

# VARIATIONAL PSOM: DEEP PROBABILISTIC CLUSTERING WITH SELF-ORGANIZING MAPS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1       Generating visualizations and interpretations from high-dimensional data is a  
2       common problem in many fields. Two key approaches for tackling this prob-  
3       lem are clustering and representation learning. There are very performant deep  
4       clustering models on the one hand and interpretable representation learning tech-  
5       niques, often relying on latent topological structures such as self-organizing maps,  
6       on the other hand. However, current methods do not yet successfully *combine*  
7       these two approaches. We present a new deep architecture for probabilistic clus-  
8       tering, VarPSOM, and its extension to time series data, VarTPSOM, composed of  
9       VarPSOM modules connected by LSTM cells. We show that they achieve supe-  
10      rior clustering performance compared to current deep clustering methods on static  
11      MNIST/Fashion-MNIST data as well as medical time series, while inducing an  
12      interpretable representation. Moreover, on the medical time series, VarTPSOM  
13      successfully predicts future trajectories in the original data space.

## 14 1 INTRODUCTION

15 Information visualization techniques are essential in areas where humans have to make decisions  
16 based on large amounts of complex data. Their goal is to find an interpretable representation of  
17 the data that allows the integration of humans into the data exploration process. This encourages  
18 visual discoveries of relationships in the data and provides guidance to downstream tasks. In this  
19 way, a much higher degree of confidence in the findings of the exploration is attained (Keim, 2002).  
20 An interpretable representation of the data, in which the underlying factors are easily visualized, is  
21 particularly important in domains where the reason for obtaining a certain prediction is as valuable  
22 as the prediction itself. However, finding a meaningful representation of complex data that can be  
23 understood by humans is challenging.

24 Clustering is one of the most natural ways for retrieving interpretable information from raw data.  
25 Long-established methods such as K-means (MacQueen, 1967) and Gaussian Mixture Models  
26 (Bishop, 2006) represent the cornerstone of cluster analysis. Their applicability, however, is of-  
27 ten constrained to simple data and their performance limited in high-dimensional, complex, real  
28 world data-sets, which do not exhibit a clustering-friendly structure.

29 Deep generative models have recently achieved tremendous success in representation learning.  
30 Some of the most commonly used and efficient approaches are Autoencoders (AEs), Variational  
31 Autoencoders (VAEs) and Generative Adversarial Networks (GANs) (Kingma & Welling, 2013;  
32 Goodfellow et al., 2014). The compressed latent representation, generated by these models, has  
33 been proven to ease the clustering process (Aljalbout et al., 2018). As a result, the combination of  
34 deep generative models for feature extraction and clustering results in a dramatic increase of the  
35 clustering performance (Xie et al., 2015). Although very successful, most of these methods do not  
36 investigate the relationship among clusters and the clustered feature points live in a high-dimensional  
37 latent space that cannot be easily observed or interpreted by humans.

38 The Self-Organizing Map (SOM) (Kohonen, 1990) is a clustering method that provides such an  
39 interpretable representation. It arranges the obtained centroids in a topologically meaningful order,  
40 inducing a flexible neighbourhood structure. If the chosen topological structure is a 2-dimensional  
41 grid, it facilitates visualization. Alas, its applicability is often constrained to simple data-sets similar  
42 to other classical clustering methods.

43 To resolve the above issues, we propose a novel deep architecture, the Variational Probabilistic SOM  
44 (VarPSOM), that jointly trains a VAE and a SOM to achieve an interpretable discrete representation  
45 while exhibiting state-of-the-art clustering performance. Instead of hard assignment of data points

46 to clusters, our model uses a centroid-based probability distribution. It minimizes its Kullback-  
 47 Leibler divergence against an auxiliary target distribution, while enforcing a SOM-friendly space.  
 48 To highlight the importance of an interpretable representation for different purposes, we extended  
 49 this model to deal with temporal data, yielding VarTPSOM. We discuss related work in Section  
 50 2. Extensive evidence of the superior clustering performance of both models, on MNIST/Fashion-  
 51 MNIST images as well as real-world medical time series is presented in Section 4.

52 Our main contributions are:

- 53 • A novel architecture for deep clustering, yielding an interpretable discrete representation  
 54 through the use of a probabilistic self-organizing map.
- 55 • An extension of this architecture to time series, improving clustering performance on this  
 56 data type and enabling temporal predictions.
- 57 • A thorough empirical assessment of our proposed models, showing superior performance  
 58 on benchmark tasks and challenging medical time series from the intensive care unit.

## 59 2 RELATED WORK

60 Self-Organizing Maps have been widely used as a means to visualize information from large  
 61 amounts of data (Tirunagari et al., 2014) and as a form of clustering in which the centroids are  
 62 connected by a topological neighborhood structure (Flexer, 1999). Since their early inception, sev-  
 63 eral variants have been proposed to enhance their performance and scope. The adaptive subspace  
 64 SOM, ASSOM (Kohonen, 1995), for example, proposed to combine PCA and SOMs to map data  
 65 into a reduced feature space. Tokunaga & Furukawa (2009) combine SOMs with multi-layer percep-  
 66 trons to obtain a modular network. Liu et al. (2015) proposed Deep SOM (DSOM), an architecture  
 67 composed of multiple layers similar to Deep Neural Networks. There exist several methods tailored  
 68 to representation learning on time series, among them (Fortuin & Rättsch, 2019; Fortuin et al., 2019),  
 69 which are however not based on SOMs. Extensions of SOM optimized for temporal data include  
 70 the Temporal Kohonen map (Chappell & Taylor, 1993) and its improved version Recurrent SOM  
 71 (McQueen et al., 2004) as well as Recursive SOM (Voegtlin, 2002). While SOM and its variants  
 72 are particularly effective for data visualization (Liu et al., 2015), it was rarely attempted to combine  
 73 their merits in this respect with modern state-of-the-art clustering methods, which often use deep  
 74 generative models in combination with probabilistic clustering.

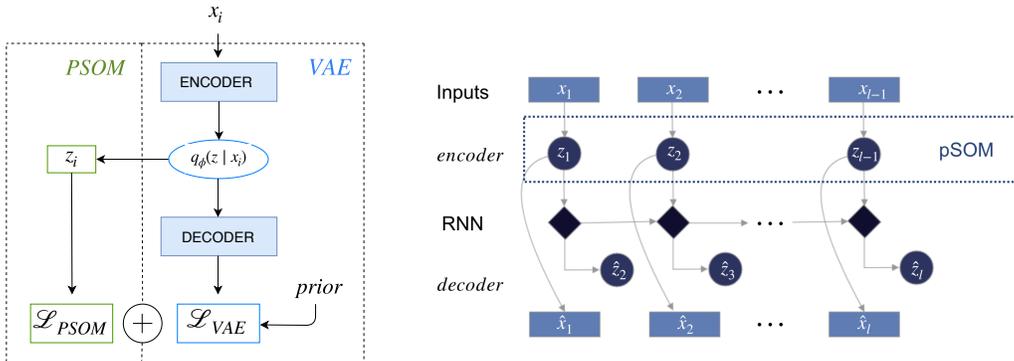
75 In particular, recent works on clustering analysis have shown that combining clustering algorithms  
 76 with the latent space of AEs greatly increases the clustering performance (Aljalbout et al., 2018). Xie  
 77 et al. (2015) proposed DEC, a method that sequentially applies embedding learning using Stacked  
 78 Autoencoders (SAE), and the *Clustering Assignment Hardening* method on the obtained representa-  
 79 tion. An improvement of this architecture, IDEC, (Guo et al., 2017), includes the decoder network  
 80 of the SAE in the learning process, so that training is affected by both the clustering loss and the  
 81 reconstruction loss. Similarly, DCN (Yang et al., 2016) combines a K-means clustering loss with the  
 82 reconstruction loss of SAE to obtain an end-to-end architecture that jointly trains representations and  
 83 clustering. These models achieve state-of-the-art clustering performance but they do not investigate  
 84 the relationship among clusters. An exception is the work by Li et al. (2018), in which they present  
 85 an unsupervised method that learns latent embeddings and discovers multi-facet clustering structure.  
 86 Relationships among clusters were discovered, however, they do not provide a latent space that can  
 87 be easily interpreted and which eases the process of analytical reasoning.

88 To the best of our knowledge, only two models used deep generative models in combination with a  
 89 SOM structure in the latent space. The SOM-VAE model (Fortuin et al., 2018), inspired by the VQ-  
 90 VAE architecture (van den Oord et al., 2017), uses an AE to embed the input data points into a latent  
 91 space and then applies a SOM-based clustering loss on top of this latent representation. It features  
 92 hard assignments of points to centroids, as well as the use of a Markov model for temporal data,  
 93 which both reduces modeling power compared to our method. The Deep Embedded SOM, DESOM  
 94 (Forest et al., 2019) improved the previous model by using a Gaussian neighborhood window with  
 95 exponential radius decay and by learning the SOM structure in a continuous setting. Both methods  
 96 feature a topologically interpretable neighborhood structure and yield promising results in visualiz-  
 97 ing state spaces. However their clustering quality is likely limited by the absence of techniques used  
 98 in state-of-the-art clustering methods like IDEC or DCN.

99 **3 PROBABILISTIC CLUSTERING WITH VARIATIONAL PSOM**

100 Given a set of data samples  $\{x_i\}_{i=1,\dots,N}$ , where  $x_i \in \mathbb{R}^M$ , the goal is to partition the data into a set  
 101 of clusters  $\{S_i\}_{i=1,\dots,K}$  while retaining a topological structure over the cluster centroids.

102 The proposed architecture for static data is presented in Figure 1a. The input vector  $x_i$  is embedded  
 103 into a latent representation  $z_i$  using a VAE. This latent vector is then clustered using *PSOM*, a  
 104 new SOM clustering strategy that extends the *Clustering assignment hardening* method (Xie et al.,  
 105 2015). The VAE and PSOM are trained jointly to learn a latent representation with the aim to boost  
 106 the clustering performance. To prevent the network from outputting a trivial solution, the decoder  
 107 network reconstructs the input from the latent embedding, encouraging it to be as similar as possible  
 108 to the original input. The obtained loss function is a linear combination of the clustering loss and  
 109 the reconstruction loss. To deal with temporal data, we propose another model variant, which is  
 110 depicted in Figure 1b.



(a) VarPSOM architecture for clustering of static data. Data points  $x_i$  are mapped to a continuous embedding  $z_i$  using a VAE (parameterized by  $\Phi$ ). The loss function is the sum of a SOM-based clustering loss and the ELBO.

(b) VarTPSOM architecture, composed of VarPSOM modules connected by LSTMs across the time axis, which predict the continuous embedding  $z_{t+1}$  of the next time step. This architecture allows to unroll future trajectories in the latent space as well as the original data space by reconstructing the  $x_t$  using the VAE.

Figure 1: Model architectures of VarPSOM / VarTPSOM

111 **3.1 BACKGROUND**

112 A Self-Organizing Map is comprised of  $k$  nodes connected to form a grid  $M \in \mathbb{N}^2$ , where the node  
 113  $m_{i,j}$ , at position  $(i, j)$  of the grid, corresponds to a centroid vector,  $\mu_{i,j}$  in the input space. The  
 114 centroids are tied by a neighborhood relation  $N(\mu_{i,j}) = \{\mu_{i-1,j}, \mu_{i+1,j}, \mu_{i,j-1}, \mu_{i,j+1}\}$ . Given a  
 115 random initialization of the centroids, the SOM algorithm randomly selects an input  $x_i$  and updates  
 116 both its closest centroid  $\mu_{i,j}$  and its neighbors  $N(\mu_{i,j})$  to move them closer to  $x_i$ . For a complete  
 117 description of the SOM algorithm, we refer to the appendix (A).

The *Clustering Assignment Hardening* method has been recently introduced by the DEC model (Xie et al., 2015) and was shown to perform well in the latent space of AEs (Aljalbout et al., 2018). Given an embedding function  $z_i = f(x_i)$ , it uses a Student’s t-distribution ( $S$ ) as a kernel to measure the similarity between an embedded data point  $z_i$ , and a centroid  $\mu_j$ :

$$s_{ij} = \frac{\left(1 + \|z_i - \mu_j\|^2 / \alpha\right)^{-\frac{\alpha+1}{2}}}{\sum_{j'} \left(1 + \|z_i - \mu_{j'}\|^2 / \alpha\right)^{-\frac{\alpha+1}{2}}}$$

It improves the cluster purity by enforcing the distribution  $S$  to approach a target distribution,  $T$ :

$$t_{ij} = \frac{s_{ij}^\gamma / \sum_i s_{ij}}{\sum_{j'} s_{ij'}^\gamma / \sum_i s_{ij'}}.$$

By taking the original distribution to the power  $\gamma$  and normalizing it, the target distribution puts more emphasis on data points that are assigned a high confidence. We follow (Xie et al., 2015) in choosing  $\gamma=2$ , which leads to larger gradient contributions of points close to cluster centers, as they show empirically. The resulting clustering loss is defined as:

$$\mathcal{L} = KL(T\|S) = \sum_i \sum_j t_{ij} \log \frac{t_{ij}}{s_{ij}}. \quad (1)$$

### 118 3.2 PROBABILISTIC SOM (PSOM) CLUSTERING

Our proposed clustering method, called PSOM, expands *Clustering Assignment Hardening* to include a SOM neighborhood structure over the centroids. We add an additional loss to (1) to achieve an interpretable representation. This loss term maximizes the similarity between each data point and the neighbors of the closest centroids. For each embedded data point,  $z_i$ , and each centroid  $\mu_j$  the loss is defined as the negative sum of all the neighbors of  $\mu_j$ ,  $\{e : \mu_e \in N(\mu_j(x_i))\}$ , of the probability that  $z_i$  is assigned to  $e$ , defined as  $s_{ie}$ . This sum is weighted by the similarity between  $z_i$  and the centroid  $\mu_j$  ( $s_{ij}$ ):

$$\mathcal{L}_{\text{SOM}} = -\frac{1}{N} \sum_i \sum_j s_{ij} \sum_{e: \mu_e \in N(\mu_j(x_i))} s_{ie}.$$

The complete PSOM clustering loss is then:

$$\mathcal{L}_{\text{PSOM}} = KL(T\|S) + \beta \mathcal{L}_{\text{SOM}}.$$

119 We note that for  $\beta = 0$  it becomes equivalent to Clustering assignment hardening.

### 120 3.3 VARPSOM: VAE FOR FEATURE EXTRACTION

In our method the nonlinear mapping between the input  $x_i$  and embedding  $z_i$  is realized by a VAE. Instead of directly embedding the input  $x_i$  into a latent embedding  $z_i$ , the VAE learns a probability distribution  $q_\phi(z | x_i)$  parametrized as a multivariate normal distribution whose mean and variance are  $(\mu_\phi, \Sigma_\phi) = f_\phi(x_i)$ . Similarly, it also learns the probability distribution of the reconstructed output given a sampled latent embedding,  $p_\theta(x_i | z)$  where  $(\mu_\theta, \Sigma_\theta) = f_\theta(z_i)$ . Both  $f_\phi$  and  $f_\theta$  are neural networks, called respectively encoder and decoder. The ELBO loss is:

$$\mathcal{L}_{\text{ELBO}} = \sum_i [-\mathbb{E}_z(\log p_\theta(x_i | z)) + D_{KL}(q_\phi(z | x_i) \| p(z))], \quad (2)$$

where  $p(z)$  is an isotropic Gaussian prior over the latent embeddings. The second term can be interpreted as a form of regularization, which encourages the latent space to be compact. For each data point  $x_i$  the latent embedding  $z_i$  is sampled from  $q_\phi(z | x_i)$ . Adding the ELBO loss to the PSOM loss from the previous subsection, we yield the overall loss function of VarPSOM:

$$\mathcal{L}_{\text{VarPSOM}} = \mathcal{L}_{\text{PSOM}} + \mathcal{L}_{\text{ELBO}}. \quad (3)$$

121

122 To the best of our knowledge, no previous SOM methods attempted to use a VAE to embed the inputs  
 123 into a latent space. There are many advantages of a VAE over an AE for realizing our goals. Most  
 124 importantly, learning a probability distribution over the embedding space improves interpretability  
 125 of the model. For example, points with a higher variance in the latent space could be identified as  
 126 potential outliers and therefore treated as less precise and trustworthy. Moreover, the regularization  
 127 term of the VAE prevents the network from scattering the embedded points discontinuously in the  
 128 latent space, which naturally facilitates the fitting of the SOM. To test if the use of CNNs can boost  
 129 clustering performance on image data, we introduce another model variant called VarCPSOM, which  
 130 uses convolutional filters as part of the VAE.

131 3.4 VARTPSOM: EXTENSION TO TIME SERIES DATA

To extend our proposed model to time series data, we add a temporal component to the architecture. Given a set of  $N$  time series of length  $T$ ,  $\{x_{t,i}\}_{t=1,\dots,T;i=1,\dots,N}$ , the goal is to learn interpretable trajectories on the SOM grid. To do so, the VarPSOM could be used directly but it would treat each time-step  $t$  of the time series independently, which is undesirable. To exploit temporal information and enforce smoothness in the trajectories, we add an additional loss to (3):

$$\mathcal{L}_{\text{smooth}} = -\frac{1}{NT} \sum_i \sum_t u_{i_t, i_{t+1}}, \quad (4)$$

132 where  $u_{i_t, i_{t+1}} = g(z_{i_t, t}, z_{i_{t+1}, t+1})$  is the similarity between  $z_i$  and  $z_j$  using a Student's t-distribution  
 133 and  $z_{i_t, t}$  refers to the embedding of time series  $x_i$  at time index  $t$ . It maximizes the similarity between  
 134 latent embeddings of adjacent time steps, such that large jumps in the latent state between time points  
 135 are discouraged.

One of the main goals in time series modeling is to predict future data points, or alternatively, future embeddings. This can be achieved by adding a long short-term memory network (LSTM) across the latent embeddings of the time series, as shown in Fig 1b. Each cell of the LSTM takes as input the latent embedding of time-step  $t$  ( $z_t$ ), and predicts a probability distribution over the next latent embedding,  $p_\omega(z_{t+1} | z_t)$ . We parametrize this distribution as a Multivariate Normal Distribution whose mean and variance are learnt by the LSTM. The prediction loss is the log-likelihood between the learned distribution and a sample of the next embedding  $z_{t+1}$ :

$$\mathcal{L}_{\text{pred}} = -\sum_i \sum_t \log p_\omega(z_{t+1} | z_t) \quad (5)$$

The final loss of VarTPSOM, which is trainable in a fully end-to-end fashion, is

$$\mathcal{L}_{\text{VarTPSOM}} = \mathcal{L}_{\text{VarPSOM}} + \mathcal{L}_{\text{smooth}} + \eta \mathcal{L}_{\text{pred}}. \quad (6)$$

136

137 4 EXPERIMENTS

138 First, we evaluate VarPSOM and VarCPSOM and compare them with state-of-the-art classical/SOM-  
 139 based clustering methods on MNIST (Lecun et al., 1998) and Fashion-MNIST (Xiao et al., 2017)  
 140 data. Hereby, particular focus is laid on the comparison of VarPSOM and the clustering models DEC  
 141 and IDEC, to investigate the role of the VAE and the SOM loss. We then present visualizations of the  
 142 obtained 2D representations, to illustrate how our method could ease visual reasoning about the data.  
 143 Finally, we present extensive evidence of the performance of VarTPSOM on real-world complex  
 144 time series from the eICU data set (Pollard et al., 2018), and illustrate how it allows visualization of  
 145 patient health state trajectories in an easily understandable 2D domain. For details on the data-sets,  
 146 we refer to the appendix (B.1).

147 **Baselines** We used two different types of baselines. The first category contains clustering methods  
 148 that do not provide any interpretable discrete latent representation. Those include K-means, the DEC  
 149 model, as well as its improved version IDEC, whose clustering methods are related to ours. We also  
 150 include a modified version of IDEC that we call VarIDEC, in which we substitute the AE with a  
 151 VAE, to investigate the role of the VAE in our method. For all these methods we use 64 clusters. In  
 152 the second category, we include state-of-the-art clustering methods based on SOMs. Here, we used  
 153 a standard SOM (minisom), AE+SOM, an architecture composed of an AE and a SOM applied on  
 154 top of the latent representation (trained sequentially), SOM-VAE and DESOM. For all SOM-based  
 155 methods we set the SOM grid size to  $(8 \times 8)$ .

156 **Implementation** In implementing our models we focused on retaining a fair comparison with the  
 157 baselines. Hence we decided to use a standard network structure, with fully connected layers of  
 158 dimensions  $d - 500 - 500 - 2000 - l$ , to implement both the VAE of our models and the AE of  
 159 the baselines. The latent dimension,  $l$ , is set to 100 for the VAE, and to 10 for the AEs. Since the  
 160 prior in the VAE enforces the latent embeddings to be compact, it also requires more dimensions  
 161 to learn a meaningful latent space. On the other hand, setting the AEs with a higher latent space,  
 162 needed for the VAE, resulted in a dramatic decrease of performance (see appendix B.2). VarCPSOM  
 163 is composed of 4 convolutional layers of feature maps [32, 64, 128, 256] and kernel size  $3 \times 3$  for all

Table 1: Clustering performance of VarPSOM using 64 clusters arranged in a  $8 \times 8$  SOM map, compared with baselines. The methods are grouped into approaches with no topological structure in the discrete latent space and interpretable methods using a SOM-based structure in the latent space, as well as an extension of our method using convolutional filters. Means and standard errors across 10 runs with different random model initializations are displayed.

	MNIST		fMNIST	
	<i>pur</i>	<i>nmi</i>	<i>pur</i>	<i>nmi</i>
Kmeans	$0.845 \pm 0.001$	$0.581 \pm 0.001$	$0.716 \pm 0.001$	$0.514 \pm 0.000$
DEC	$0.944 \pm 0.002$	$0.682 \pm 0.001$	$0.758 \pm 0.002$	$0.562 \pm 0.001$
IDEC	$0.950 \pm 0.001$	$0.681 \pm 0.001$	-	-
VarIDEC (ours)	<b><math>0.961 \pm 0.002</math></b>	<b><math>0.698 \pm 0.001</math></b>	<b><math>0.765 \pm 0.003</math></b>	<b><math>0.569 \pm 0.002</math></b>
SOM	$0.701 \pm 0.005$	$0.539 \pm 0.002$	$0.667 \pm 0.003$	$0.525 \pm 0.001$
AE+SOM	$0.874 \pm 0.004$	$0.646 \pm 0.001$	$0.706 \pm 0.002$	$0.543 \pm 0.001$
SOM-VAE	$0.868 \pm 0.004$	$0.595 \pm 0.004$	$0.739 \pm 0.005$	$0.520 \pm 0.003$
DESOM	0.939	0.657	0.752	0.538
VarPSOM (ours)	<b><math>0.964 \pm 0.001</math></b>	<b><math>0.705 \pm 0.001</math></b>	<b><math>0.764 \pm 0.003</math></b>	<b><math>0.571 \pm 0.001</math></b>
VarCPSOM (ours)	<b><math>0.980 \pm 0.001</math></b>	<b><math>0.726 \pm 0.001</math></b>	<b><math>0.783 \pm 0.003</math></b>	<b><math>0.574 \pm 0.001</math></b>

164 layers. For all architectures, no greedy layer-wise pretraining was used to tune the VAE. Instead we  
 165 simply run the VAE without the clustering loss for a few epochs for initialization. A standard SOM  
 166 was then used to produce an initial configuration of the centroids/neighbourhood relation. Finally,  
 167 the entire architecture is trained for 100,000 iterations. To avoid fine-tuning hyperparameters, given  
 168 the unsupervised setting,  $\alpha$  is set to 10 for all experiments while the other hyperparameters are  
 169 modified accordingly to maintain the same order of magnitude of the different loss components.

170 **Clustering Evaluation** Table 1 shows the clustering quality results of VarPSOM and VarCPSOM  
 171 on MNIST and Fashion-MNIST data, compared with the baselines. Purity and Normalized Mutual  
 172 Information are used as evaluation metrics. We observe that our proposed models outperform the  
 173 baselines of both categories and reach state-of-the-art clustering performance.

174 **VarPSOM vs. IDEC** VarIDEC shows superior clustering performance compared to DEC and  
 175 IDEC (Table 1). We conclude that the VAE indeed succeeds in capturing a more meaningful latent  
 176 representation compared to a standard AE. Regarding the second difference, the SOM structure was  
 177 expected to slightly decrease the clustering performance, due to a trade-off between interpretability  
 178 and raw clustering power. However, we do not observe this in our results. Adding the SOM loss  
 179 rather leads to an increase of the clustering performance. We suspect this is due to the regularization  
 180 effect of the SOM’s topological structure. Overall, VarPSOM outperforms both DEC and IDEC.

181 **Improvement over Training** After obtaining the initial configuration of the SOM structure, both  
 182 clustering and feature extraction using the VAE are trained jointly. To illustrate that our architecture  
 183 improves clustering performance over the initial configuration, we plotted NMI and Purity against  
 184 the number of training iterations in Figure 2. We observe that performance is stable when increasing  
 185 number the number of epochs and no overfitting is visible.

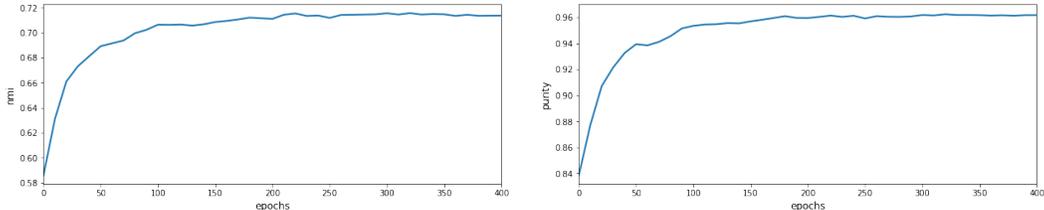


Figure 2: NMI (left) and Purity (right) performance of VarPSOM over iterations on MNIST test set.

186 **Role of the SOM loss** To investigate the influence of the SOM loss component, we plot the clus-  
 187 tering performance of VarPSOM against the weight ( $\beta$ ) of  $\mathcal{L}_{\text{SOM}}$  in Fig. 3, using MNIST dataset.

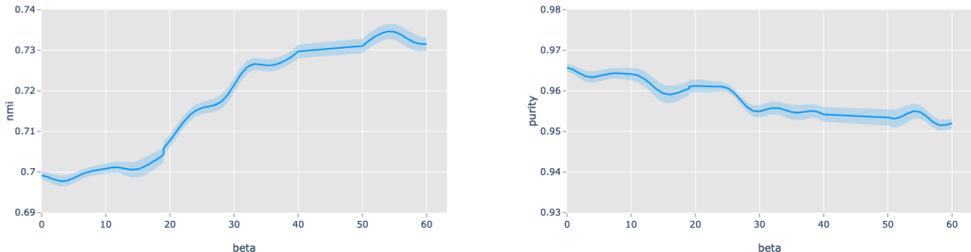
Table 2: Mean NMI and standard error of cluster enrichment vs. current/future APACHE physiology scores, using a 2D ( $8 \times 8$ ) SOM map, across 10 runs with different random model initializations.

Model	Apache12	Apache6	Apache0
SOM-VAE	0.0444 $\pm$ 0.0006	0.0474 $\pm$ 0.0005	0.0510 $\pm$ 0.0005
VarPSOM	0.0631 $\pm$ 0.0008	0.0639 $\pm$ 0.0008	0.0730 $\pm$ 0.0009
VarTPSOM ( $\eta = 0$ )	0.0710 $\pm$ 0.0005	0.0719 $\pm$ 0.0006	0.0818 $\pm$ 0.0006
VarTPSOM	<b>0.0719 <math>\pm</math> 0.0004</b>	<b>0.0733 <math>\pm</math> 0.0004</b>	<b>0.0841 <math>\pm</math> 0.0005</b>

Table 3: MSE for predicting the time series of the last 6 hours before ICU dispatch, given the prior time series since ICU admission.

Model	LSTM	SameState	VarTPSOM
MSE	0.0386 $\pm$ 0.0049	0.0576 $\pm$ 0.0012	<b>0.0297 <math>\pm</math> 0.0009</b>

188 With  $\beta = 30$ , the  $KL$  term (responsible for improving clustering purity) and the  $\mathcal{L}_{SOM}$  term (respon-  
 189 sible for enforcing a SOM structure over the centroids) are almost equal. It is interesting to observe  
 190 the different trends in NMI and Purity. The NMI performance increases for increasing values of  
 191  $\beta$  while Purity slightly decreases. Overall, enforcing a more interpretable latent space results in a  
 192 more robust clustering model with higher NMI clustering performance.

Figure 3: NMI (left) and Purity (right) performance of VarPSOM, with standard error, over  $\beta$  values on MNIST test set.

193 **Time Series Evaluation** We evaluate the clustering performance of our proposed models on the  
 194 eICU dataset, comprised of complex medical time series. We compare them against SOM-VAE,  
 195 as this is the only method among the baselines that is suited for temporal data. Table 2 shows the  
 196 cluster cell enrichment in terms of NMI for three different labels, the current (APACHE-0) and worst  
 197 future (APACHE-6/12 hours) physiology scores. VarTPSOM clearly achieves superior clustering  
 198 performance compared to SOM-VAE. This, we hypothesize, is due to the better feature extraction  
 199 using a VAE as well as the improved treatment of uncertainty using PSOM, which features soft  
 200 assignments, whereas SOM-VAE contains a deterministic AE and hard assignments. Moreover,  
 201 both the smoothness loss and the prediction loss seem to increase the clustering performance. More  
 202 results on ICU time series are contained in the appendix (B.3).

203 To quantify the performance of VarTPSOM in unrolling future trajectories, we predict the final  
 204 6 latent embeddings of each time series. For each predicted embedding we reconstruct the input  
 205 using the decoder of the VAE. Finally we measure the MSE between the original input and the  
 206 reconstructed inputs for the last 6 hours of the ICU admission. As baselines, we used an LSTM that  
 207 takes as input the first 66 hours of the time series and then predicts the next 6 hours. Since most  
 208 of the trajectories tend to stay in the same state over long periods of time, another strong baseline  
 209 is obtained by duplicating the last seen embedding over the final 6 hours. The results (Table 3)  
 210 indicate that the joint training of clustering and prediction used by VarTPSOM clearly outperforms  
 211 the 2 baselines.

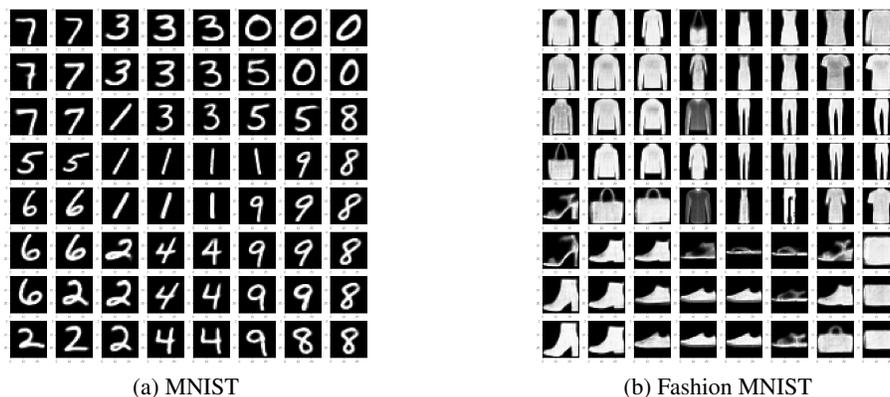


Figure 4: Reconstructions of MNIST / Fashion MNIST data from SOM cells in the  $8 \times 8$  grid learned by VarPSOM, illustrating the topological neighbourhood structure induced by our method, which aids interpretability.

212 **Interpretability** To illustrate the topological structure in the latent space, we present reconstructions of the VarPSOM centroids, arranged in a  $(8 \times 8)$  grid, on static MNIST/Fashion-MNIST data in Figure 4. On the real-world ICU time series data, we show example trajectories for one patient dying at the end of the ICU stay, as well as two control patients which are dispatched healthily from the ICU. We observe that the trajectories are located in different parts of the SOM grid, and form a smooth and interpretable representation (Fig. 5). For further results, including a more quantitative evaluation using randomly sampled trajectories as well as an illustration of how the uncertainty generated by the soft assignments can help in data visualization, we refer to the appendix (B.4).

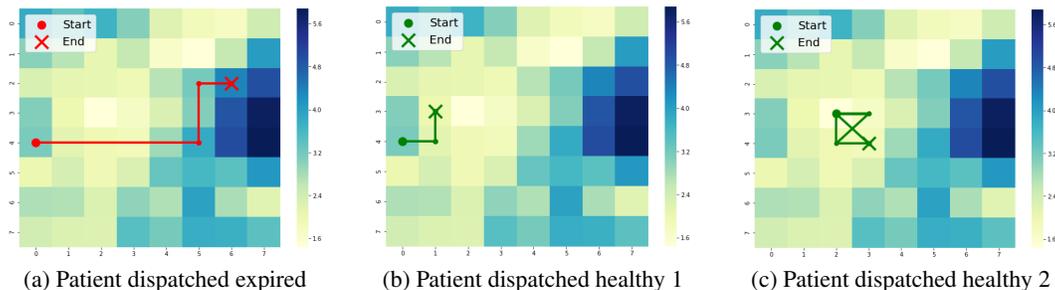


Figure 5: Illustration of 3 example patient trajectories between ICU admission and ICU dispatch, in the 2D SOM grid of VarTPSOM. The heatmap shows the enrichment of cells for the current APACHE physiology score. We observe qualitative differences in the trajectories of dying/healthy patients.

## 220 5 CONCLUSION

221 We presented two novel methods for interpretable unsupervised clustering, VarPSOM and VarTP-  
 222 SOM. Both models make use of a VAE and a novel clustering method, PSOM, that extends the  
 223 classical SOM algorithm to include a centroid-based probability distribution. Our models achieve  
 224 superior clustering performance compared to state-of-the-art deep clustering baselines on bench-  
 225 mark data sets and real-world medical time series. The use of a VAE for feature extraction, instead  
 226 of an AE, used in previous methods, and the use of soft assignments of data points to clusters results  
 227 in an interpretable model that can quantify uncertainty in the data.

## 228 REFERENCES

- 229 Elie Aljalbout, Vladimir Golkov, Yawar Siddiqui, and Daniel Cremers. Clustering with deep learn-  
230 ing: Taxonomy and new methods. *CoRR*, abs/1801.07648, 2018. URL [http://arxiv.org/](http://arxiv.org/abs/1801.07648)  
231 [abs/1801.07648](http://arxiv.org/abs/1801.07648).
- 232 Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statis-*  
233 *tics)*. Springer-Verlag, Berlin, Heidelberg, 2006. ISBN 0387310738.
- 234 Geoffrey J. Chappell and John G. Taylor. The temporal kohonen map. *Neural Netw.*, 6(3):441–445,  
235 March 1993. ISSN 0893-6080. doi: 10.1016/0893-6080(93)90011-K. URL [http://dx.doi.](http://dx.doi.org/10.1016/0893-6080(93)90011-K)  
236 [org/10.1016/0893-6080\(93\)90011-K](http://dx.doi.org/10.1016/0893-6080(93)90011-K).
- 237 Arthur Flexer. On the use of self-organizing maps for clustering and visualization. In Jan M. Żytkow  
238 and Jan Rauch (eds.), *Principles of Data Mining and Knowledge Discovery*, pp. 80–88, Berlin,  
239 Heidelberg, 1999. Springer Berlin Heidelberg. ISBN 978-3-540-48247-5.
- 240 Florent Forest, Mustapha Lebbah, Hanene Azzag, and Jérôme Lacaille. Deep embedded som: Joint  
241 representation learning and self-organization. 04 2019.
- 242 Vincent Fortuin and Gunnar Rätsch. Deep mean functions for meta-learning in gaussian processes.  
243 *arXiv preprint arXiv:1901.08098*, 2019.
- 244 Vincent Fortuin, Matthias Hüser, Francesco Locatello, Heiko Strathmann, and Gunnar Rätsch.  
245 Som-vae: Interpretable discrete representation learning on time series. *arXiv preprint*  
246 *arXiv:1806.02199*, 2018.
- 247 Vincent Fortuin, Gunnar Rätsch, and Stephan Mandt. Multivariate time series imputation with  
248 variational autoencoders. *arXiv preprint arXiv:1907.04155*, 2019.
- 249 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
250 Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. *arXiv e-prints*, art.  
251 [arXiv:1406.2661](https://arxiv.org/abs/1406.2661), Jun 2014.
- 252 Xifeng Guo, Long Gao, Xinwang Liu, and Jianping Yin. Improved deep embedded clustering with  
253 local structure preservation. In *Proceedings of the Twenty-Sixth International Joint Conference*  
254 *on Artificial Intelligence, IJCAI-17*, pp. 1753–1759, 2017. doi: 10.24963/ijcai.2017/243. URL  
255 <https://doi.org/10.24963/ijcai.2017/243>.
- 256 D. A. Keim. Information visualization and visual data mining. *IEEE Transactions on Visualization*  
257 *and Computer Graphics*, 8(1):1–8, Jan 2002. ISSN 1077-2626. doi: 10.1109/2945.981847.
- 258 Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv e-prints*, art.  
259 [arXiv:1312.6114](https://arxiv.org/abs/1312.6114), Dec 2013.
- 260 T. Kohonen. The self-organizing map. *Proceedings of the IEEE*, 78(9):1464–1480, Sep. 1990. ISSN  
261 0018-9219. doi: 10.1109/5.58325.
- 262 Teuvo Kohonen. The adaptive-subspace som (assom) and its use for the implementation of invariant  
263 feature detection. 1995.
- 264 Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document  
265 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998. ISSN 0018-9219. doi:  
266 10.1109/5.726791.
- 267 Xiaopeng Li, Zhoung Chen, and Nevin L. Zhang. Latent tree variational autoencoder for joint  
268 representation learning and multidimensional clustering. *CoRR*, abs/1803.05206, 2018. URL  
269 <http://arxiv.org/abs/1803.05206>.
- 270 N. Liu, J. Wang, and Y. Gong. Deep self-organizing map for visual classification. In *2015 Interna-*  
271 *tional Joint Conference on Neural Networks (IJCNN)*, pp. 1–6, July 2015. doi: 10.1109/IJCNN.  
272 2015.7280357.

- 273 J. MacQueen. Some methods for classification and analysis of multivariate observations. In  
274 *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol-*  
275 *ume 1: Statistics*, pp. 281–297, Berkeley, Calif., 1967. University of California Press. URL  
276 <https://projecteuclid.org/euclid.bsmmsp/1200512992>.
- 277 T. A. McQueen, A. A. Hopgood, J. A. Tepper, and T. J. Allen. A recurrent self-organizing map for  
278 temporal sequence processing. In Ahamad Lotfi and Jonathan M. Garibaldi (eds.), *Applications*  
279 *and Science in Soft Computing*, pp. 3–8, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.  
280 ISBN 978-3-540-45240-9.
- 281 Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi.  
282 The eicu collaborative research database, a freely available multi-center database for critical care  
283 research. *Scientific data*, 5, 2018.
- 284 S. Tirunagari, N. Poh, K. Aliabadi, D. Windridge, and D. Cooke. Patient level analytics using  
285 self-organising maps: A case study on type-1 diabetes self-care survey responses. In *2014 IEEE*  
286 *Symposium on Computational Intelligence and Data Mining (CIDM)*, pp. 304–309, Dec 2014.  
287 doi: 10.1109/CIDM.2014.7008682.
- 288 Kazuhiro Tokunaga and Tetsuo Furukawa. Modular network som. *Neural Networks*, 22(1):82 –  
289 90, 2009. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2008.10.006>. URL <http://www.sciencedirect.com/science/article/pii/S0893608008002335>.
- 291 Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learn-  
292 ing. *CoRR*, abs/1711.00937, 2017. URL <http://arxiv.org/abs/1711.00937>.
- 293 Thomas Voegtlin. Recursive self-organizing maps. *Neural Networks*, 15(8):979 – 991, 2002.  
294 ISSN 0893-6080. doi: [https://doi.org/10.1016/S0893-6080\(02\)00072-2](https://doi.org/10.1016/S0893-6080(02)00072-2). URL <http://www.sciencedirect.com/science/article/pii/S0893608002000722>.
- 296 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-  
297 ing machine learning algorithms. *CoRR*, abs/1708.07747, 2017. URL <http://arxiv.org/abs/1708.07747>.
- 299 Junyuan Xie, Ross B. Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering  
300 analysis. *CoRR*, abs/1511.06335, 2015. URL <http://arxiv.org/abs/1511.06335>.
- 301 Bo Yang, Xiao Fu, Nicholas D. Sidiropoulos, and Mingyi Hong. Towards k-means-friendly  
302 spaces: Simultaneous deep learning and clustering. *CoRR*, abs/1610.04794, 2016. URL  
303 <http://arxiv.org/abs/1610.04794>.

## 304 APPENDIX

## 305 A SELF-ORGANIZING MAPS

306 Among various existing interpretable unsupervised learning algorithms, Kohonen’s self-organizing  
 307 map (SOM) (Kohonen, 1990) is one of the most popular models. It is comprised of  $K$  neurons  
 308 connected to form a discrete topological structure. The data are projected onto this topographic map  
 309 which locally approximates the data manifold. Usually it is a finite two-dimensional region where  
 310 neurons are arranged in a regular hexagonal or rectangular grid. Here we use a grid,  $M \in \mathbb{N}^2$ ,  
 311 because of its simplicity and its visualization properties. Each neuron  $m_{ij}$ , at position  $(i, j)$  of the  
 312 grid, for  $i, j = 1, \dots, \sqrt{K}$ , corresponds to a centroid vector,  $\mu_{i,j}$  in the input space. The centroids  
 313 are tied by a neighborhood relation, here defined as  $N(\mu_{i,j}) = \{\mu_{i-1,j}, \mu_{i+1,j}, \mu_{i,j-1}, \mu_{i,j+1}\}$ .

314 Given a random initialization of the centroids, the SOM algorithm randomly selects an input  $x_i$  and  
 315 updates both its closest centroid  $\mu_{i,j}$  and its neighbors  $N(\mu_{i,j})$  to move them closer to  $x_i$ . The  
 316 algorithm (1) then iterates these steps until convergence.

---

**Algorithm 1** Self-Organizing Maps

---

**Require:**  $0 < \alpha(t) < 1$ ;  $\lim_{t \rightarrow \infty} \sum \alpha(t) \rightarrow \infty$ ;  $\lim_{t \rightarrow \infty} \sum \alpha^2(t) < \infty$ ;

**repeat**

At each time  $t$ , present an input  $x(t)$  and select the winner,

$$\nu(t) = \arg \min_{k \in \Omega} \|\mathbf{x}(t) - \mathbf{w}_k(t)\|$$

Update the weights of the winner and its neighbours,

$$\Delta \mathbf{w}_k(t) = \alpha(t) \eta(\nu, k, t) [\mathbf{x}(t) - \mathbf{w}_\nu(t)]$$

**until** the map converges

---

317 The range of SOM applications includes high dimensional data visualizations, clustering, image  
 318 and video processing, density or spectrum profile modeling, text/document mining, management  
 319 systems and gene expression data analysis.

## 320 B EXPERIMENTAL AND IMPLEMENTATION DETAILS

## 321 B.1 DATASETS

- 322 • **MNIST:** It consists of 70000 handwritten digits of 28-by-28 pixel size. Digits range from  
 323 0 to 9, yielding 10 patterns in total. The digits have been size-normalized and centered in a  
 324 fixed-size image Lecun et al. (1998).
- 325 • **Fashion MNIST:** A dataset of Zalando’s article images consisting of a training set of  
 326 60,000 examples and a test set of 10,000 examples Xiao et al. (2017). Each example is  
 327 a  $28 \times 28$  grayscale image, associated with a label from 10 classes.
- 328 • **eICU:** For temporal data we use vital sign/lab measurements of intensive care unit (ICU)  
 329 patients resampled to a 1-hour based grid. The last 72 hours of these time series were  
 330 used for the experiments. As labels we use a variant of the current dynamic APACHE  
 331 physiology score (APACHE-0) as well as the worst APACHE score in the next 6 and 12  
 332 hours (APACHE-6/12).

333 Each dataset is divided into training, validation and test sets for both our models and the baselines.

## 334 B.2 LATENT SPACE DIMENSION

335 We evaluated the DEC model for different latent space dimensions. Table S1 shows that the AE,  
 336 used in the DEC model, performs better when a lower dimensional latent space is used.

Table S1: Mean/Standard error of NMI and Purity of DEC model on MNIST test set, across 10 runs with different random model initializations. We use 64 clusters and different latent space dimensions.

Latent dimension	Purity	NMI
$l = 10$	$0.950 \pm 0.001$	$0.681 \pm 0.001$
$l = 100$	$0.750 \pm 0.001$	$0.573 \pm 0.001$

### 337 B.3 LEARNING HEALTH STATE REPRESENTATIONS IN THE ICU

338 By enforcing a SOM structure, VarPSOM, as well as VarTPSOM, project the cluster centroids onto  
 339 a discrete 2D grid. Such a grid is particularly suited for visualization purposes and relations between  
 340 centroids become immediately intuitive. In Fig. S1 a heat-map (colored according to enrichment  
 341 in the current APACHE score, as well as future mortality risk in the next 24 hours) shows compact  
 342 enrichment structures. Clusters with similar enrichment for mortality risk and current APACHE  
 343 score, respectively, are often close to each other on the SOM grid. Our model thus succeeds in creat-  
 344 ing a meaningful neighbourhood structure over the centroids with respect to these clinical quantities,  
 345 even though it is learned purely unsupervised. The two heat-maps (S1a and S1b) show different enrich-  
 346 ment patterns. Clusters which are enriched in patients with higher APACHE scores often do not  
 347 correspond exactly to clusters with a higher mortality risk. This suggests that traditional represen-  
 348 tations of physiologic values, such as the APACHE score, fail to fully use all complex multivariate  
 349 relationships present in the ICU recordings.

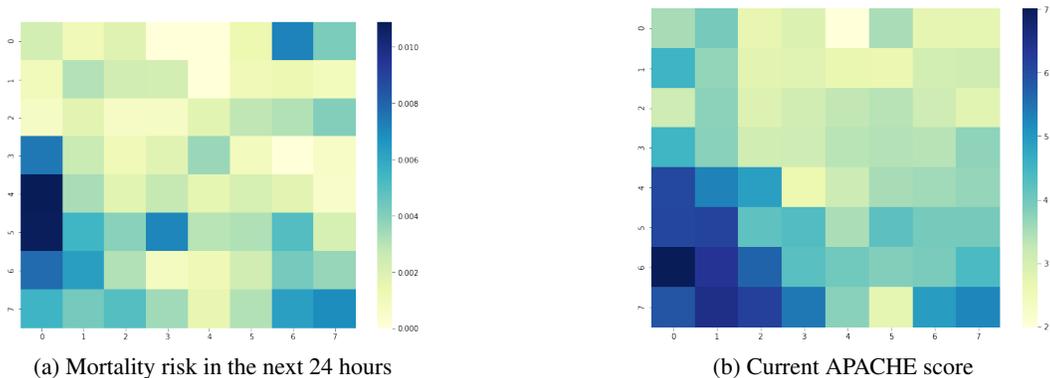


Figure S1: Heat-maps on enrichment in mortality risk in the next 24 hours as well as the current dynamic APACHE score, superimposed on the discrete 2D grid learned by VarTPSOM.

### 350 B.4 VISUALIZING HEALTH STATE TRAJECTORIES IN THE ICU

351 To analyze the trend of the patient pathology, VarTPSOM induces trajectories on the 2D SOM grid  
 352 which can be easily visualized. Fig. S2 shows 20 randomly sampled patient trajectories obtained  
 353 by our model. Trajectories ending in the death of the patient are shown in red, healthily dispatched  
 354 patients are shown in green.

355 One of the main advantage of VarTPSOM over the traditional SOM algorithm is the use of soft  
 356 assignments of data points to clusters which results in a better ability to quantify uncertainty in the  
 357 data. For visualizing health states in the ICU, this property is very important. In Fig S3 we plot an  
 358 example patient trajectory, where 6 different time-steps (in temporal order) of the trajectory were  
 359 chosen. Our model yields a soft centroid-based probability distribution which evolves with time and  
 360 which allows estimation of likely discrete health states at a given point in time. For each time-step  
 361 the distribution of probabilities is plotted using a heat-map, whereas the overall trajectory is plotted  
 362 using a black line. The circle and cross indicate ICU admission and dispatch, respectively.

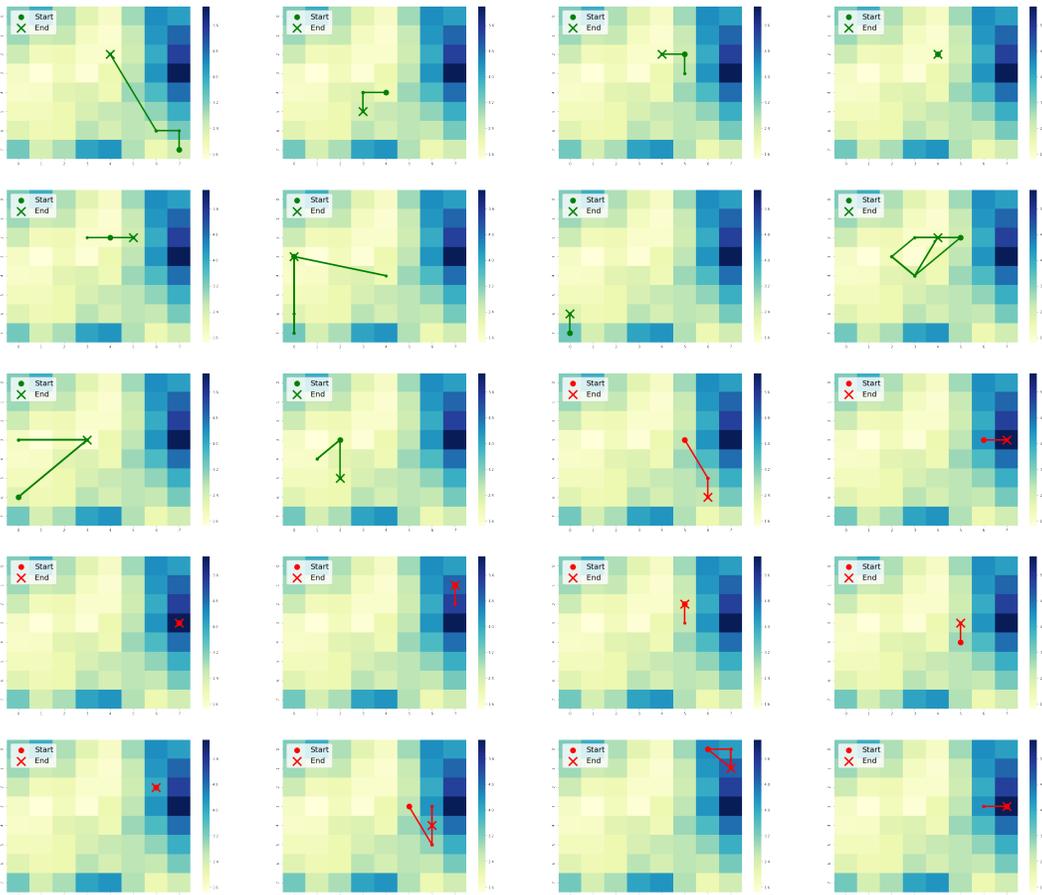


Figure S2: Randomly sampled VarTPSOM trajectories, from patients expired at the end of the ICU stay, as well as healthily dispatched patients. Superimposed is a heatmap which displays the cluster enrichment in the current APACHE score. We observe that trajectories of dying patients are often in different locations of the map as healthy patients, in particular in those regions enriched for high APACHE scores, which corresponds with clinical intuition.

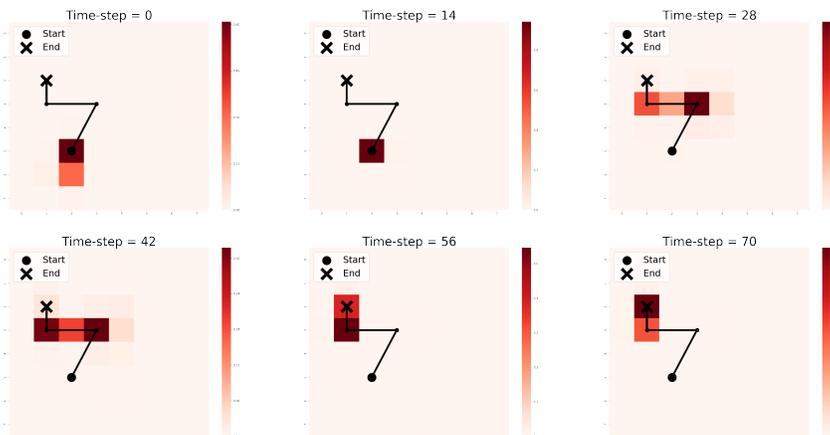


Figure S3: Probabilities over discrete patient health states for 6 different time-steps of the selected time series.