

Identifying Decision Points for Safe and Interpretable Reinforcement Learning in Hypotension Treatment

Kristine Zhang*
Yuanheng Wang*
Jianzhun Du*
Brian Chu
Harvard University

KRISTINE_ZHANG@COLLEGE.HARVARD.EDU
 YUANHENG_WANG@G.HARVARD.EDU
 JZDU@G.HARVARD.EDU
 BRIANCHU@G.HARVARD.EDU

Leo Anthony Celi
Massachusetts Institute of Technology

LEOANTHONYCELI@YAHOO.COM

Ryan Kindle
Massachusetts General Hospital

KINDLR@GMAIL.COM

Finale Doshi-Velez
Harvard University

FINALE@SEAS.HARVARD.EDU

Abstract

Many batch RL health applications first discretize time into fixed intervals. However, this discretization both loses resolution and forces a policy computation at each (potentially fine) interval. In this work, we develop a novel framework to compress continuous trajectories into a few, interpretable decision points—places where the batch data support multiple alternatives. We apply our approach to create recommendations from a cohort of hypotensive patients dataset. Our reduced state space results in faster planning and allows easy inspection by a clinical expert.

Keywords: batch RL

and robust off-policy evaluation before policies gleaned from batch data are applied to patients ([Gottesman et al., 2019](#)).

In this paper, we focus on the issue of time discretization in this batch RL context. To date, most batch RL applications first discretized time into bins and then solved for a policy decision for each bin. Not only are learnt policies often sensitive to the binning method, but decisions must be solved for every time step, affecting both computation and policy evaluation ([Tallec et al., 2019](#)).

We ask and answer the question: what if we didn't try to solve at every time point? Indeed, especially with batch data, the *only* times when one might be able to give a recommendation are in states where clinicians treat similar patients differently. We call these *decision points* (DP). Compressing trajectories down to high-ambiguity *decision regions*, or areas of the state space with many DPs, has advantages not only for efficient planning but also for interpretability. Summarizing a continuous time and space problem in 20-40 discrete clusters enables clinicians to manually inspect the recommended

1. Introduction

Many works have investigated reinforcement learning for sequential decision-making in health ([Yu et al., 2019](#)). A standard approach is to take already-collected (batch) data and use that to train a policy. However, batch RL has many subtleties, and the high stakes of healthcare settings necessitate special attention to safety, interpretability,

action for each cluster and thus globally validate a policy. Our contributions are: (1) developing a general pipeline for learning this decision region MDP and (2) demonstrating its interpretability on a hypotension management task.

2. Related Work

Many works have considered RL for health, including for hypotension management in the ICU (Komorowski et al., 2018; Futoma et al., 2020; Srinivasan and Doshi-Velez, 2020). Unlike our work, all of these approaches choose some time discretization with a constant time step size and optimize decisions at each of these points. Gottesman et al. (2020) propose an interpretable RL approach to handling treatment changes or patient measurements with irregular observations, which is similar to our methods, though we leverage the uncertainty of the behavior policy for both state and time compression. Du et al. (2020) take a model-based approach to continuous-time RL with demonstrations in the healthcare environment, whereas it cannot be applied to the batch RL setting.

More broadly, state abstraction (e.g. Li et al. (2006)) is a popular RL problem. However, these works have been largely in online settings and focused on aggregation based on policy-relevant information like Q -values or transitions; they also only compress the state space and the trajectory lengths remain the same. In the context of batch RL, we instead aim to identify important states for decision-making and filter out regions where the experts have consensus.

In this way, our methods enable efficient planning and policy validation. Komorowski et al. (2018) also employ clustering to discretize a continuous environment; however, their state clustering methods do not consider practice variation and thus produce large state spaces (750 vs. 32 for us). The

issue of time discretization also remains. Other works use deep models for state representations and policies, but these are difficult to interpret, while ours are easy to inspect.

3. Background and Notation

A Markov Decision Process (MDP) is a tuple $\langle \mathcal{X}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma \rangle$. We assume a continuous state space $\mathcal{X} \in \mathbb{R}^d$, discrete action space \mathcal{A} , and discount factor $\gamma \in (0, 1]$. $T(x'|x, a)$ and $R(x, a)$ denote the state transition function and reward function, respectively. A policy $\pi : (\mathcal{X}, \mathcal{A}) \rightarrow [0, 1]$ defines the probability of an action given a state. The objective is to learn a policy that maximizes the expected return, or sum of discounted rewards $V^\pi = E [\sum_t \gamma^t r_t | a_t \sim \pi]$.

In batch reinforcement learning, we cannot interact with the environment during learning, but instead begin with a pre-collected dataset of trajectories $\mathcal{D} = \{(x_i, a_i, x'_i, r_i)\}_{i=1}^N$. These trajectories are sampled by following a (possibly unknown) behavior policy π_b and used to learn a new evaluation policy π_e . MDPs with small discrete state and action spaces can be solved using value or policy iteration (Sutton and Barto, 2018).

4. Method

The innovation of our method is that we do not consider every state to be a potential decision point. Rather, we only consider those points with high behavior policy variability, that is, states whose similar neighbors are administered different treatments. While this idea is simple, its realization requires formalizing many subtleties, which we describe below.

Decision Point Identification We operationalize the notion of states with high clinical variation by first training a kernel-based classifier to *predict* clinician actions.

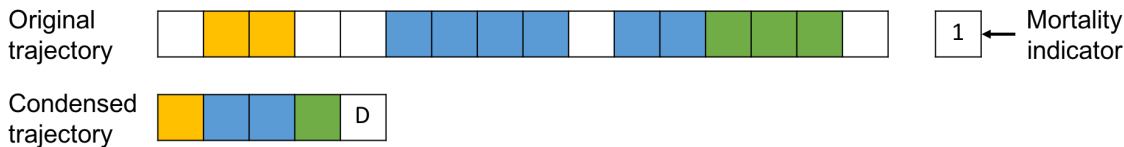


Figure 1: Compression process from a trajectory in the original state space to decision clusters. White squares are non-DP states, colors denote DPs in the same cluster. The mortality outcome is captured via two more states: A (alive) or D (dead).

States with sufficient neighbors (according to the kernel) whose actions still cannot be predicted well are then considered *decision points*.

Specifically, we consider a weighted Gaussian kernel:

$$k(x, x') = \exp(-\|(x - x')\|_2^2) \quad (1)$$

where $k(x, x')$ is the estimated similarity between states, and \mathbf{w} can be interpreted as an importance weighting over state dimensions and \odot represents an element-wise multiplication between vectors. We learn the parameters \mathbf{w} by backpropagating through the kernel space with optimizing the cross-entropy loss of predicting behavior actions. Details can be found in Appendix B.1.

The second step is to determine a set of neighbors $\{x'\}$ for each state x in the dataset. We introduce two hyperparameters: the minimum kernel similarity threshold $\delta \in (0, 1]$ to be considered a neighbor, and the minimum number of neighbors n for an action to be considered valid. That is, if a state x has at least n neighbors within δ performing action a , then action a is considered a valid action for x . If a state has more than one valid action, it is considered a decision point.

Clustering Decision Points The process above identifies whether each state in the dataset is a decision point. Next, we cluster only the decision points to identify discrete *decision regions* that will form our reduced MDP. A top-down hierarchical clustering approach is applied, where the Distances

between decision points are calculated using the standardized state features (dimensions of x). For each intermediate cluster, we examine the difference of mean values for each feature across actions. If the difference of these means exceed specified thresholds, this implies the region is not homogeneous and we further split clusters (we truly want decision regions where the optimal action is unclear). Secondly, we examine if the current clustering is creating loops—defined as trajectories that leave a cluster and coming back to the same cluster within the next three time steps—as these (artificial) loops can make the agent believe that it can “freeze time” and avoid any final consequence. If loop percentage is high, we split that cluster.

Getting a Policy Finally, we are ready to create the MDP and solve for a policy. We convert each patient trajectory of the form $\{(x, a, x')\}$ into a new, compressed trajectory $\{(\bar{x}, \bar{a}, \bar{x}')\}$. We use these transitions to learn a discrete, maximum-likelihood MDP in which only valid actions are allowed, and then solve.

The compression process from $\{(x, a, x')\}$ to $\{(\bar{x}, \bar{a}, \bar{x}')\}$ is summarized by Figure 1 and formalized in Appendix B.4. At a high level, we append a new cluster \bar{x}_t each time a transition is observed from a *different cluster or non-DP state*. Condensing contiguous DPs in the same cluster and ignoring non-DPs is based on the idea that a “decision” is an action taken in one decision region that causes movement to another decision region. In the

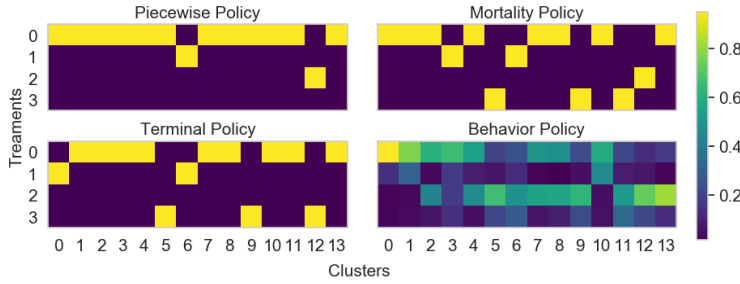


Figure 2: Three policies obtained using three different rewards functions, in comparison to the behavior policy.

Policy	WIS Score	ESS
Behavior	-1.11	3,884
Piecewise	-0.93	1,945
Mortality	-0.94	1,895
Terminal	-0.94	1,921

Table 1: OPE results from different policies. “ESS” means effective sample size.

final tuple of a trajectory, we include a transition to the mortality state of the patient to help track their outcome.

5. Case Study for Hypotension Management

Setting Now, we apply our method to recommend treatments to manage hypotension in the ICU. We select a cohort of 11,739 patients records from the MIMIC-III dataset (See Appendix A for details). Our goal is to (1) recommend a set of treatment policies that are reasonable, and (2) offer a toolkit with visualizations for clinicians to validate and refer to when making a treatment decision for hypotensive patients. Since this is a proof-of-concept, we only consider four action classes: do nothing, use fluids, use vasopressors, and use both.

To check for sensitivity to reward formulation, we solved with three reward functions: (1) Piecewise: based on the MAP and urine volume for each patient at each timestamp. Lower map and urine values correspond to lower rewards, (2) Mortality probability: inverse of the mortality probability of each cluster-treatment pair in the constructed MDP, and (3) Terminal state: Based on whether patients reach Alive ($R = 1$) or Death ($R = 0$) at the end of the ICU admission. We trained the policy using value iteration and use a discount rate of 0.98.

Results Running our algorithm on these data results in 33 clusters, where the largest cluster includes 117,730 decision points and the smallest cluster includes 13 decision points. Figure 4 in Appendix B.3 shows the mean value of each feature for clusters containing high mortality rate ($\geq 10\%$) and enough data (≥ 10 patients treatment points for at least two types of treatments); this figure was easily inspected by intensivists.

Figure 2 shows the optimal policy for each reward function. We also include the behavior policy for comparison, which is the empirical proportion of each action taken in the historical treatment data. The fact that each policy from different rewards functions is slightly different but largely similar to behavior policy demonstrates the robustness of our methods, and the fact it can be so succinctly summarized allowed us to discuss each case with clinicians; in each case, they could posit whether the policy recommendation made sense or could have been due to some confound. Such discussion would not be possible with a black-box policy.

To quantitatively evaluate the evaluation policy learned using the aforementioned rewards, we use *weighted importance sampling* (WIS) for off policy evaluation (OPE), whose estimate is biased but consistent (asymptotically correct) with much lower variance. We apply standard weight clipping by capping weights at 95 percentile. Table 1 reports the

values of WIS estimates. It shows all three evaluation policies achieve similar rewards, which is higher than the reward of the behavior policy.

6. Conclusion

In this work, we present a pipeline for identifying the decision regions in an RL environment and employing them to construct a summary MDP. This summary MDP can then be used for efficient planning and policy validation by human experts. We demonstrate the value of this new approach to interpretability in healthcare settings using the example of hypotension data. Our proposed framework is highly generalizable and can be fine-tuned for other applications in healthcare and beyond.

References

Jianzhun Du, Joseph Futoma, and Finale Doshi-Velez. Model-based reinforcement learning for semi-markov decision processes with neural odes. *arXiv preprint arXiv:2006.16210*, 2020.

Joseph Futoma, Michael C Hughes, and Finale Doshi-Velez. Popcorn: Partially observed prediction constrained reinforcement learning. *arXiv preprint arXiv:2001.04032*, 2020.

Omer Gottesman, Fredrik Johansson, Matthieu Komorowski, Aldo Faisal, David Sontag, Finale Doshi-Velez, and Leo Anthony Celi. Guidelines for reinforcement learning in healthcare. *Nat Med*, 25(1): 16–18, 2019.

Omer Gottesman, Joseph Futoma, Yao Liu, Sonali Parbhoo, Leo Anthony Celi, Emma Brunskill, and Finale Doshi-Velez. Interpretable off-policy evaluation in reinforcement learning by highlighting influential transitions, 2020.

Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3:160035, 2016.

Matthieu Komorowski, Leo A. Celi, Omar Badawi, Anthony C. Gordon, and A. Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature Medicine*, 24(11):1716–1720, 2018. doi: 10.1038/s41591-018-0213-5. URL <https://doi.org/10.1038/s41591-018-0213-5>.

Lihong Li, Thomas J. Walsh, and Michael L. Littman. Towards a unified theory of state abstraction for mdps. In *In Proceedings of the Ninth International Symposium on Artificial Intelligence and Mathematics*, pages 531–539, 2006.

Tu Dinh Nguyen, Trung Le, Hung Bui, and Dinh Phung. Large-scale online kernel learning with random feature reparameterization. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI’17*, page 2543–2549. AAAI Press, 2017. ISBN 9780999241103.

Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.

Srivatsan Srinivasan and Finale Doshi-Velez. Interpretable batch irl to extract clinician goals in icu hypotension management. *AMIA Jt Summits Transl Sci Proc*, pages 636—645, 2020.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018. ISBN 0262039249.

Corentin Tallec, Léonard Blier, and Yann Olivier. Making deep q-learning methods robust to time discretization, 2019.

Chao Yu, Jiming Liu, and Shamim Nemati. Reinforcement learning in healthcare: A survey. *arXiv preprint arXiv:1908.08796*, 2019.

80% of the states receiving no treatment. To analyze treatment impact on outcome, we define a patient’s mortality to be true if they died in the hospital within 30 days of ICU admission.

Appendix A. Data Processing

A.1. MIMIC Data Cleaning

Our hypotension data is obtained from the publicly available MIMIC-III dataset for ICU patients (Johnson et al., 2016), with records from 15,653 unique ICU stays. Patients are considered hypotensive if their systolic blood pressure (SBP) drops below 90 mmHg.

Preprocessing was performed and is described in greater detail by Futoma et al. (2020). Trajectories were initially discretized into hourly bins starting at $t = 1$ hour into ICU admission and were truncated either at discharge or at $t = 72$ hours (as patients staying beyond 72 hours tend to be abnormal cases).

The initial state space consists of 111 clinical features including patient measurements, vital signs, and records of past treatment actions. These dimensions contain some redundancy with indicators and quantitative measure for the same variables.

For actions, we track the amount of either fluids or vasopressors administered to each patient at each time step. For this pipeline, we primarily consider the four general actions corresponding to “no treatment”, “fluids only”, “vasopressors only”, or “fluids and vasopressors both given”. The learned policy is intended to help physicians decide whether to give each type of treatment at all, as there are existing guidelines for the quantity to give once a choice has been made. Actions are unbalanced in this dataset, with around

A.2. Dataset Split

Because we lack access to a simulator for MIMIC, the efficacy of our methods must be tested on the batch data. We split the trajectories in the original dataset 75-25 into train and test sets, yielding 555,196 transition tuples from 11,739 trajectories in the train set and 185,671 from 3,914 trajectories in the test set.

The train set is used for initial data analysis, MDP compression, and planning. The test set is reserved for OPE.

A.3. State Feature Selection

Due to the high dimensionality of the starting MIMIC state space, we preemptively select for features that are most relevant to action prediction. The objective was to increase the efficiency of kernel-learning and avoid collinearity between redundant features. We construct a random forest predictor from states to actions, implemented using a random forest classifier with 100 estimators, depth 3, and balanced class weights. We employ only the top 20 features ranked by random forest feature importance in the kernel learning algorithm. Figure 4 lists all continuous features within this top 20 and the features clinicians consider to be important.

Appendix B. Implementation Details

B.1. Kernel Learning

We are interested in evaluating the Gaussian kernel distance between states, defined as

$$k(x, x') = \exp(-\|(x - x')\|_2^2) \quad (2)$$

Naively computing the pairwise distance between N states of dimension d is computationally expensive, given $\mathcal{O}(N^2d)$ complexity and the large MIMIC-III dataset. [Rahimi and Recht \(2008\)](#) showed that this process can be accelerated through the use of *Random Fourier Features* (RFF). After applying a randomized feature map $z : \mathbb{R}^d \rightarrow \mathbb{R}^D$ from state vectors to a low-dimensional Euclidean inner product space, kernel evaluation can be estimated by the inner product of the randomized features:

$$k(x, x') \approx z(x)^\top z(x') \quad (3)$$

We can therefore approximately train kernel machines by projecting to the Fourier space and applying fast linear methods such as regression. We evaluate the efficacy of a set of transformations \mathbf{w} by performing multi-class logistic regression to predict actions. We define $z(x_i) = \text{RFF}(\mathbf{w} \odot x_i)$ to be the covariates such that $k(x, x') \approx z(x)^\top z(x')$. The targets are a_i with regression coefficients \mathbf{V} . The objective function is the corresponding cross-entropy loss in Equation 4.

$$\min_{\mathbf{w}, \mathbf{V}} - \sum_{i=1}^N \sum_{a \in \mathcal{A}} \mathbb{I}(a_i = a) \hat{p}(a|x_i) \quad (4)$$

$$\text{with } \hat{p}(a|x) = \text{softmax}(z(x)^\top \mathbf{V})_a$$

We propose Algorithm 1 that uses gradient descent to simultaneously optimize \mathbf{w} and \mathbf{V} . [Nguyen et al. \(2017\)](#) showed that the gradient of the cross-entropy loss can be calculated over \mathbf{w} despite the random feature mapping. We also apply a reparameterization trick, optimizing over $\mathbf{u} = \log \mathbf{w}$ to ensure that \mathbf{w} is positive.

Algorithm 1: Learning Weighted Gaussian Kernel

Input : Set of state-action tuples $\{x_i, a_i\}_{i=1}^N$; learning rate η , number of epochs M , mini-batch size B .

Output: Learned kernel weights \mathbf{w} . Initialize $\mathbf{u} \leftarrow \mathbf{0}^d, \mathbf{V} \leftarrow \mathbf{1}^{D \times |\mathcal{A}|}$; Generate randomized parameters $\omega \in \mathbb{R}^d, b \in \mathbb{R}$ according to the RFF algorithm;

for $l \leftarrow 1$ **to** M **do**

for $j \leftarrow 1$ **to** $\lfloor N/B \rfloor$ **do**

Select mini-batch of states X_j and actions A_j ;

$\mathbf{w} \leftarrow \exp(\mathbf{u})$

Project $z(x) \leftarrow \text{RFF}((\mathbf{w} \odot x; \omega, b)$ for all $x \in X_j$;

$L \leftarrow$ cross-entropy loss when predicting A_j with features Z_j and weights \mathbf{V} ;

$\mathbf{V} \leftarrow \mathbf{V} - \eta \frac{\partial L}{\partial \mathbf{V}}$;

$\mathbf{u} \leftarrow \mathbf{u} - \eta \frac{\partial L}{\partial \mathbf{u}}$;

end

end

Return $\mathbf{w} = \exp(\mathbf{u})$;

B.2. DP Hyperparameter Choice

To identify decision points, we have two main hyperparameters: δ , the minimum similarity from a point x to neighbors, and n the number of neighbors with the same action (Figure 3). Corresponding actions for points present in this ϵ -ball are considered to be allowable actions, and decision points are defined as points that exist in this ϵ -ball where neighbors do not dominantly take a sole action. Since we are more interested to measure the similarity between every point in the kernel space, we instead consider points within the similarity threshold δ . Also, we require a decision point must have a minimum number of neighbors n in this δ -

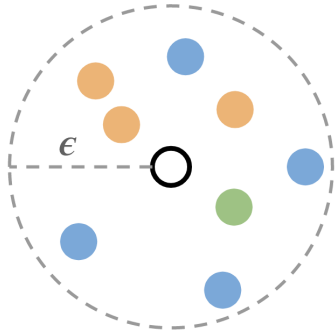


Figure 3: ϵ conceptually represents distance between points (which inversely approximate similarity δ). Higher similarity threshold δ translates to a smaller ϵ -ball, which is equivalent to a smaller neighborhood.

specified region. We perform the grid search over a set of candidate values on a sampled validation dataset. For each pair of δ, n , we evaluate the overall AUC score for a sampled subset of points, whose labels are determined by the empirical distribution of labels present in the region. Conducting such experiment multiple rounds, we finally choose the similarity threshold $\delta = 0.95$ and $n = 20$ that achieves the highest AUC score.

B.3. Decision Region Clustering

In order to form decision regions efficiently, we utilize the hierarchical clustering implementations in `scipy`. We take the output from `scipy` linkage function call, which is a matrix specifying the order of points being merged at each step. We iterate through the matrix in a top-down manner, and determine if the children of each split meets the specified condition, i.e., the differences of the means of feature values for points in each child are smaller than the predetermined threshold. If not, we continue iterate through the matrix until conditions are satisfied or maximum number of splits are reached.

Algorithm 2: MDP Trajectory Compression

Input : Trajectory $\{(x_i, a_i)\}_{i=1}^L$ with corresponding cluster labels $\{c_i\}_{i=1}^L$ (with $c_i = 0$ for non-DPs); mortality state I_m , action summary function h .

Output: Lists of abstracted states and actions $\{\bar{x}_j\}, \{\bar{a}_j\}$.

Initialize $\bar{X} \leftarrow [], \bar{A} \leftarrow [], c_0 \leftarrow 0, A \leftarrow []$;

for $l \leftarrow 1$ **to** L **do**

if $c_l > 0$ **then**

$A.append(a_l)$;

 // entering new decision region

if $c_l \neq c_0$ **then**

$\bar{X}.append(c_l)$;

end

 // leaving previous decision region

if $l = L$ or $c_{l+1} \neq c_l$ **then**

$\bar{A}.append(h(A))$;

$A \leftarrow []$;

end

end

end

$\bar{X}.append(I_m)$;

Return \bar{X}, \bar{A} ;

B.4. MDP Compression Details

Algorithm 2 formalizes the MDP compression method. Recall that we convert each patient trajectory $\{(x, a, x')\}$ into a compressed trajectory $\{(\bar{x}, \bar{a}, \bar{x}')\}$, where all $\bar{x}, \bar{x}' \in \mathcal{C}$ are decision clusters and $\bar{a} \in \bar{\mathcal{A}}$ are summarized actions. The summary function h should compress multiple actions in \mathcal{A} that take place while the patient state remains in the same cluster.

Given a set of compressed trajectories, we estimate a behavior policy $\bar{\pi}_b$ and transition function $\bar{T}(\bar{x}'|\bar{x}, \bar{a})$ from the new dataset, as formalized in Equation 5 and 6. Following standard state abstraction protocol (Li et al.,

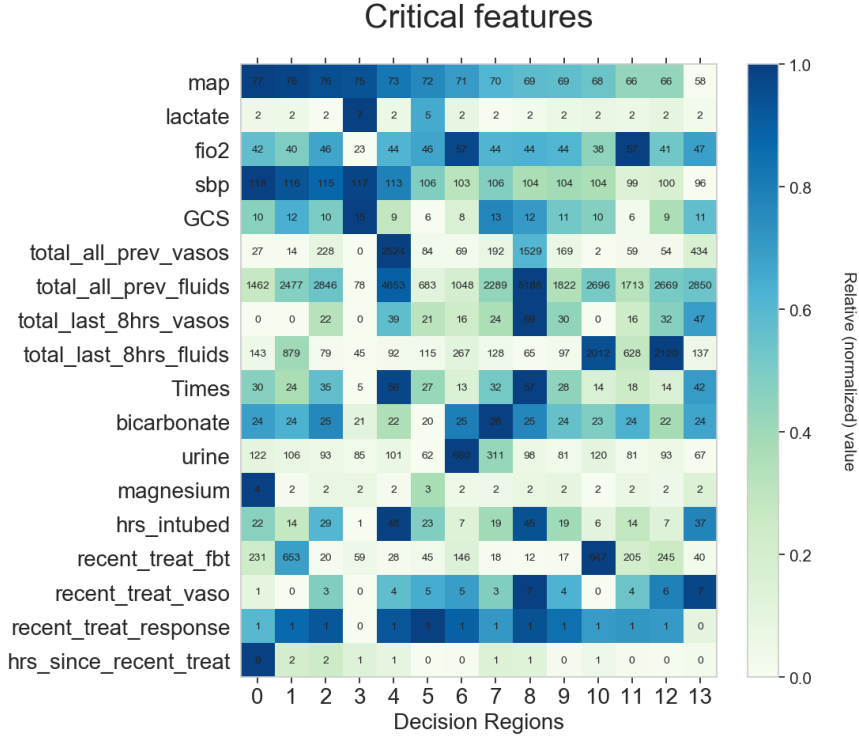


Figure 4: The feature means for points in each decision region. Sorted by descending order of Mean Arterial Pressure (MAP) value.

2006), the behavior policy is estimated as the proportion of times an action was taken from a given cluster, and the transition function is estimated as the proportion of times a given cluster-action pair led to another cluster.

$$\bar{\pi}_b(\bar{a}^*|\bar{x}^*) = \frac{\sum_i \mathbb{I}(\bar{x}_i = \bar{x}^*, \bar{a}_i = \bar{a}^*)}{\sum_i \mathbb{I}(\bar{x}_i = \bar{x}^*)} \quad (5)$$

$$\bar{T}(\bar{x}'^*|\bar{x}^*, \bar{a}^*) = \frac{\sum_i \mathbb{I}(\bar{x}_i = \bar{x}^*, \bar{a}_i = \bar{a}^*, \bar{x}'_i = \bar{x}'^*)}{\sum_i \mathbb{I}(\bar{x}_i = \bar{x}^*, \bar{a}_i = \bar{a}^*)} \quad (6)$$

This process can be thought of as simultaneous state and temporal abstraction. We perform state abstraction to map decision points into high-uncertainty clusters, then temporal abstraction to condense all actions that do not cause transitions out of decision regions. In this way, we reduce both the size of the state space and the length of trajectories.