. These data were collected from patients with thoracic and abdominal cancer, in Georgetown University Hospital, who were treated with Cyberknife Synchrony [18], [11]. The Synchrony subsystem (Figure 4) tracks tumor motion by estimating the tumor position (using an internal algorithm) and produces 4D data, of 3D position information (in mm)

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing info@vlldb.org.

Proceedings of the VLDB Endowment, Vol. 10, No. 12
Copyright 2017 VLDB Endowment 2150-8097/17/08.

with time (in seconds). The tracking process requires the implantation of gold internal fiducial markers, making data collection expensive and difficult [7]. In this dataset, one entire clinical recording is referred to as a **fraction**, and thus, one patient can have more than one fraction. There are 46 patients and 160 fractions with the length of each fraction ranging between 7810 record samples (over 5 minutes) and 165592 record samples (over 110 minutes). So even though the dataset is small in terms of the number of patients, it is large in terms of samples collected (approximately, 84 hours of time series data with 10^6 recorded samples).

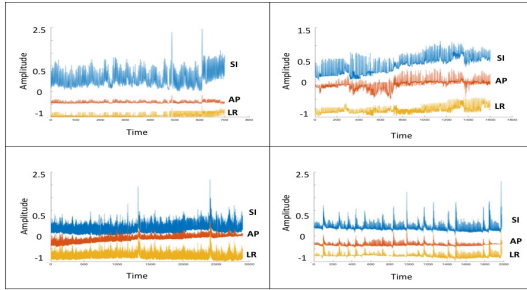


Figure 2: 4D respiration-induced tumor-motion data of four different fractions from the dataset, showing inter and intra fraction variability (translated along the vertical axis for better readability), with SI showing maximum variation. SI refers to the superior-interior axis, AP to Anterior-Posterior axis and LR to Left-Right axis.

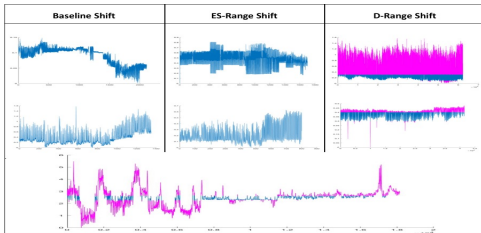


Figure 3: Few varieties of Baseline Shift, ES-Range Shift and D-Range Shift as seen in the dataset along the SI axis. Along the gating column, the part of the signal colored in magenta illustrates the inhalation (and sometimes the exhalation) instances where gating is necessary.

The inter-patient variations can be explained by the dependency of respiratory pattern on factors like body size, gender, age, life style, etc. The intra-patient variation can be explained by the fact that breathing function is controlled consciously and unconsciously and is prone to be affected by the person’s emotions, instantaneous thought process, etc. The errors due to these uncertainties can cause geometric misalignments between the radiation beam (which is directed based on computed tomographic imaging acquired a few days before treatment) and the instantaneous position of the tumor target. Such misalignment can reduce the efficacy of radiation therapy due to the fact that tumors may receive less than prescribed dose (thereby not achieving the desired cell-kill) or normal tissue and critical organs may

receive more than intended dose (thereby causing excessive radiation-related toxicity).

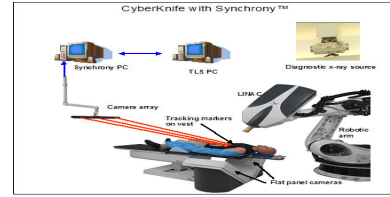


Figure 4: CyberKnife with Synchrony Tracking System

These variations can arise from three types of effects that cause cycle-to-cycle changes in respiratory motion: (a) amplitude changes, (b) frequency changes and (c) baseline shifts. For a given time interval within a respiratory trace, these three types of changes may occur separately or in combination (any two or all three). In this paper, we are concerned with **baseline shifts** and two forms of amplitude changes:

1. **Baseline Shifts** are essentially changes in the trend lines and refer to permanent or persistent changes in the mean position of the tumor.
2. **ES-Range Shift**, a form of amplitude change, describes a permanent or persistent change (**expansion/shrinkage**) in the range of the amplitude fluctuation.
3. **D-Range Shift**, another form of amplitude change, is a **deviation** in amplitude from an allowable range of motion, as pre-defined by the physician. This amplitude change, often requires the radiation beam to be turned off during inhalation or exhalation.

Knowing if the patient to be treated is prone to display a particular variation on treatment day, will allow the physicians to take appropriate preventive steps to regularize breathing patterns and negate the three variations and reduce radiation-related toxicity. The respiratory motion management technique involves: 1) **abdominal compression** for baseline shifts, where the patient is required to wear a belt around stomach to control the amplitude of motion, 2) **marginal expansion** for ES-range shifts, which involves defining the tumor and target volume for radiation therapy, and 3) **gating** for amplitude deviations, where the tumor motion is tracked so that the radiation can be turned off when/if required. The three major patterns that characterize the dataset are shown in Figure 3.

1.2 Our Approach and Contributions

We observe that a patient can exhibit either a gradual baseline shift, a sudden permanent baseline shift or a sudden transient change in mean position. These baseline shifts or changes in the trend line of the signal, constitute a change in the mean of a signal segment. However, a change in mean is also possible by uneven change in the amplitude range resulting in ES-range shift. Thus, we hypothesize that the clinically/medically relevant patterns are a combination of simple, yet unknown, patterns. As such, we propose the creation of group profiles for the patients for analysis. The group profiles are developed based on distribution profiles of patients as inspired by the creation of class profiles in [1]

via clustering. The creation of these profiles is preceded by our own signal segmentation and followed by medical characterization of the profiles based on our hypothesis of simple pattern combination.

Our main **contributions** come in the form of studying an unimodal, unannotated dataset of relatively few patients (due to the difficulty and expense of collecting data) but a large number of samples, with lots of intra and inter patient variability. In order to obtain useful and medically relevant information from this dataset we propose an adaptable segmentation method, which bypasses the need for concrete thresholds and thus, gives reasonable results with all fractions or signals without any parameter adjustments. And even though we use the data structures (that facilitates the clustering of the fractions) and cluster centroids in the same manner as described in [1], we propose a different way to initialize the structures, a different seeding technique, and a different distance metric for clustering to suit our application. We also incorporate a way to update or learn parameters to suit our choice of feature vectors. Finally, we provide an objective and quantitative medical interpretation of the clustering results.

2. RELATED WORK

To the best of our knowledge, this is the first attempt to create specialized profiles for patients using unimodal, time series data.

As for segmenting signals, [8], [9] and [14] suggest using different variations of Modified Varri, a standard adaptive segmentation technique. [8] proposes to smooth the signal using moving average or Savitzky-Golay filter, which avoids causing shifts in the original signal, prior to using Modified Varri. [14], on the other hand, suggests to replace the original feature used in the original Modified Varri with fractal dimension approximated by Katz’s algorithm. While both methods show improvement in performance, both suffer from the three parameter problem associated with Modified Varri. If the size of the two windows, their percentage overlap and the value of a sensitivity threshold are not set correctly, the Modified Varri fails to segment properly. [9] suggests an approximate solution to the parameter selection problem by using optimization via genetic algorithm. However, the fact remains that the solution is approximate and adds a computational overhead to the segmentation algorithm. Also, all three variations of Modified Varri were tested on EEG signals, which exhibit non-stationery variations. On a very different note, [6] simplifies the BIC model, by using constant window size, to obtain a generalized likelihood ratio model to segment audio signals. However, this assumes that the data points within a window follow a Gaussian distribution, which invariably imposes a large window size. [17] describes the creation of DSTree, which facilitates the adaptive and dynamic indexing of time series data. They prove an upper bound on the distance between two time series and use the upper and lower bounds to restrict their search space without calculating the distance between actual time series. It provides two methods for query based search. While, it does prove to be efficient segmentation technique on large dataset, the DSTree does not match our application requirements and adds extra overhead of creating the tree.

In terms of clustering, the technique for fuzzy clustering [12], is very similar to frequency based clustering used

in this paper. However, the frequency based technique described in [1], which is modified in this paper, does not assign one datapoint to multiple centroids with a membership score. Rather, the frequency based technique assigns different data points *belonging to a single set* to different pre-chosen centroids in a non-fuzzy way. The *set of these data points* is then defined as a frequency vector, which measures how many of its data point belongs to each centroid. Since [1] provides a way for clustering multiple sets, where each set is a collection of data points, the clustering method fits our paradigm where each fraction (or a set) is a collection of segments (or data points). However, due to challenges faced with the dataset, we have to make changes to the algorithm described in [1] and the details of these changes can be found in Section III. In terms of similar application, [15] creates call profiles of customers to optimize tariff using clustering. However, the crux of the paper is the use of big data tools like HBase and Storm and [15] provides strategies for re-partitioning and aggregation of data to achieve their goals. Another example of patient profiling can be found in [19], where they create a healthcare system that gathers data from different patient textual databases, creates profile trees for patients, and provides platform for doctor intervention during the analysis stage.

3. CHALLENGES AND OVERVIEW

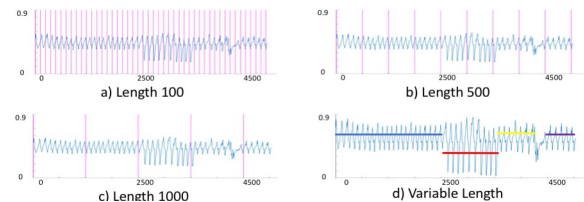


Figure 5: Fixed length segmentation of a section of fraction 117, setting the fixed length window at a) 100, b) 500, and c) 1000. d) Shows the ideal segmentation on this fraction fragment into 5 different segments of varying length.

As mentioned earlier, the dataset has no annotations and no control dataset. Complicating matters, our dataset is small (in terms of the number of patients, but large in terms of total hours) and displays a multitude of inter- and inpatient variations. Thus, it is impossible to find a standard mode of variation, against which we can compare or define anomalies or novel patterns. Moreover, it is an unimodal dataset, consisting of only the signals representing the respiration-induced tumor motion without any information on the patient’s age, life style details, changing heart rate, etc. Also, to the best of our knowledge, the three medically relevant patterns that are of interest to the oncologists, have not been defined in computational or mathematical terms. In our work, we rely on expert opinion for the validation of our results. Thus, in order to handle this dataset, we rely on segmentation and clustering techniques. Our workflow consists of three main stages of processing (described in Sections IV, V and VI, respectively):

1. **Stage 1- Segmentation:** In time series processing, the segmentation methodology determines how a single signal is divided into multiple segments for ease

of processing. In a fixed window segmentation, the signals are divided into segments of fixed length, L . However, for any choice of L , it is highly likely that the predetermined point of segmentation breaks a pattern of interest into two halves, affecting the nature of feature vectors derived from the segments used for further analysis (Figure 5). Thus the challenge here is to capture consistent behavior in this volatile dataset. Hence, we propose an adaptive segmentation process, the *variable length segmentation* (VL segmentation), which should entail the following characteristics: **1)** the segments produced should encapsulate only one form of or no interesting variation. This ensures that a feature vector associated with a segment is characteristic of that particular variation. This is necessary because the profiling process largely relies on clustering the segments, **2)** the segmentation should not be affected by the starting point of the algorithm. That is to say, regardless of whether the algorithm is initialized at time point t_i or t_{i+c} , where c is positive constant, both runs of the algorithm should result in the same segments for time points following t_{i+c} , **3)** the segmentation process should overlook *trivial variations* (explained fully in the next section), **4)** the algorithm should produce acceptable results for all patients and fractions, without manual intervention of parameters or thresholds.

2. **Stage 2- Creation of distribution profiles:** In this stage, we cluster the segments produced in *Stage 1*, for each fraction and develop their distribution profiles. At the end of this stage, we expect to describe every fraction in terms of the frequencies of few selected initial patterns (or the distribution profiles) and group the fractions into groups of similar frequency behaviors. To accomplish this, we rely on the multi-set or setwise clustering paradigm described in [1]. However, this paradigm assumes a large volume (over 10000 data points) of incoming data that changes over time; our data, changes over time but has a small size and does not involve streaming. As such, the selection of initial patterns, on which the frequency behavior is formed, requires to be handled differently, without compromising the stability of the group profiles. Also, grouping the fractions, based on their distribution profiles, requires a different method for initializing the group profile structure to ensure purity of the resulting groups. In addition, [1] suggests us to use euclidean distance as the similarity matrix when forming the group profiles. While euclidean distance is quite common and robust, it does not capture the differences (and similarities) in the frequency behaviors as required by the application (example is provided in description of Phase 3 in Section V). Thus, we had to incorporate a simple but non-trivial penalty term to euclidean distance metric to capture the required differences in the frequency behaviors. Due to the need to use discrete feature vector, it is worth mentioning that we need to incorporate a way of updating the discretization parameters during the formation of the distribution and group profiles.
3. **Stage 3- Medical characterization of group profiles:** During this stage, we provide quantitative analysis of the distribution and group profiles, by inter-

preting the selected initial patterns (from *Stage 2*) in terms that are found to be medically relevant. The challenge here is finding and describing the bridge that connects the group profiles to the medically relevant patterns. This is accomplished by defining *templates* for the medical behaviors in terms of the selected initial patterns and assigning three descriptive scores to each group profile.

4. VARIABLE LENGTH SEGMENTATION

The variable length (VL) segmentation is designed to catch and prioritize persistent amplitude changes or variations in the amplitude patterns, lasting longer than few milliseconds. A persistent change results in a change in pattern, whereas a transient change is a momentary variation in the data and thus, trivial by definition.

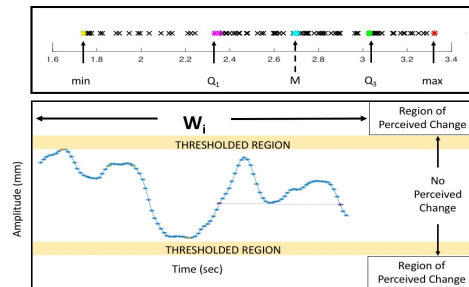


Figure 6: The five-number summary of a particular W_i , where min denotes the minimum, Q_1 denotes the first quartile, M denotes the median, Q_3 denotes the third quartile, max denotes the maximum. The figure also illustrates the defined regions for perceived change.

VL segmentation achieves this aim by pausing at *every* possible point of *perceived* change as it loops over the signal using a sliding window that moves one time step at a time. A perceived change (Figure 6) is defined as a change, which is caused by an amplitude value that goes beyond a *thresholded region* defined for the current sliding window, W_i . The *thresholded region* for the upper envelope is given by two perfectly horizontal lines, $y = max(W_i)$ and $y = max(W_i) + \delta_{W_i}$, and for the lower envelope by $y = min(W_i)$ and $y = min(W_i) - \delta_{W_i}$. Note, a perceived change can be either persistent or transient. And hence, when a perceived change is detected, VL executes a series of tests and verifies that the change is persistent and hence, a point of segmentation. If there is no perceived change, then VL segmentation moves on to the next iteration after updating its variables.

Next, to search for and determine persistent changes, we use the **five-number summary** [5] for W_i , an example of which is shown in Figure 6. (min, Q_1, M, Q_3, max) gives information about the location (from the median), spread (from the quartiles) and range (from minimum and maximum) of the samples in W_i . Computationally, *any* change in the pattern of amplitude (whether it is a spike or a permanent increase/decrease in the size of amplitude or a change in mean position of the signal) will invariably alter the maximum and minimum values of W_i . A *persistent change* on the other hand, is likely to have a *deeper* effect in the neighborhood of the change and likely to affect not only the maximum

(max) and minimum (min) values of W_i but also the first (Q_1) and third (Q_3) quartile values of W_i . And as such, VL segmentation needs to look for changes in maximum and minimum of the amplitudes, followed by the changes in quartile values.

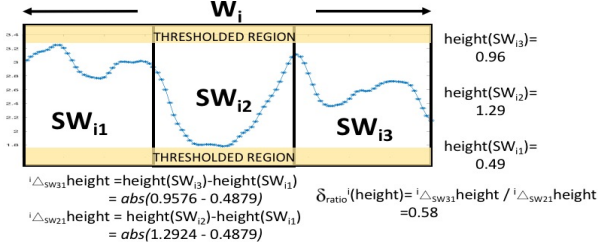


Figure 7: An example of Checkpoint 1

Note, change is a relative concept; that is, in general terms, to detect a change we need to compare at least *two* objects separated by time or space or both using a threshold value. If the threshold value cannot be set to any fixed value, we need to calculate a varying threshold. This can be done by introducing a *third* object, whose comparison with the other two objects, can be used to gauge the appropriate threshold level, at any given time. Thus, to detect changes with varying thresholds, every W_i is divided into and associated with three sub-windows, SW_{i1} , SW_{i2} and SW_{i3} . As shown in Figure 7, SW_{i3} encapsulates the most recent points in W_i , whereas SW_{i1} encapsulates the earliest points in W_i .

4.1 Checkpoint 1: Intra-Window Amplitude Change Ratio (Intra-ACR)

We combine our tools and strategy by defining $height(SW_{ij})$ to be the $abs(maximum(SW_{ij}) - minimum(SW_{ij}))$, where abs is the absolute value operation, i refers to the corresponding W_i and $j \in [1, 2, 3]$ subscripts the sub-windows. Then, a change can be detected as a ratio of height changes between the sub-windows for any particular W_i . The height change ratio for W_i is given by:

$$\begin{aligned} \delta_{intra-ACR}^i &= \frac{{}^i\Delta_{SW31}}{{}^i\Delta_{SW21}} \text{ where,} \\ {}^i\Delta_{SW31} &= abs[height(SW_{i3}) - height(SW_{i1})] \\ {}^i\Delta_{SW21} &= abs[height(SW_{i2}) - height(SW_{i1})] \end{aligned} \quad (1)$$

If Intra-ACR is zero or less than one or both ${}^i\Delta_{SW31}$ and ${}^i\Delta_{SW21}$ equal to zero, then we can conclude that there is no overall change in the maximum and minimum values of W_i . It can also mean that here is no change in SW_{i3} , which contains the most recent points. Challenges arise if ${}^i\Delta_{SW21}$ is zero and ${}^i\Delta_{SW31}$ is non-zero (case 1) or if the intra-ACR itself is equal to or greater than 1 (case 2). In general, at the occurrence of case 1, there is probably a sign of impending change; however this change might be a spike or a trivial transient change. Case 2 can result when both ${}^i\Delta_{SW31}$ and ${}^i\Delta_{SW21}$ are either very high or very low, which can be due to transient changes. Thus, to rule out transient changes and as a measure of persistence, we put forward two more tests: an *Inter-Window Amplitude Change Ratio checkpoint* using heights, followed by a *Quartile checkpoint* using quartile comparison.

4.2 Checkpoint 2: Inter-Window Amplitude Change Ratio (Inter-ACR)

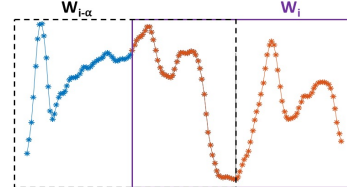


Figure 8: The two windows for Inter-ACR

To verify that the change detected by the Intra-ACR, is relevant along the time frame, we compare the change in sub-window SW_{i3} to the change in a window that lags behind W_i by α time steps and is denoted by $W_{i-\alpha}$ (Figure 8). More precisely, we compare:

$$\begin{aligned} {}^i\Delta_{SW31}^{normalized} &= \frac{{}^i\Delta_{SW31}}{height(W_i)} \\ \Delta_{W_{i\alpha}}^{normalized} &= \frac{abs(height(W_{i-\alpha}) - height(W_i))}{height(W_{i-\alpha})} \quad (2) \end{aligned}$$

checkpoint :

$${}^i\Delta_{SW31}^{normalized} > \Delta_{W_{i\alpha}}^{normalized}$$

Note, if the comparison, ${}^i\Delta_{SW31}^{normalized} > \Delta_{W_{i\alpha}}^{normalized}$, is true then change detected in SW_{i3} is larger than the change detected in ($W_{i-\alpha}$), showing evidence for impending persistent change.

4.3 Checkpoint 3: Quartile Change Ratio

The *quartile checkpoint* is only brought into action if the Inter-ACR has been successfully satisfied and is the main tool for determining persistent changes. Unlike a transient change, such as a spike, if a persistent change is taking place in the vicinity of W_i , then in addition to the minimum and maximum values of W_i , the first and third quartile values of W_i is likely to be affected as well. For example, take a *signal* = [1, 1, 1, 1, 1, 1, 10, 1, 1, 1, 1,], which clearly shows a spike as the transient change. A sliding window moving over *signal* will show an obvious change in maximum value; however, its first (Q_1) and third (Q_3) quartile values would remain unchanged. As such, the *quartile checkpoint* uses the interquartile range, iqr , of W_i , where $iqr(W_i) = Q_3(W_i) - Q_1(W_i)$, where $Q_3(W_i) \geq Q_1(W_i)$ by definition. The *quartile checkpoint* sets up tests in a manner that is perfectly analogous to the Intra-ACR and Inter-ACR checkpoints, except instead of using *height*, they use the *interquartile range*. Also, if the condition described in checkpoint 2 (using *iqr*) is satisfied, then we finally mark the point as a segmentation point. The VL segmentation algorithm is given in Table 1.

Table I. Pseudo code for VL Segmentation

```

TR = Thresholded Region
Declare window size,  $W_{size}$ 
Initialize, the  $W$ , height and IQR
for each point,  $p$ 
  Add  $p$  to  $W_{i-1}$  and obtain  $W_i$ 
  if  $p$  lies within TR
    Remove  $p_1$  in  $W_i$ 
    Update the height and IQR of  $W_i$ 
  else
    if Intra-ACR satisfied then
      if Inter-ACR satisfied then
        if Quartile Checkpoint satisfied then
          Mark Segmentation Point.
    end for
end for

```

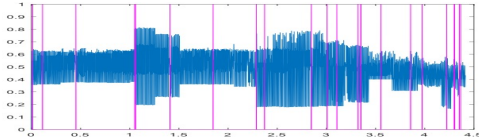


Figure 9: Segmentation of ‘whole’ fraction 117 as achieved via VL Segmentation

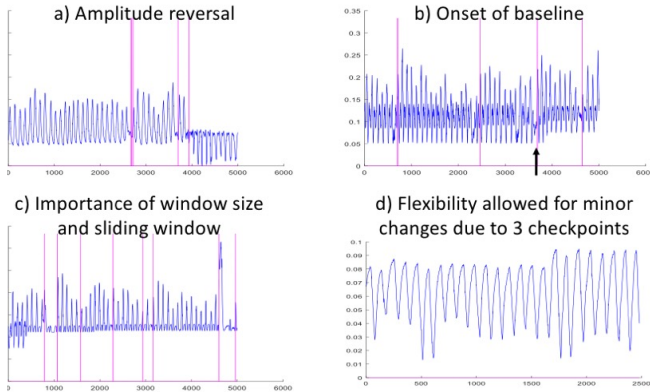


Figure 10: The benefits of using VL Segmentation

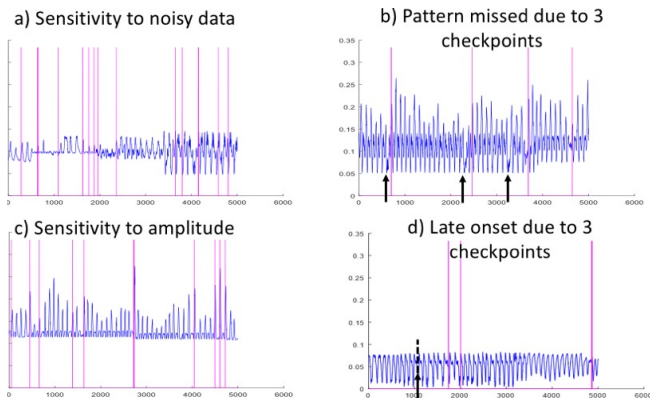


Figure 11: Some trade-off of using VL Segmentation

VL Segmentation is robust to starting point:- We mentioned earlier that with fixed length segmentation can produce segments that vary depending on the starting point of segmentation. So, we empirically verified that VL

segmentation is more robust than fixed length segmentation. For verification, we went over each fraction in the dataset, and for each fraction, we ran VL algorithm multiple times, by removing a record from the left each time until no records were left to be removed. Each time we recorded the segmentation mark point. To clarify our finding, let us denote c_j as a marked segmentation point and $Area_j$ as the region of 200 sample points preceding c_j . Then as long as the algorithm is initialized, such that the starting point is not in $Area_j$, c_j will be detected for the current set of fractions.

4.4 Visualizing VL Segmentation

To summarize the description of variable segmentation, we present the results in Figures 9, 10 and 11. Figure 9 shows all the segmentation points detected for the entire fraction 117. Figure 10 demonstrates that VL captures the important patterns, such as, reversal or reduction of amplitude (Figure 10a and 10c) and changes in mean position (Figure 10b); while Figure 10d shows the flexibility allowed by VL. Figure 11a shows that VL is still sensitive to noisy or aberrant data to some extent, while Figure 11c shows the delicate balance between catching important and trivial expansion/reduction in amplitude. Figure 11b shows that some trivial but repeating patterns (black arrows) are ignored by the technique, while Figure 11d focuses on the fact that sometimes, VL marks a point for segmentation after the change in already underway (ideal mark indicated by the black arrow). Also, it is worth mentioning that by the end of VL segmentation, we end up with 9665 segments.

5. ADOPTING SETWISE CLUSTERING FOR RESPIRATION-INDUCED TUMOR MOTION DATA

As mentioned before, the medically relevant patterns are a combination of simpler patterns. And given the requirements and challenges, it is useful to computationally learn or identify those simple patterns (or **base behaviors**) and obtain a distribution profile for each patient in terms of the base behaviors. The framework for such an analysis is provided by *multi-set stream clustering*, as described in [1]. The segmentation process allows each patient to be described as a sequence of temporally ordered segments, thus creating *entities* [2] or *multi-sets* [1] of records/segments. However, with the respiration-induced dataset, it is not possible to implement the framework exactly as it is implemented in [1]. The three main reasons are:

- **R1:** The theory proposed in [1], assumes a large dataset, with incoming data points via streaming. Our dataset, though large in the number of samples, is small in terms of number of patients, and so, we need to optimize the use of data at each stage of the algorithm.
- **R2:** [1] is proposed for applications (such as, the sale’s behavior of different stores of a large super market chain) that do not *necessarily* require a specific kind of base behaviors e.g. they do not *need* to (but might) have a semantic meaning. For our application, the base behaviors need to map to or describe complex medical patterns and hence, requires to be chosen explicitly in terms of baseline shifts and amplitude and

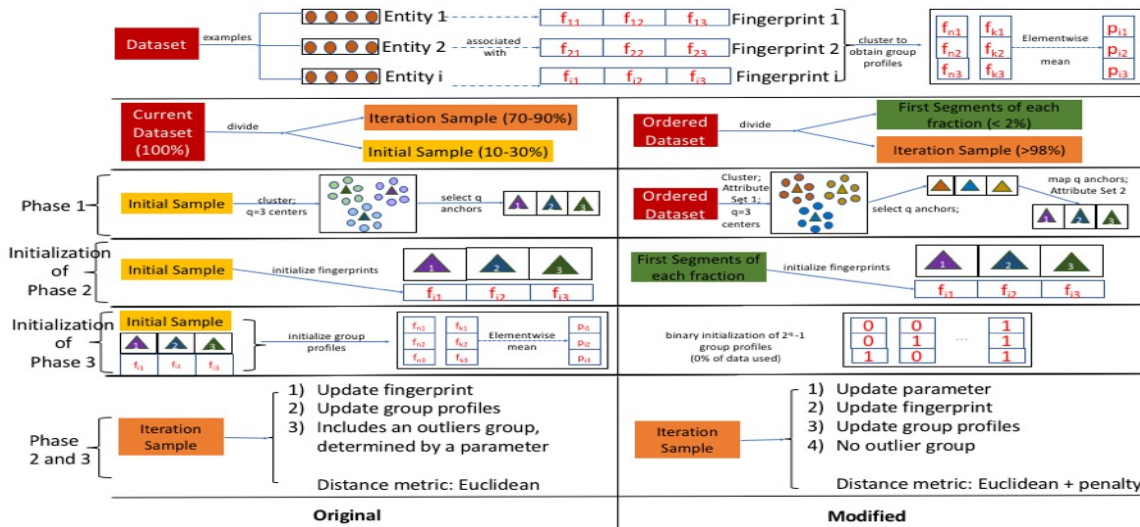


Figure 12: The original and modified setwise clustering framework. f_{ij} is the j^{th} anchor frequency of the i^{th} entity. p_{kj} is the j^{th} anchor average frequency for the k^{th} class profile.

frequency variations. On the other hand, for distribution profile creation, we require segments to be described succinctly in lower dimension.

- **R3:** The definition of a match, in terms of two distribution profiles for our application, is different from a match described in [1].

We schematically present the original and the modified frameworks used for analysis in this paper in Figure 12. For clarity, the process described in Figure 12, is presented in three phases: **Phase 1**-Selection of base behaviors (or *anchors*), **Phase 2**-Creation of distribution profiles. **Phase 3**-Creation of group profiles.

5.1 Phase 1: Selection of base behaviors

The exact base behaviors, though unknown to us, are present within the collection of segments produced by VL segmentation. Ideally, the base behaviors should be represented by a minimal set of segments that differ from one another. In unsupervised learning, clustering is the most straightforward way for finding class boundaries and the final centroids, thus obtained, are expected to be as different from each other as possible. Thus, in order for us to find the different base behaviors, we can **cluster the segments** and use the description of the centroids as a guide to finding the base behaviors (Figure 12). This process is completely analogous to the process of finding *anchors* in multi-set stream clustering (Figure 12). However, the modification of the original algorithm comes in the form of three crucial decisions: 1) the **population or sample of segments to use for clustering**, 2) **attributes for defining the segments**, and 3) **the number of centroids/anchors**.

5.1.1 Deciding the sample

[1] divides the entire dataset into two samples: the *initial* sample and the *iteration* sample. The initial sample (10%-30% of the actual population) is then clustered to obtain the anchors, under the assumption that the sample is perfectly representative of the actual population of data points.

This assumption, however, only holds if the dataset is large enough. Using an unrepresentative initial sample will lead to false distribution profiles, as they are defined in terms of the anchors or base behaviors. On the other hand, if we have to use over 50% of the population to obtain a representative initial sample, then the size of the iteration sample goes down and makes the final result unstable.

As such, we decide to probe the whole population as initial sample and find the segments most suitable to act as base behaviors. This also means that we need to use the entire population as the iteration sample, which may result in an unexpected systematic error. Thus, to avoid any systematic error and for reasons stated in R2, we use an attribute set (to find anchors) that is different than the attribute set used to find the distribution profiles.

5.1.2 Attribute Set for finding base behaviors, AS_{BB}

A signal or its segment can also be treated as a wave that can be described by its amplitude and frequency. If unlike a regular wave, the signal segment shows a lot of amplitude variation then a single value/coefficient is not sufficient to describe the change in amplitude. In such cases, a more comprehensive change in amplitude is obtained by calculating and maintaining a sequence crest to trough distance. However, it is possible that each segment will have its own range of amplitude and that two different segments with different amplitude ranges can end up embodying the same pattern. At this stage of analysis, we are not interested in the actual crest to trough distances but rather how it changes relative to one another with progressing time. Thus, within a given segment, s_k , we define the relative change between two consecutive crest to trough distances, $\Delta_{CT_{j-1}}$ and Δ_{CT_j} as:

$${}^{s_k} \delta_{CT}^j = \frac{{}^{s_k} \Delta_{CT_j}}{{}^{s_k} \Delta_{CT_{j-1}}} \text{ where,} \quad (3)$$

$${}^{s_k} \Delta_{CT_j} = {}^{s_k} \text{ crest}_j - {}^{s_k} \text{ trough}_j$$

Hence, instead of maintaining a sequence of crest to trough distances, we maintain a sequence of ${}^{s_k} \delta_{CT}^j$ for each segment and define the first six attributes of AS_{BB} to be $[mean(\delta_{CT}),$

$var(\delta_{CT}), min(\delta_{CT}), max(\delta_{CT}), iqr(\delta_{CT}), median(\delta_{CT})$]. The seventh attribute of AS_{BB} represents the variations in frequency and defined as the number of peaks present in a segment, normalized by length of the segment. It is denoted by F . The eight and final attribute measures the correlation of the raw segment with any positive sloping line (to capture changes in trendline) and is denoted by T . The complete $AS_{BB} = [mean(\delta_{CT}), var(\delta_{CT}), min(\delta_{CT}), max(\delta_{CT}), iqr(\delta_{CT}), median(\delta_{CT}), F, T]$.

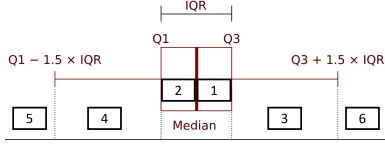


Figure 13: Feature Discretion Cutoff

In order to incorporate an insight of the distribution of each attribute across the entire dataset into AS_{BB} , each attribute is discretized using statistical quartiles and outliers, in accordance to the cut off shown in Figure 13. This discretization process reduces noise in the set of feature vectors obtained for base behavior selection.

5.1.3 Number of centroids/anchors

As mentioned earlier we want the base behaviors or centroids to be a minimal set. Referring to the number of base behaviors as q , in an experimental setup, we varied q (see Figure 13) between values of 1 and 10. We found that if $q > c$ and $c \in [1, 10]$, then we ended up with similar patterns, resulting in degeneration of the base behaviors. Thus we set $q = c$.

5.1.4 Back tracing the centroids to the segment

Since we are using different feature set for base behavior selection and for distribution profile creation (as explained in R2), we need to map the centroids/anchors, which are found in the AS_{BB} , back to the original segments. During the mapping process, we look for segments in their AS_{BB} format that is closest to the centroids. Once those segments are identified, they are marked to be used during the generation of distribution profiles.

Note, that this mapping is not necessarily one to one as there can be more than one segment (in their AS_{BB} format) that are closest and the same distance from a particular centroid. However, that does not affect our analysis because segments, which are equidistant from the centroids, present similar variability and are all equally suitable to be an anchor or a base behavior.

5.2 Phase 2: Creation of distribution profiles

In order to obtain the distribution profiles in terms of the base behavior for each entity/signal: **1)** we need to define the attribute set, and **2)** track the proportion of each of the q base behaviors or similar behaviors, present in the entity's data/segments.

5.2.1 Attribute Set for finding distribution profiles, AS_{DP}

The choice of feature vector, for this stage of analysis, is determined by the need to identify changes in patterns, which are often caught by "the shape parameters". Since

these shape parameters are captured succinctly by moments [10] we use the first five moments of the amplitude distribution. The *first moment* accounts for the DC (in the frequency domain) value or mean of the the signal segment; the *second moment* gives the average AC power; the *third moment* measures the skewness; the *fourth moment* or kurtosis, accounts for fluctuation in the power; and the *fifth moment* measures the heaviness of tail and mode of the distribution, given the skewness. The moment feature vectors are then discretized according to the cut off shown in Figure 13.

5.2.2 Tracking proximity of entity data to base behaviors

Once the base behaviors have been identified, we can iterate over the entity data or feature vectors in temporal order and match the feature vectors to the base behaviors.

$$f_{ij}^{t+1} = \frac{(f_{ij}^t \times n_i^t) + 1}{n_i^t + 1} \quad (4)$$

$$n_i^{t+1} = n_i^t + 1$$

At each iteration, we need to track the number of individual feature vector we have seen from a particular entity and also, record the number of times a feature vector from a particular entity is assigned to a particular base behavior, based on proximity measure. This tracking process is conveniently simplified by the use of *fingerprints* by the multi-set stream clustering framework. A fingerprint is a q dimensional vector, associated with each entity, such that the q bins track the frequency of the base behavior proximity assignment and is denoted as $[f_1, \dots, f_q]$. When a fingerprint is associated with a set of data points to track the number of feature vectors seen, it can be denoted as $[f_1, \dots, f_q, n]$. For example- suppose the $(t + 1)^{th}$ data point is a feature vector from $entity_i$ and is closest to base behavior, b_j , where $j \in [1, q]$. Then, f_{ij} and n_i is updated as shown in equation 4. Since, sum of f_{ij} over j is required to be 1 at any point in time, the completed fingerprints can be used as distribution profiles for the corresponding entities.

5.3 Phase 3: Creation of group profiles

Since it is possible to have a large number of entities, some of which are similar to one another, studying the entities as a group would provide a better insight. The multi-set streaming clustering framework suggests to dynamically group the fingerprints even as they are updated during each iteration. This allows the groups to change and evolve in a way that tracks progression over time. The groups are described by the fingerprint average of the members and are referred to as class profiles. For the purpose of this paper, we denote a particular class profile as $[p_{k1}, \dots, p_{kq}]$, where k identifies the group and the values $[1, \dots, q]$ refer to the base behaviors. For full description of class profiles as micro-clusters, please refer to [1].

To obtain the best separation of entities, k is allowed to run from 1 to $2^q - 1$ because it allows us to initialize the profiles with all possible binary combination of q zeros and ones. When p_{kj} is initialized to be non-zero, it attracts fingerprints that have a high frequency for the corresponding, b_j .

Note, the class profiles are similar to histogram and calculating similarity based only on euclidean distance results in a subtle discrepancy. For example, suppose we have three different histograms with four bins: $hist1 = [0.8, 0.2, 0, 0]$,

$hist2 = [0.6, 0.2, 0.2, 0]$ and $hist3 = [0.6, 0.4, 0, 0]$. The distances between the histograms are: $dis(hist1 - hist2) = 0.2828$, between $dis(hist1 - hist3) = 0.2828$ and between $dis(hist2 - hist3) = 0.2828$. However, for our application, it is required to have $hist1$ and $hist3$ to more similar to one another than to $hist2$ because their non-zero bins have a perfect match even though the bin proportions are different. As a result, for the similarity metric, we use the the Euclidean distance with a *penalty_term* for histograms with mismatched bins.

6. MEDICAL CHARACTERIZATION OF THE GROUP PROFILES

As mentioned earlier, the data for respiration-induced tumor motion is of importance to the radiation oncology community because it can be used to study three clinically significant patterns, *baseline*, *ES-range* and *D-range shifts*.

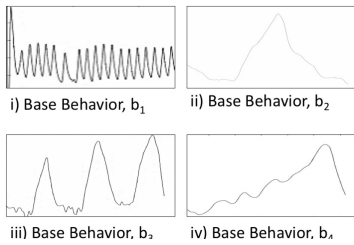


Figure 14: Selected base behaviors

The process of medical characterization of class profiles involves three steps: **1)** understanding the base behaviors, **2)** defining clinically significant patterns in terms of the base behaviors, and **3)** for ease of analysis, quantifying the presence of the clinical patterns in each group profile as an assigned score.

In order, to identify these three patterns in the class profiles, which is significantly of lower dimension than the actual entity signal, we need to re-examine the base behaviors that were selected during the clustering process. Figure 14 shows that b_1 sets up the comparison for high amplitude, while, b_2 represents abnormal breathing cycles. b_3 captures increasing sudden changes in mean position, coupled with widening range of amplitude. And finally, b_4 captures gradual changes in mean position and/or irregular changes in amplitude.

One obvious characteristic of the class profile, which needs to be taken into consideration, is that just the presence of a base behavior, b_j , by itself, does not give a comprehensive description of the entity. The complete description of the entity, is given by the presence of a particular b_j , along with the relative magnitudes of the base behaviors present. Since $q = 4$, there are six unique pairwise combinations of the base behaviors. Thus, based on the description of the three clinically significant patterns, obtained from our co-author, who is a medical physician, we define a *6-vector* in terms of b_j , which allows us to create templates for clinically significant patterns, without explicitly defining a combination ratio cutoff for the base behaviors. The *6-vector* is defined as: $[(b_1 - b_2), (b_1 - b_3), (b_1 - b_4), (b_2 - b_3), (b_2 - b_4), (b_3 - b_4)]$. For example, if we are looking for changes in mean position, we want the proportions of b_3 and b_4 to exceed the proportions of b_1 and b_2 and require the term, $(b_2 - b_3)$ to

be negative. If we are looking for D-range shifts, then we need the term, $(b_1 - b_3)$ to be positive. For mean position changes, the sign for the term, $(b_1 - b_2)$ does not matter; whereas for D-range shifts, the term should be positive and negative for ES-range shifts. Following a similar approach, we define the templates for the clinically significant patterns, using the base behaviors as follows:

- *Baseline shifts*: $\vec{b} = [0/\pm, -, -, -, -, 0/\pm]$
- *ES-range shifts*: $\vec{e}s = [-, 0/\pm, -, +, 0/-, -]$
- *D-range shifts*: $\vec{d} = [+ , + , + , -, 0/\pm, 0/\pm]$

Suppose, the *6-vector* for a particular class profile m is given by $v_m = [v_{m1}, v_{m2}, v_{m3}, v_{m4}, v_{m5}, v_{m6}]$. Then to quantify the clinical patterns, each class profile is matched to all three templates and three separate scores, $[score_b, score_{es}, score_d]$, are assigned to the class profile. If $j \in \{b, es, d\}$ and there is a sign match between v_{mi} and b_i or es_i or d_i , where i refers to a particular index in \vec{b} or $\vec{e}s$ or \vec{d} , then $score_j$ is updated by $|v_{mi}|$ amount. Thus, the final scores are defined as:

$$score_j = \frac{\sum_{i=1}^6 f(v_{mi}, j_i)}{\sum_{i=1}^6 |v_{mi}|}, \text{ where} \quad (5)$$

$$f(v_{mi}, j_i) = |v_{mi}|, \text{ if there is a sign match;} \\ \text{otherwise, } f(v_{mi}, j_i) = 0$$

7. RESULTS AND ANALYSIS

7.1 Implementation Details

VL Segmentation: δ_{W_i} is set to be proportional to $height(SW_{i3})$. Also in $W_{i-\alpha}$, α is set to be 50 in the current implementation. The size of the sliding window is set to 100 sample points because based on the dataset used for experimental analysis in this study, the sample rate of tumor motion is 25 Hz, and an average respiratory cycle has a duration of approximately 4 seconds and 100 samples points approximates 4 seconds. Also, fractions that showed only trivial variations, were segmented using fixed length segmentation as trivial variations are ignored by VL.

Setwise Clustering: During Phase 1, we find $q = 4$ using complete linkage clustering, with $q > 4$ resulting in degeneration of the anchors. Also, during back tracing the centroids to the segment, we empirically verified that for our dataset, the mapping is one to one. In Phase 2, the fingerprints are initialized using the first segments of each entity. In Phase 3, we end up with $2^q - 1 = 2^4 - 1 = 15$ different class profiles that are explained and analyzed in Sections 7.2, 7.3 and 7.4.

Result Validation: As mentioned earlier, one of our co-authors is a medical physician and the results presented in the following sections have been verified by his expertise.

7.2 Base behavior analysis

We described b_1 , b_2 , b_3 and b_4 in Section VI. The fifteen class profiles are presented as stacked bar charts in Figures 15 and 16, along with some examples of their respective entity members. The class profiles in Figures 15 and 16 are grouped according to a *kmeans++* clustering, which is performed on the profiles for ease of analysis and presentation and is explained later in this section. The amplitude of the signals in Figures 15 and 16 has been normalized to be between $[-1, 1]$ for clear comparison.

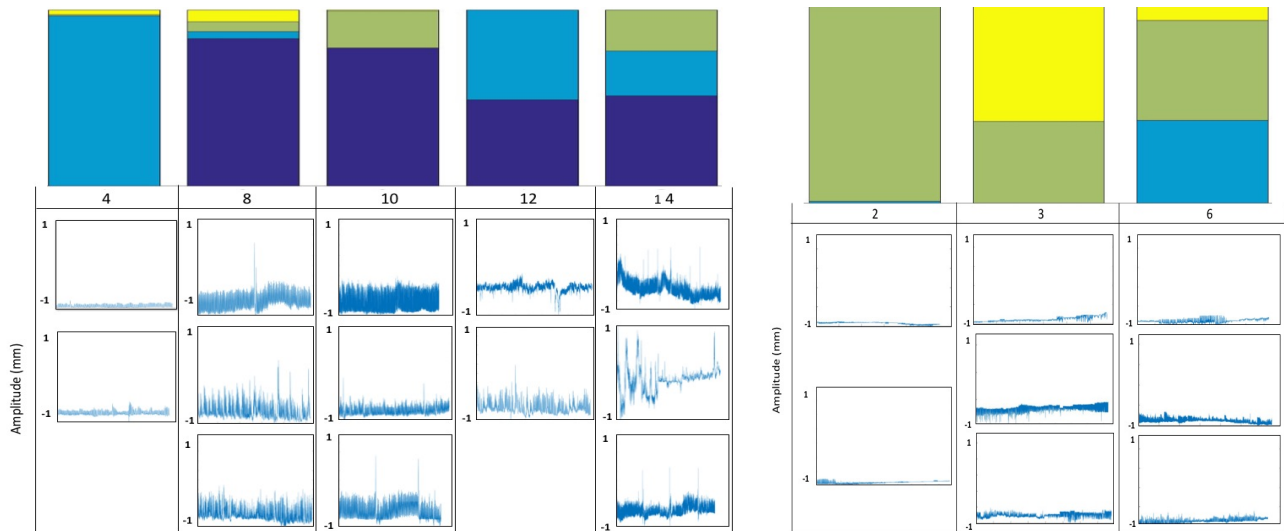


Figure 15: From left to right, meta-classes 1 and 2 respectively. The example of the members of each class profile is presented beneath the respective bar graph. In the reading the bar charts, dark blue represents the proportion of b_1 , light blue represents the proportion of b_2 , light green the proportion of b_3 and yellow represents the proportion of b_4 .

It is very clear from these two figures that any presence of b_3 predicts *simultaneous* changes in the mean position within a very narrow amplitude range and the proportion of b_3 determines whether or not a baseline shift is imminent. Presence of b_4 (without any presence of b_3) is indicative of a signal displaying gradual changes in mean position *and/or* amplitude variation. b_1 is representative of signals with a maximum amplitude higher than 1mm and its proportion distinguishes between spikes and consistent high amplitude. Finally, b_2 represents high variability in amplitude range and pattern.

7.3 Quantifying the Clinical Patterns

The evidence, found from base behavior analysis, further supports our definition of the 6-*vector* as: $[(b_1 - b_2) - (b_1 - b_3), (b_1 - b_4), (b_2 - b_3), (b_2 - b_4), (b_3 - b_4)]$ and *baseline shift* template as $[0/\pm, -, -, -, -, 0/\pm]$, *ES-range shift* as $[-, 0/\pm, -, +, 0/-, -]$ and *D-range shift* $[+, +, +, -, 0/\pm, 0/\pm]$. The score obtained from this template comparison (obtained by using the technique described in Section VI) is summed up in Table II. The values in brackets in Table II is the relative ranks of the scores. To obtain the relative rank, we sort the scores and rank the unique positions. There are fifteen different class id because we set $q = 4$ as described earlier.

The score for each pattern can be seen as fuzzy labeling/classification and can be interpreted as a probability measure of how likely the particular pattern will show up on any member fraction of the group. For example- member fractions of Class ID 2 will mostly display changes in mean position often leading to drastic baseline shifts and show amplitude variation because it has a high score for baseline shift and relatively low values (more prominently seen from the high relative ranking) for the other two patterns. Class ID 5 and 7, on the other hand, will show both high amplitudes and changes in mean position because both baseline shift and D-Range shift have high scores. However,

compared to members of Class ID 5, the changes in mean position will be more prominent for members of Class ID 7 because Class ID 7 has a lower relative rank for baseline shift than Class ID 5. Also, since Class ID 7 got a lower relative rank for ES-Range Shift than Class ID 7, the members of Class ID 7 are more likely to show ES-Range Shift.

In summary, we find that classes 1, 2, 3, 5, 7, 9, 11 and 15 show high similarity to baseline shift template. ES-Range Shift appears to be less polar between the groups with class 9 showing the highest similarity to the respective template. Classes 1, 5, 7, 9, 11, 13 and 15 show high similarity for the D-range shift template. From this table, we see that Classes 12 and 14 have the lowest similarity to baseline shift template; classes 2 and 6 does not show strong likeness to D-range shift template; and class 4 is most dissimilar with ES-range shift template. Furthermore, the relative ranking provides evidence that the class profiles can be grouped further. For example- classes 2 and 3 and classes 7 and 9 have very similar scores for all patterns.

7.4 Interpretation of the Quantification of Clinical Patterns

Table II shows us that classes 13 and 15 have very similar scores for marginal expansion and gating but their scores for baseline shift are very different. We want to know whether it is possible to group classes like 13 and 15 together into meta-classes, such that the meta-classes retain medical context and meaning.

This leads us to our next step of our analysis, which is to cluster the class profiles using `kmeans++` based on the scores presented in Table II. We try different values of k and find that $k > 3$ characterizes trivial variations only. With $k = 3$, we get three meta-classes:

- *Meta-Class 1* contains the class profiles with class id 4, 8, 10, 12, and 14, all of which display high amplitude. With a total of 29 fractions, the amplitude of

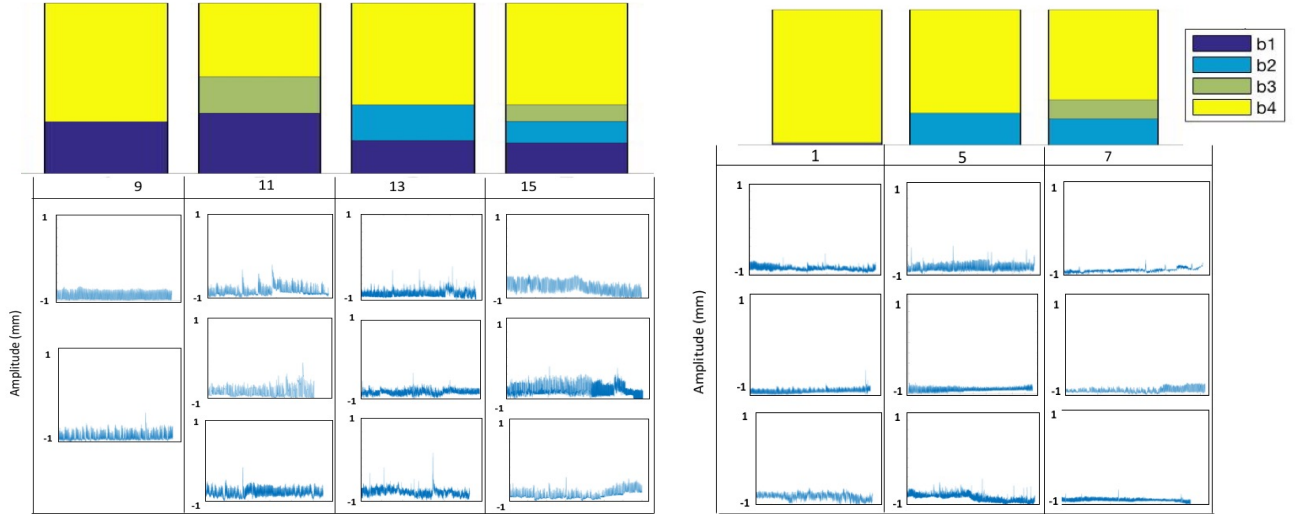


Figure 16: Meta-class 3 with profiles 1, 5, 7, 9, 11, 13, 15.

this meta-class, ranges from 1.2mm (class id 4) to 6mm (with one exception where the amplitude is 0.45mm). The best way to interpret this result is that these 29 fractions in group 2 dominantly display high amplitude (different from spikes), which might mask/devalue the presence of baseline and ES-range shifts.

- *Meta-Class 2* contains the class profiles with class id 2, 3, 6. With a total of 45 fractions, this meta-class dominantly exhibits baseline shifts with likely displays of some ES-range shifts as well. Class profile with class id 2 displays drastic baseline shifts. Class profile with class id 3, contains signals that display baseline shifts punctuated by spikes that are uncharacteristic of the fractions. Class profile with class id 6, is more likely to display shifts in mean position frequently for shorter duration, when compared to class profile with class id 2.
- *Meta-Class 3* contains the class profiles with class id 1, 5, 7, 9, 11, 13, 15. With a total of 86 fractions, represents a group where we the fractions show infrequent changes in baseline (gradual) *and/or* ES-range shifts that is present only for a short duration and is punctuated by spikes. None of the characters seem to dominate these group and any conclusion on dominance requires a closer look at individual class profiles. Figures 15 and 16 organize the profiles according to this grouping.

In order to test our methodology, we present our results on the data along SI axis only, because majority of the treatment fractions show dominant SI (Superior-Inferior) motion [18] [11]. We have also tested the methodology separately on data along AP and LR axes, which produces similar results. AP data shows similar anchor degeneration as SI data but displays a much narrower range of amplitude motion than SI motion (the mean is less than half of SI amplitude range). As a result, due to our method of discretizing feature vectors, the AP data does not pick an anchor high amplitudes. Data along the LR axis, picks similar anchors as SI data but

the anchors do not degenerate at $q = 5$; this is due to the cutoffs used for feature vector discretization. As result, LR data picks one extra anchor that adds resolution to changes in mean position.

The importance of **meta-classes 1, 2 and 3** arises due its ability to highlight the dominant character their members, which can be used to explore the breathing intervention techniques. Meta-class 1 makes exploring gating as a relevant breathing intervention treatment for the group. Members of meta-class 2 might benefit from the consideration of abdominal compression as the intervention technique. Meta class 3 might benefit from abdominal compression and marginal expansion and maybe, even gating.

Table II. Percentage Scores and Relative Rank. The class ID corresponds to the stacked bar charts in Figures 15-16. The numbers in bracket shows the relative ranking of the entire score table, with 1 being the highest rank.

Table II			
Class ID	Baseline Shift	ES-Range Shift	D-Range Shift
1	99.9 (1)	66.7 (17)	99.80 (1)
2	99.9 (1)	66.6 (17)	33.4 (29)
3	100 (1)	73.0 (13)	62.3 (20)
4	34.9 (27)	34.0 (28)	66.6 (17)
5	90.8 (5)	60.0 (21)	90.8 (5)
6	81.9 (10)	54.6 (22)	45.4 (26)
7	97.6 (3)	65.0 (19)	83.7 (9)
8	35.5 (27)	66.8 (17)	67.4 (17)
9	87.4 (7)	70.9 (14)	100 (1)
10	47.2 (25)	61.3 (21)	61.3 (21)
11	90.3 (6)	69.9 (15)	85.1 (8)
12	1.2 (31)	49.4 (24)	74.7 (12)
13	77.8 (11)	66.0 (18)	99.0 (2)
14	31.4 (30)	50.7 (23)	67.1 (17)
15	92.9 (4)	67.8 (16)	100 (1)

7.5 A Note on Method Performance

As we mentioned earlier, one of the biggest challenges with this dataset is that it neither has a control dataset nor is it annotated. If we had the control dataset, then we could use supervised learning methods to learn different patterns and

class boundaries. If the dataset was annotated, then even in the absence of control dataset, we could have reduced our result to a single accuracy value. Our previous works, [3], [16] and [4], are centered around this dataset. Even though [3] and [16] use very similar approaches as used in this paper, due to difference in target application, we cannot provide a comparison. [4] provides an automatic annotation technique and explores the readiness with which learning methods can learn the patterns if a proper annotation was provided. As a result, the only way we can validate our methods and analysis is through the expertise of a medical physician. One of our co-authors is a medical physician and at every step of the process, we relied on his expertise to determine medical relevance of our paper.

8. CONCLUSION

We conclude by saying that through a series of analysis we created low dimensional patient profiles with medical context from an uni-modal respiratory-induced tumor data. Our results were verified by the expertise of one of our co-authors, who is a medical physician. To obtain the result, we proposed a new adaptive segmentation technique that makes the correct compromises to capture persistent change in a time series signal. We then modify an existing multi-set clustering in various ways to create the patient class profiles. We follow this up by the quantification of the clinical patterns and its analysis within medical context. In our next work, we would like to combine information from all three axes and create the profiles. For future work, we could try to map the profiles to breathing treatment interventions which are subject to change over time. We could also attempt to apply this technique to other bio-sensor dataset (once available) and test the generalizability of the technique.

9. ACKNOWLEDGEMENTS

This material is based upon work supported by the National Science Foundation (NSF) under Grant No. 1626586. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF.

10. REFERENCES

- [1] C. Aggarwal. The multi-set stream clustering problem. *Society for Industrial and Applied Mathematics. Proceedings of the SIAM International Conference on Data Mining*, pages 59–69, 2012.
- [2] C. Aggarwal. The setwise stream classification problem. *20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 432–441, 2014.
- [3] A. S. Arvind Balasubramanian, Balakrishnan Prabhakaran. Mining pattern sequences in respiratory tumor motion data. *IEEE Engineering in Medicine and Biology Society*, pages 5262 – 5265, 2012.
- [4] Y. C. A. S. B. P. Arvind Balasubramanian, Rittika Shamsuddin. Exploring baseline shift prediction in respiration induced tumor motion. *IEEE International Conference on Healthcare Informatics*, pages 155 – 160, 2014.
- [5] J. T. D. Caster, F. Mosteller. *Understanding robust and explanatory data analysis*. 2001.
- [6] M. M. D. Wang, R. Vogt and S. Sridharan. Automatic audio segmentation using the generalized likelihood ratio. *International Conference on Signal Processing and Communication Systems*, pages 1–5, 2008.
- [7] T. K. R. C. J. R. D. B. L. I. DP. Steinfort, S. Siva. Multimodality guidance for accurate bronchoscopic insertion of fiducial markers. *Journal of Thoracic Oncology*, pages 324–330, 2015.
- [8] K. M. H. Azami and B. Bozorgtabar. An improved signal segmentation using moving average and savitzky-golay filter. *Journal of Signal and Information Processing*, pages 39–44, 2012.
- [9] K. M. H. Azami and H. Hassanpour. An improved signal segmentation method using genetic algorithm. *International Journal of Computer Applications*, pages 5–9, 2011.
- [10] W. M. Hartmann. *Signals, Sound, and Sensation*. Springer, 1998.
- [11] M. M. J. Adler, S. Chang. The cyberknife: A frameless robotic system for radiosurgery. *Stereotact Funct Neurosurg*, pages 124–8, 1997.
- [12] W. F. J. Bezdek, R. Ehrlich. Fcm: The fuzzy c-means clustering algorithm. *Computer and Geosciences*, pages 191–203, 1984.
- [13] M. and A. Gaya. Stereotactic body radiotherapy: a review. *Clinical Oncology*, pages 157–72, 2010.
- [14] S. K. G. G. G. H. M. Kirlangic, D. Perez and G. Ivanova. Fractal dimension as a feature for adaptive electroencephalogram segmentation in epilepsy. *Engineering in Medicine and Biology Society*, pages 1573–1576, 2001.
- [15] H. F. Mehmet Ali Abbasoglu, Bugra Gedik. Aggregate profile clustering for telco analytics. *Journal Proceedings of the VLDB Endowment*, pages 1234–1237, 2013.
- [16] A. S. R. Shamsuddin, A. Balasubramanian and B. Prabhakaran. Calculating patient similarity based on respiration induced tumor motion. *IEEE International Conference on Healthcare Informatics 2015*, pages 122–129, 2015.
- [17] Y. Wang, P. Wang, J. Pei, W. Wang, and S. Huang. A data-adaptive and dynamic segmentation index for whole matching on time series. *Journal Proceedings of the VLDB Endowment*, pages 793–804, August 2013.
- [18] B. C. Y. Suh, S. Dieterich and P. Keall. An analysis of thoracic and abdominal tumour motion for stereotactic body radiotherapy patients. *Physics in medicine and biology*, page 3623, 2008.
- [19] J. F. G. C. H. K. T. N. C. S. T. J. W. L. Y. M. Z. Zheng Jye Ling, Quoc Trung Tran. Gemini: an integrative healthcare analytics system. *Journal Proceedings of the VLDB Endowment*, pages 1766–1771, 2014.