
Medical Dead-ends and Learning to Identify High-risk States and Treatments

Mehdi Fatemi
Microsoft Research
mehdi.fatemi@microsoft.com

Taylor W. Killian
University of Toronto, Vector Institute
twkillian@cs.toronto.edu

Jayakumar Subramanian
Media and Data Science Research, Adobe India
jayakumar.subramanian@gmail.com

Marzyeh Ghassemi
Massachusetts Institute of Technology
mghassem@mit.edu

Abstract

Machine learning has successfully framed many sequential decision making problems as either supervised prediction, or optimal decision-making policy identification via reinforcement learning. In data-constrained offline settings, both approaches may fail as they assume fully optimal behavior or rely on exploring alternatives that may not exist. We introduce an inherently different approach that identifies possible “dead-ends” of a state space. We focus on the condition of patients in the intensive care unit, where a “medical dead-end” indicates that a patient will expire, regardless of all potential future treatment sequences. We postulate “treatment security” as avoiding treatments with probability proportional to their chance of leading to dead-ends, present a formal proof, and frame discovery as an RL problem. We then train three independent deep neural models for automated state construction, dead-end discovery and confirmation. Our empirical results discover that dead-ends exist in real clinical data among septic patients, and further reveal gaps between secure treatments and those that were administered.

1 Introduction

Off-policy Reinforcement Learning (RL) was designed as the way to isolate behavioural policies, which generate experience, from the target policy, which aims for optimality. It also enables learning multiple target policies with different goals from the same data-stream or from previously recorded experience [1]. This algorithmic approach is of particular importance in safety-critical domains such as robotics [2], education [3] or healthcare [4] where data collection should be regulated as it is expensive or carries significant risk. Despite significant advances made possible by off-policy RL combined with deep neural networks [5–7], the performance of these algorithms degrade drastically in fully *offline* settings [8], without additional interactions with the environment [9, 10]. These challenges are deeply amplified when the dataset is limited and exploratory new data cannot be collected for ethical or safety purposes. This is because robust identification of an optimal policy requires exhaustive trial and error of various courses of actions [11, 12]. In such fully offline cases, naively learned policies may significantly overfit to data-collection artifacts [13–15]. Estimation errors due to limited data may further lead to mistimed or inappropriate decisions with adverse safety consequences [16].

Even if optimality is not attainable in such constrained cases, negative outcomes in data can be used to identify behaviors to avoid, thereby guarding against overoptimistic decisions in safety-critical domains that may be significantly biased due to reduced data availability. In one such domain, healthcare, RL has been used to identify optimal treatment policies based on observed outcomes of

past treatments [17]. These policies correspond to advising *what treatments to administer*, given a patient’s condition. Unfortunately, exploration of potential courses of treatment is not possible in most clinical settings due to legal and ethical implications; hence, RL estimates of optimal policies are largely unreliable in healthcare [18].

In this paper, we develop a novel RL-based method, Dead-end Discovery (DeD), to identify *treatments to avoid* as opposed to what treatment to select. Our paradigm shift avoids pitfalls that may arise from constraining policies to remain close to possibly suboptimal recorded behavior as is typical in current state of the art offline RL approaches [10, 19–21]. When the data lacks sufficient amounts of exploratory behavior, these methods fail to attain a reliable policy. We instead use this data to constrain the scope of the policy, based on retrospective analysis of observed outcomes, a more tractable approach when data is limited. Our goal is *to avoid future dead-ends* or regions in the state space from which negative outcomes are inevitable (formally defined in Section 3.2). DeD identifies dead-ends via two complementary Markov Decision Processes (MDPs) with a specific reward design so that the underlying value functions will carry special meaning (Section 3.4). These value functions are independently estimated using Deep Q-Networks (DQN) [5] to infer the likelihood of a negative outcome occurring (D-Network) and the reachability of a positive outcome (R-Network). Altogether DeD formally connects the notion of value functions to the dead-end problem, learned directly from offline data.

We validate DeD in a carefully constructed toy domain, and then evaluate real health records of septic patients in an intensive care unit (ICU) setting [22]. Sepsis treatment and onset is a common task in medical RL [23–26] because the condition is highly prevalent [27, 28], physiologically severe [29], costly [30] and poorly understood [31]. Notably, the treatment of sepsis itself may also contribute to a patient’s deterioration [32, 33], thus making treatment avoidance a particularly well-suited objective. We find that DeD confirms the existence of dead-ends and demonstrate that 12% of treatments administered to terminally ill patients reduce their chances of survival, some occurring as early as 24 hours prior to death. The estimated value functions underlying DeD are able to capture significant deterioration in patient health 4 to 8 hours ahead of observed clinical interventions, and that higher-risk treatments possibly account for this delay. Early identification of suboptimal treatment options is of great importance since sepsis treatment has shown multiple interventions within tight time frames (10 to 180 minutes) after suspected onset decreases sepsis mortality [34].

While motivated by healthcare, we propose the use of DeD in safety-critical applications of RL in most data-constrained settings. We introduce a formal methodology that outlines how DeD can be implemented within an RL framework for use with real-world offline data. We construct and train DeD in a generic manner which can readily be used for other data-constrained sequential decision-making problems. In particular, we emphasize that DeD is well suited to analyze high-risk decisions in real-world domains.

2 Related Work

RL in Health: RL has been the subject of much focus in health [17], with particular emphasis on sepsis seeking to develop optimal treatment recommendation policies [23–26, 35–38]. However, with fixed retrospective medical data, an optimal policy that maximizes a patient’s chance of recovery is both computationally and experimentally infeasible. To our knowledge, we are the first to target improved treatment recommendations by avoiding high-risk treatments in a fully offline manner.

Safety in RL: RL has a rich history in safety [39], with recent work attempting to limit high risk actions by constraining parametric uncertainty [40], through alignment between agent and human objectives [41, 42], by directly constraining the agent optimization process to avoid unsafe actions [43], or by improving over a baseline policy [44]. In these settings model performance is evaluated in online settings where more data can be acquired or models can be tested against new cases as well as known baselines. We focus on the more challenging offline setting with limited and non-exploratory data, reflecting the reality of healthcare settings.

Dead-ends: The concept of dead-ends and the corresponding security condition that we build from was proposed by Fatemi et al. [45] in the context of *exploration*. In their work an online RL agent needs to experience various courses of actions from each state, through which it learns optimal behavior. We adapt this approach and expand the theoretical results to an offline RL setting as is found within healthcare—where exploration is untenable—to determine which treatments increase the likelihood of entering a dead-end, based on the patient’s current health state.

Related concepts to dead-ends were introduced by Irpan et al. [46], focused primarily on policy evaluation. The authors introduce a notion of *feasible* states as those that are not *catastrophic* and from which an agent will not immediately fail. Whether or not a state is feasible is determined via positive-unlabeled classification. This inherently differs from our approach where we formally characterize dead-ends and a corresponding security condition through which we can identify treatments to avoid that likely lead to dead-ends¹. Our formalization is discussed in the next section.

3 Methods

3.1 Preliminaries

Our pipeline isolates state construction from value estimation with RL. Therefore we consider episodic Markov Decision Processes (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, T, r, \gamma)$, where \mathcal{S} and \mathcal{A} are the discrete sets of states and treatments²; $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a function that defines the probability of transitioning from state s_t to s_{t+1} if treatment a_t is administered; $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [r_{min}, r_{max}]$ is a finite reward function and $\gamma \in [0, 1]$ denotes a scalar discount factor.

A *policy* $\pi(s, a) = \mathbb{P}[A_t = a | S_t = s]$ defines how treatments are selected, given a state. A *trajectory* is comprised of sequences of tuples (S_t, A_t, R_t, S_{t+1}) with S_0 being the initial state of the trajectory. Sequential application of the policy is used to construct trajectories. The reward collected over the course of a trajectory induces the *return* $G_t = \sum_{j=0}^{\infty} \gamma^j R_{t+j+1}$. We assume that all the returns are finite and bounded. A trajectory is considered *optimal* if its return is maximized. A state-treatment value function $Q^\pi(s, a) = \mathbb{E}^\pi[G_0 | S_0 = s, A_0 = a]$ is defined in conjunction with a policy π to evaluate the expected return of administering treatment a at state s and following π thereafter. The optimal state-treatment value function is defined as $Q^*(s, a) = \max_\pi Q^\pi(s, a)$, which is the maximum expected return of all trajectories starting from (s, a) . We define state value and optimal state value as $V^\pi(s) = \mathbb{E}_{a \sim \pi} Q^\pi(s, a)$ and $V^*(s) = \max_a Q^*(s, a)$.

3.2 Special States

We define a terminal state as the final observation of any recorded trajectory. We focus on two types of terminal state that correspond to positive or negative outcomes. Our goal is to identify all *dead-end* states, from which negative outcomes are unavoidable (happening w.p.1), regardless of future treatments. In safety-critical domains, it is crucial to avoid such states *and* identify the probability with which any available treatment will lead to a dead-end. We also introduce the complementary concept of *rescue* states, from which a positive outcome is *reachable* with probability one. If an agent is in a rescue state, there exists at least one treatment at each time step afterwards which leads to either another rescue state or the eventual positive outcome. The fundamental contrast between dead-end and rescue states is that if the agent enters a rescue state, it *does not* mean the treatment process is done; it rather means that at each time step afterwards there exists at least one treatment to be found and administered until the positive outcome occurs. There might be trajectories starting from a rescue state which include non-rescue states. This is not the case for a dead-end state.

Formally, we augment \mathcal{M} with a non-empty termination set $\mathcal{S}_T \subset \mathcal{S}$, which is the set of all terminal states. Mathematically, a terminal state is absorbing (self-transition w.p.1) with zero reward afterwards. All terminal states are by definition zero-valued, but the transitions to them may be associated with a non-zero reward. We require that, from all states, there exists at least one trajectory with non-zero probability arriving at a terminal state. In an offline setting with limited and non-exploratory data, inducing an optimal policy *is not feasible* [12]; hence, we do not specify the reward function of \mathcal{M} for which standard RL would optimize cumulative rewards, but in later sections present a specific design of reward (and discount factor) to assist in identifying dead-end/rescue states. Finally, the sets of dead-end and rescue states are denoted respectively by \mathcal{S}_D and \mathcal{S}_R . We formally distinguish dead-end/rescue states from the outcome, asserting that $\mathcal{S}_D, \mathcal{S}_R \not\subset \mathcal{S}_T$.

¹A more in depth discussion on the differences between Irpan et al. [46] and this work can be found in Appendix Section A2

²Our results can easily be extended to continuous state-spaces by properly replacing summations with integrals. For brevity, we only present formal proofs for the discrete case. Additionally, as our primary motivating domain lies within healthcare, we use the term “treatment” in place of “action”.

3.3 Treatment Security

When dealing with data-constrained offline scenarios, a core distinction is necessary: Realization of an optimal treatment at a given state requires knowledge of all future outcomes for all possible treatments, which is not feasible. However, the data may contain enough information to estimate the possible outcome of a certain treatment at a similar state. If such an outcome is negative with high probability, then we should advise against the treatment, even if an optimal treatment still remains unknown. This distinction leads to a paradigm shift from finding the best possible treatment to mindful avoidance of dangerous ones. This shift further motivates a different design space to make use of the limited, yet available data.

We adapt the security condition from Fatemi et al. [45] and formalize the treatment avoidance problem with a more generalized *treatment security condition*. We note that the chance of a negative outcome is best described by the probability of falling into a dead-end or immediate negative termination. The security condition therefore constrains the scope of a given behavioral policy π if *any* knowledge exists about dead-ends or negative termination. Formally, if at state s , treatment a leads to a dead-end with probability $P_D(s, a)$ or immediate negative termination with probability $F_D(s, a)$ with a level of certainty $\lambda \in [0, 1]$, then π must avoid selecting a at s with the same certainty:

$$P_D(s, a) + F_D(s, a) \geq \lambda \implies \pi(s, a) \leq 1 - \lambda. \quad (1)$$

E.g., if a treatment leads to a dead-end or termination with probability more than 80%, then that treatment should be selected for administration no more than 20% of the time. While we would like (1) to hold for the maximum λ , inferring such maximal values is intractable for all state-treatment pairs. Moreover, directly computing P_D and F_D would require explicit knowledge of all dead-end and negative terminal states as well as all transition probabilities for future states. These make the application of (1) nearly impossible. We next develop a learning paradigm to enable (1) from data.

3.4 Dead-end Discovery (DeD)

In order to identify and confirm the existence of dead-end states, we construct two Markov Decision Processes (MDPs) \mathcal{M}_D and \mathcal{M}_R to be identical to \mathcal{M} , with $\gamma = 1$ for both. We also define the following reward functions: \mathcal{M}_D returns -1 with any transition to a negative terminal state (and zero with all other transitions) whereas \mathcal{M}_R returns $+1$ with any transition to a positive terminal state (zero otherwise). Let Q_D^* , Q_R^* , V_D^* and V_R^* denote the optimal state-treatment and state value functions of \mathcal{M}_D and \mathcal{M}_R , respectively. Note that due to the reward functions of these MDPs, for all states and treatments, $Q_D^*(s, a) \in [-1, 0]$ and $Q_R^*(s, a) \in [0, 1]$.

Having selected treatment a at state s , using the Bellman equation, we prove³ that

$$-Q_D^*(s, a) = P_D(s, a) + F_D(s, a) + M_D(s, a) \quad (2)$$

In addition to the quantities defined previously, $M_D(s, a)$ denotes the probability of circumstances in stochastic environments where a negative terminal state ultimately occurs despite receiving optimal treatments at all steps in the future. Equation (2) therefore reveals that $-Q_D^*$ carries special physical meaning: it corresponds to the *minimum probability of a negative outcome*, because future treatments may not necessarily be optimal. Equivalently, $1 + Q_D^*(s, a)$ can be seen as the *maximum hope of a positive outcome* if treatment a is administered at state s .

Building from Fatemi et al. [45], we show that V_D^* of all dead-end states will be precisely -1 . By extension, $Q_D^*(s, a) = -1$ for all treatments a at state s if and only if s is a dead-end. In fact, $1 + Q_D^*(s, a)$ provides an appropriate threshold to secure any given policy $\pi(s, a)$. More formally, the following statement guarantees treatment security as presented in (1) for all values of λ :

$$\pi(s, a) \leq 1 + Q_D^*(s, a) \quad (3)$$

In short, for treatment security it is sufficient to abide by the maximum hope of a positive outcome. This construction directly connects the RL concept of value functions to dead-end discovery. While $V_D^*(s)$ enables detecting dead-end states, (3) leverages Q_D^* for treatment avoidance. We establish parallel results for rescue states similarly. The following theorem summarizes the theory and shapes the basis of DeD. See Appendix A1 for the proof and further details.

³All proofs to the theoretical claims presented in this paper can be found in Appendix A1

Theorem 1. Let treatment a be administered at state s , and $P_D(s, a)$ and $P_R(s, a)$ denote the probability of transitioning to a dead-end or rescue state. Similarly, let $F_D(s, a)$ and $F_R(s, a)$ denote the probability of transitioning to either a negative or positive terminal state. The following hold:

- T1 $P_D(s, a) + F_D(s, a) = 1$ if and only if $Q_D^*(s, a) = -1$.
- T2 $P_R(s, a) + F_R(s, a) = 1$ if and only if $Q_R^*(s, a) = 1$.
- T3 There exists a threshold $\delta_D \in (-1, 0)$ independent of states and treatments, such that $Q_D^*(s, a) \geq \delta_D$ for all s and a , unless $P_D(s, a) + F_D(s, a) = 1$.
- T4 There exists a threshold $\delta_R \in (0, 1)$ independent of states and treatments, such that $Q_R^*(s, a) \leq \delta_R$ for all s and a , unless $P_R(s, a) + F_R(s, a) = 1$.
- T5 For any policy π , state s , and treatment a , if $\pi(s, a) \leq 1 + Q_D^*(s, a)$ and $\lambda \in [0, 1]$ exists such that $P_D(s, a) + F_D(s, a) \geq \lambda$, then $\pi(s, a) \leq 1 - \lambda$.
- T6 For any policy π , state s , and treatment a , if $\pi(s, a) \geq Q_R^*(s, a)$ and $\lambda \in [0, 1]$ exists such that $P_R(s, a) + F_R(s, a) \geq \lambda$, then $\pi(s, a) \geq \lambda$.

It is immediate from (T1) and (T2) that Q_D^* and Q_R^* incorporate complete information when transitioning to a dead-end state or to a rescue state as a result of administering treatment a at s . (T3) assures that a threshold δ_D exists to separate treatments that lead immediately to dead-ends from alternatives. (T4) allows us to confirm a dead-end by examining if Q_R^* is also smaller than some threshold δ_R . No dead-end can violate δ_R due to (T4) and such a threshold exists. If Q_D^* is available and δ_D is known, then this step is redundant. However, without access to Q_D^* and an accurate δ_D , (T4) helps to confirm any presumed dead-end. Finally, (T5) provides the means by which the treatment policy is guided to avoid dangerous treatments. (T6) is used to also confirm whether the treatment should be avoided. We explain how to practically select the thresholds δ_D and δ_R in Sec. 5.

Of note, by definition, value functions encompass long-term consequences and are not myopic to possible immediate events, as opposed to supervised learning from immediate observation of an outcome. This inherent characteristic of value functions indeed yields the theoretical result presented by Lemma 2 (Appendix Sec. A1), one result of which is that $-Q_D$ corresponds to the minimum probability of a negative future outcome. Supervised learning from immediate outcomes, on the other hand, lacks this formal property [47]; hence, it is not expected to provide parallel results with DeD.

3.5 Neural Network Based State Construction and Identification

State construction (SC-Network). In domains where solitary observations do not carry salient information for learning the decision-making process, states may need to be constructed from data using a neural network. In these circumstances a separate SC-Network can be used to transform a single or possible sequence of observations into a fixed embedding, considered the state s at time t .

Identification (D-Network and R-Network). In order to approximate Q_D^* and Q_R^* , two separate neural networks can be used to compute Q_D and Q_R for all treatments given a state constructed by the SC-Network. With trained Q_D and Q_R networks, we can then apply thresholds δ_D and δ_R as specified in Theorem 1. As data is limited and non-exploratory, approximation error is inevitable. To mitigate this limitation, the method’s sensitivity can be adjusted by adapting the thresholds δ_D and δ_R (additionally, see Proposition 1 and Remarks 1-4 in Appendix A1). Smaller thresholds result in more false negative and less false positive cases. Of note, value-overestimation, a known limitation of deep RL models, will often cause Q_D and Q_R to be larger than Q_D^* and Q_R^* respectively. This naturally reduces false positives while increasing false negatives.

3.6 Toy Problem Validation: Life-Gate

We briefly provide a tabular toy-example (Life-Gate), which involves dead-ends, to empirically illustrate the merit of Theorem 1 by learning Q_D^* and Q_R^* (Figure 1). This toy set-up comprises an interesting case, where the agent faces an environment to examine with no knowledge of possible dangers. Importantly, once a dead-end state (yellow) is reached, it may take some random number of steps before reaching a “death gate” (red). All along such trajectories of dead-end states, the agent still has to choose actions with the (false) hope of reaching a “life-gate” (blue). Discovering any single dead-end state and signaling the agent when it is approached would be of significant

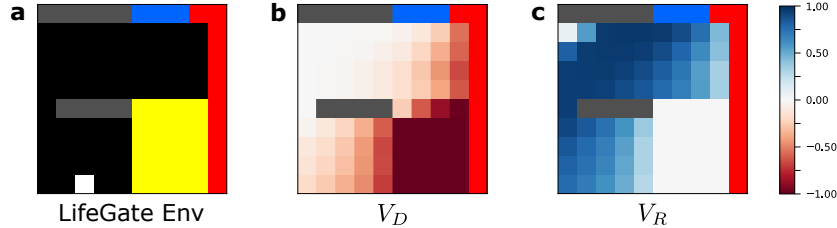


Figure 1: **The Life-Gate Example.** The tabular navigation task of life-gate is illustrated in (a). Corresponding dead-end and rescue state-value functions, V_D and V_R , are shown in (b) and (c). The value functions are learned through Q-learning and with the definition of \mathcal{M}_D and \mathcal{M}_R .

importance. On the other hand, adjacent states to dead-ends are possibly the most critical to alert, as it might be the last chance to still do something to avoid failure (see Appendix A3 for more details).

We next use the tools provided by Theorem 1. The value functions are more than 90% trained, still allowing learning errors. In this example, even with the errors due to lack of full convergence, $\delta_D = -0.7$ and $\delta_R = 0.7$ seem to clearly set the boundary for most states (with a few exceptions due to the errors). If a state is observed whose V_D and V_R values violate these thresholds, the state can be flagged as a dead-end with high probability. Setting a lower threshold can help to raise flags earlier on, when the conditions are of high-risk, but it is still not too late. We can apply the same thresholds to further flag high-risk actions (not shown). Lastly, we note (T1) from Theorem 1. It can be seen that only for all the yellow area (aside from the few erroneous states), $V_D = -1$. Clearly, no dead-end state can be a rescue, as seen by $V_R = 0$ for the yellow area too.

4 Empirical Setup for Dead-end Analysis

Data: We use DeD to identify medical dead-ends in a cohort of septic patients drawn from the MIMIC (Medical Information Mart for Intensive Care) - III dataset (v1.4) [22, 48]. This cohort totals 19,611 patients (17,730 survivors and 1,881 nonsurvivors), with 44 observation variables, and 25 treatment choices (5 discrete levels for each of IV fluid and vasopressor). We follow prior work [25] and aggregate each variable every four hours using the per-patient variable mean if data is present, or impute using the value from the nearest neighbor.

Terminal States. In our ICU setting, possible terminal states are either patient recovery (discharge from ICU) or death. We define “death” as the last recorded point in the EMR of nonsurviving patients when expiration is imminent, but may not necessarily be the biological point of death. In practice this definition of terminal state may occur hours or days before biological death and covers situations where care support devices are disconnected, when a patient requests a cessation of treatment, etc.

Our goal is to identify all *medical dead-end* states, defined as patient states from which death is unavoidable, regardless of future treatments. Relatedly, we also desire to discover all treatments that may possibly lead to a medical dead-end state in order to learn which treatments to avoid.

SC-Network. As observations of patient health are inherently partial, we need an informative latent representation of state [49], sufficient for evaluating treatment security. To form these state representations we process a sequence of observations prior to and including any time t as well as the last selected treatment to form the state s_t . We train a standalone State Construction (SC) network using Approximate Information State (AIS) [50] in a self-supervised manner for this purpose. Details of AIS and how it is used to train the SC-Network are included in Appendix A4.

D-Network and R-Network. Computed states are given as input to the D- and R-networks to approximate Q_D^* and Q_R^* . We use the double-DQN algorithm [51] to train each network (details included in Appendix A5). The outputs of trained D- and R- Networks produce value estimates of both the embedded patient state and all possible treatments to evaluate the probability of transitioning to a dead-end. This process of determining possibly high-risk treatments is central to DeD.

Training: We train the SC-, D-, and R- networks in an offline manner, using retrospective data (Fig. A2). All models are trained with 75% of the patient cohort (14,179 survivors, 1,509 nonsurvivors),

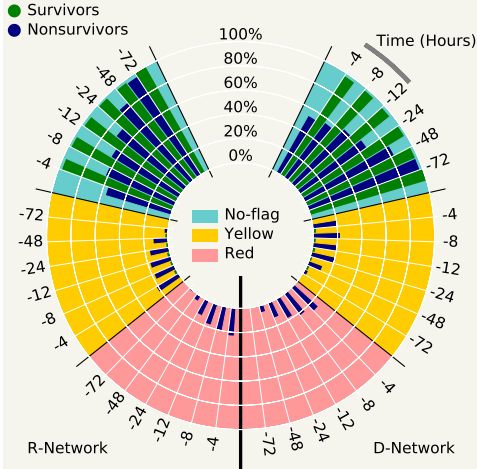


Figure 2: **Flag emergence for ICU patients.** Histograms of median Q according to the flag status, for both surviving (green) and non-surviving patients (blue) according to the R-Network (left) and D-Network (right). The bars are plotted according to the time prior to the recorded terminal state (the maximum trajectory length is 72 hours) and measure the percentage of patients whose states raise either a red, yellow or no flag. There is a clear worsening trend of state values for nonsurviving patients as they approach their terminal state, beginning as early as 48 hours prior.

validated with 5% (890 survivors, 90 nonsurvivors), and we report all results on the remaining held out 20% (2,660 survivors, 282 nonsurvivors). Further details of how the patient cohort is processed are provided in Appendix A6. Finally, to mitigate the data imbalance between surviving and non-surviving patients we use an additional data buffer that contains only *the last transition* of nonsurvivors trajectories. Thus, a stratified minibatch of size 64 is constructed of 62 samples from the main data, augmented with 2 samples from this additional buffer, all selected uniformly. This same minibatch structure is used for training each of the three networks. For the training details see Appendix A4 and A5.

5 Empirical Results

5.1 Septic Dead-End State Prediction

Experiment. In order to flag potentially non-secure treatments, we examine if Q_D and Q_R of each treatment at a given state pass certain thresholds δ_D and δ_R , respectively. To flag potential dead-end states, we need to probe the state values, for which we examine the *median* of Q (rather than *max* of Q) against similar thresholds. Using the median helps to avoid extreme approximation error due to generalization from potentially insufficient data. We found that $\delta_D = -0.25$ and $\delta_R = 0.75$ minimize both false positives and false negatives, and use these as the thresholds for raising “red” flags. We also define a second, looser threshold of $\delta_D = -0.15$ and $\delta_R = 0.85$, as raising “yellow” flags with higher sensitivity but increased false positives. This looser threshold targets an early indication of a patient’s health condition deteriorating toward a dead-end state. In Appendix Fig. A5 we report histogram of values at different quantiles, from which we established these thresholds.

Results. Using the specified thresholds, DeD identifies increasing percentages of patients raising fatal flags as nonsurvivors approach death (Figure 2 and Appendix Table A3). Note the distinctive difference between the trend of values in survivors (green bars) and nonsurvivors (blue bars). Over the course of 72 hours in the ICU, survivor trajectories raise nearly no red flag for both networks. In contrast, nonsurvivor trajectories demonstrate a steep reduction in *no-flag* zone with increasing numbers of patients flagged in the *Red* zone. The *Yellow* zone is dominated largely by the nonsurvivors, yet there are also survivors who ultimately recover. Under the red-flag threshold, more than 12% of treatments administered to non-surviving patients are identified to be detrimental 24 hours prior to death with a 0.6% false positive rate (Appendix Table A3). We further identify that 2.7% of non-surviving patient cases have entered unavoidable dead-end trajectories up to 48 hours before recorded expiration, with only a 0.2% false positive rate, i.e., patients misidentified as near death.

We find that 5% of nonsurviving patients maintain the red flag for their last 24 hours recorded in the ICU before reaching a death terminal state. This monotonically increases to 13.9% for patients who maintain a red flag through their final 8 hours of care (Appendix Fig. A6b,c). These patients likely reached a dead-end with no subsequent chance of recovery; this is as compared to 89.3% of nonsurviving test patients with no flag raised in their first 8 hours (Appendix Fig. A6d).

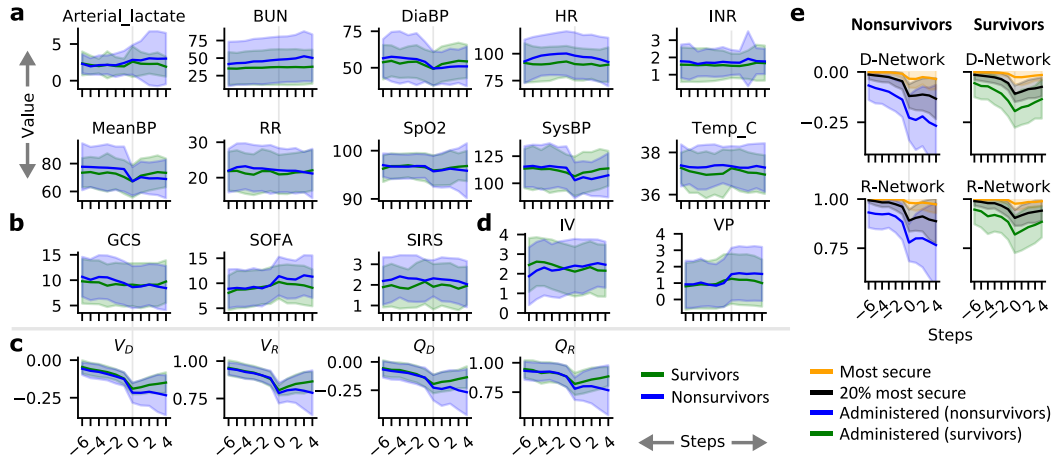


Figure 3: **Trend of measures around the first raised flag.** Various measures are shown 24 hours (6 steps) before the first (red or yellow) flag is raised and 16 hours (4 steps) afterward. All nonsurviving (blue) and surviving (green) patient trajectories that fall within this window are averaged, shaded areas represent a single standard deviation. (a) selected key vital measures and lab tests, (b) established clinical measures, and (c) DeD value measures of state (V) and administered treatment (Q) from the D- and R-Networks and, (d) administered treatments. There is a clear turning point 4 to 8 hours prior to the flag being raised, which precisely corresponds to a drastic increase of VP and IV treatments. (e) the value of the maximum, the 5th maximum (20% best) and the actually administered treatment, demonstrating that better treatments were available when the chosen treatments were administered.

There is a distinct difference between remaining on a flag for survivors and nonsurvivors (Appendix Fig. A6a). Even with our red threshold, very few survivors (0.5%) raise and remain on red-flag for more than eight hours, decreasing to nearly zero for longer periods. In contrast, 32.6% of nonsurvivors remain on red flags for similar duration with a fat tail. These results suggest that red-flag membership for long periods strongly correlates with mortality, inline with our theoretical analysis.

5.2 First Flag Analysis

Experiment. To further support our hypothesis that dead-end states exist among septic patients and may be preventable, we align patients according to the point in their care when a flag is “first raised”. We select all trajectories in the test data with at least 24 hours (6 steps) prior to the first flag and at least 16 hours (4 steps) afterwards (77 surviving and 74 nonsurviving patients). This window excludes patients with flags that occur either too early or too late. This allows for an investigation of the average trend of patient observations, administered treatments as well as the measures used in DeD over a sufficiently large window (Figure 3).

Results. The V and Q values estimated by DeD have similar behavior in survivors and nonsurvivors prior to the first flag, but values diverge after the flag is raised. Notably, the time step pinpointed by DeD to raise a flag corresponds to a similar diverging trend among various clinical measures, including SOFA and patient vitals (Figure 3a,b). This distinct behavior is also seen if looser threshold values are used for δ_D and δ_R (Appendix Fig. A8). After the flag is raised there is slight improvement in all value estimates, perhaps in response to the change in treatment. However the values of nonsurviving patient trajectories quickly collapse while survivors continue to improve.

The results of this analysis suggest two main points. First, DeD identifies a clear *critical point* in the care timeline where nonsurviving patients experiencing an irreversible deterioration in health. Second, there is a significant gap (Figure 3e) between the value of administered treatments and the 20% most secure ones (5 out of 25). The critical point appears to arise when a patient’s condition shifts towards improvement or otherwise enters a dead-end towards expiration. Perhaps most notable is that there is a clear inflection in the estimated values 4 to 8 hours prior to a flag being raised. Signaling this shift in the inferred patient response to treatment and the resulting flag may be used to provide an early indicator for clinicians (more conservative thresholds may be used to signal earlier). The trend of survivors shows that there is still hope to save the patient at this point. Note that *all* these patients

of leading to dead-ends, regions of the state space from which negative outcomes are inevitable. We establish theoretical results that expand the concept of dead-ends in RL, facilitating the notification of high-risk treatments or, as applied to healthcare, septic patient conditions with increased likelihood of leading to a dead-end. Globally, sepsis is a leading cause of mortality [27, 53, 54], and an important end-stage to many conditions. Consequently, even a slight decrease in mortality rate or improved efficacy of treatment could have a significant impact both in terms of saving lives and reducing costs.

Our work lays the groundwork for dead-end analysis in medical settings and is, to the best of our knowledge, the first use of RL to flag bad treatments rather than finding the best ones through estimating an optimal policy π^* . Our algorithm is generic, using RL methodology that is formally guaranteed to hold the security condition re-established in this paper. The discovery of dead-end states, and the treatments that likely lead to them, provides actionable insights in intensive care intervention. Further improvement of DeD’s prediction quality could target additional features from the EMR environment, such as pre-ICU admission co-morbidities. In future work, we also hope to explore the specific drugs and dosages used in treatment.

Given its general construction, DeD is well matched for safety-critical applications of RL in data-constrained settings where it may be too expensive or unethical to collect additional exploratory data. With formal guarantees of satisfying the security condition, DeD is suitable for broader adoption when developing critical insights from retrospective data. Our framework is particularly relevant to data-constrained offline RL application domains such as robotics, industrial control, and automated dialogue generation where negative outcomes can be clearly identified [55].

Limitations: While DeD is a promising framework for decision support in safety-critical domains with limited offline data, there are certain core limitations. While we use median values of Q_D or Q_R to avoid extreme extrapolation, training the D- and R- networks is still performed offline and extrapolation is likely still occurring. For simplicity we did not estimate Q_D or Q_R with contemporary offline RL methods; however, DeD is generic and replacing the DDQN learning method with more recent approaches would be straightforward, which can significantly improve the pipeline (we also note that finding Q_D or Q_R is an *exponentially smaller* problem compared to finding π^* to recommend best treatments). Additionally, we did not investigate the sensitivity of DeD to demographic information or with respect to specific features from the EMR. Thorough analysis of this sensitivity may elucidate the fairness and reliability of DeD. Finally, we did not externally validate DeD using data from a separate hospital or through investigation of suggested treatment avoidance by human clinicians. These investigations and more, concerning the causal entanglement of outcome and sequential treatments, are a focus of current and future work.

Ethical Considerations and Societal Impact: This work, or derivatives of it, should never be used in isolation to exclude patients from being treated, e.g., not admitting patients or blindly ignore treatments. The treatment-avoidance part of our proposed approach is meant to shrink the scope of possible treatment options, and help the doctors make better decisions. Signalling high-risk states is also meant to warn the clinicians for immediate attention before it possibly becomes too late. In both cases, the flags that DeD supplies are statistically tied to the training data and unavoidable sources of error and bias and should not be seen as a binary treat/don’t treat decision. In particular, even in the case of red flags, the signals should not be interpreted as mathematical dead-ends with full precision. The intention of our approach is to assist clinicians by highlighting possibly unanticipated risks when making decisions and is not to be used as a stand-alone tool nor as a replacement of a human expert. Misuse of this algorithmic solution could carry significant risk to the well-being and survival of patients placed in a clinician’s care.

The primary goal of this work is to establish a proof of concept where especially high-risk treatments can be avoided, where possible, in context of a patient’s health condition. In acute care scenarios treatments come with inherent risk profiles and potential harms. In these settings tendencies to overtreat patients have arisen in attempt of ensuring their survival, increasing the chance of clinical errors to occur [56]. Recent clinical research has sought to simplify practice to only the most necessary treatments⁴. In this spirit, we seek to infer the long-term impact of each available treatment in view of their risk of pushing the patient into a medical dead-end. The secondary goal of our work, on the other hand, is signal when the patient’s condition deteriorates, but may not be noticed by clinicians through monitoring clinical measures. This follows from the fact that DeD uses value functions, which provably enable such predictions.

⁴see <http://jamanetwork.com/collection.aspx?categoryid=6017>

Acknowledgments and Disclosure of Funding

We thank our many colleagues who contributed to thoughtful discussions and provided timely advice to improve this work. We specifically appreciate the feedback provided by Nathan Ng, Sindhu Gowda, and the RL team at MSR Montreal, as well as the encouragement and suggested improvements provided by anonymous reviewers.

This research was supported in part by Microsoft Research, a CIFAR Azrieli Global Scholar Chair, a Canada Research Council Chair, and an NSERC Discovery Grant.

Resources used in performing this research were provided, in part, by Microsoft Research, the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute www.vectorinstitute.ai/#partners.

Data and Code Availability

Our code and pretrained models to replicate the analysis (including figures) presented in this paper is located at <https://github.com/microsoft/med-deadend>.

The MIMIC-III databases (DOI: 10.1038/sdata.2016.35) that support the findings of this study are publicly available through Physionet website: <https://mimic.physionet.org>, which facilitates reproducibility of the presented results. The cohort definition, extraction and preprocessing code can be found at https://github.com/microsoft/mimic_sepsis.

Author Contributions

MF and MG designed the research. MF conceptualized theoretical ideas and developed formal results and proofs. JS developed the code for data generation and state construction and prepared initial experimental results. TK finalized the script to generate data from MIMIC, performed benchmarking of several state-construction algorithms—finalizing the decision to use AIS in the paper—and executed the experiments for the Life-Gate toy example. MF developed code for RL and made final experimental results. MF, TK and MG interpreted the results and wrote the manuscript. JS contributed to this work only during his internship at Microsoft Research. All the authors read and agreed on the final draft.

References

- [1] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998. ISBN 0262193981.
- [2] Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- [3] Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. In *AAMAS*, pages 1077–1084, 2014.
- [4] Susan A Murphy, Mark J van der Laan, James M Robins, and Conduct Problems Prevention Research Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 96(456):1410–1423, 2001.
- [5] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [6] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [7] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. IMPALA:

- Scalable distributed deep-RL with importance weighted actor-learner architectures. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1407–1416. PMLR, 10–15 Jul 2018. URL <http://proceedings.mlr.press/v80/espeho1t18a.html>.
- [8] Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. In *Reinforcement learning*, pages 45–73. Springer, 2012.
- [9] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- [10] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pages 2052–2062. PMLR, 2019.
- [11] Dimitri P. Bertsekas and John N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1st edition, 1996. ISBN 1886529108.
- [12] Harold Kushner and George Yin. *Stochastic Approximation and Recursive Algorithms and Applications*. Springer-Verlag, 2003. doi: 10.1007/b97441.
- [13] Vincent François-Lavet, Guillaume Rabusseau, Joelle Pineau, Damien Ernst, and Raphael Fonteneau. On overfitting and asymptotic bias in batch reinforcement learning with partial observability. *Journal of Artificial Intelligence Research*, 65:1–30, 2019.
- [14] Samarth Sinha and Animesh Garg. S4rl: Surprisingly simple self-supervision for offline reinforcement learning. *arXiv preprint arXiv:2103.06326*, 2021.
- [15] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*, pages 104–114. PMLR, 2020.
- [16] Ramesh Rebba, Sankaran Mahadevan, and Shuping Huang. Validation and error estimation of computational models. *Reliability Engineering & System Safety*, 91(10-11):1390–1397, 2006.
- [17] Chao Yu, Jiming Liu, and Shamim Nemati. Reinforcement learning in healthcare: a survey. *arXiv preprint arXiv:1908.08796*, 2019.
- [18] Omer Gottesman, Fredrik Johansson, Matthieu Komorowski, Aldo Faisal, David Sontag, Finale Doshi-Velez, and Leo Anthony Celi. Guidelines for reinforcement learning in healthcare. *Nature Medicine*, 25(1):16–18, January 2019. doi: 10.1038/s41591-018-0310-5. URL <https://doi.org/10.1038/s41591-018-0310-5>.
- [19] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. In *Advances in Neural Information Processing Systems*, pages 11784–11794, 2019.
- [20] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- [21] Ziyu Wang, Alexander Novikov, Konrad Żołna, Jost Tobias Springenberg, Scott Reed, Bobak Shahriari, Noah Siegel, Josh Merel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized regression. In *Advances in Neural Information Processing Systems*, 2020.
- [22] 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.
- [23] Katharine E Henry, David N Hager, Peter J Pronovost, and Suchi Saria. A targeted real-time early warning score (trewscore) for septic shock. *Science translational medicine*, 7(299):299ra122–299ra122, 2015.

- [24] Joseph Futoma, Sanjay Hariharan, Katherine Heller, Mark Sendak, Nathan Brajer, Meredith Clement, Armando Bedoya, and Cara O'Brien. An improved multi-output gaussian process rnn with real-time validation for early sepsis detection. In *Machine Learning for Healthcare Conference*, pages 243–254, 2017.
- [25] 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.
- [26] Suchi Saria. Individualized sepsis treatment using reinforcement learning. *Nature medicine*, 24(11):1641–1642, 2018.
- [27] Clifford S. Deutschman and Kevin J. Tracey. Sepsis: Current dogma and new perspectives. *Immunity*, 40(4):463 – 475, 2014. ISSN 1074-7613. doi: <https://doi.org/10.1016/j.immuni.2014.04.001>. URL <http://www.sciencedirect.com/science/article/pii/S1074761314001150>.
- [28] Mervyn Singer, Clifford S. Deutschman, Christopher Warren Seymour, Manu Shankar-Hari, Djillali Annane, Michael Bauer, Rinaldo Bellomo, Gordon R. Bernard, Jean-Daniel Chiche, Craig M. Coopersmith, Richard S. Hotchkiss, Mitchell M. Levy, John C. Marshall, Greg S. Martin, Steven M. Opal, Gordon D. Rubenfeld, Tom van der Poll, Jean-Louis Vincent, and Derek C. Angus. The third international consensus definitions for sepsis and septic shock (sepsis-3). *JAMA*, 315(8):801, feb 2016. doi: 10.1001/jama.2016.0287. URL <https://doi.org/10.1001%2Fjama.2016.0287>.
- [29] Jean-Louis Vincent, Steven M Opal, John C Marshall, and Kevin J Tracey. Sepsis definitions: time for change. *The Lancet*, 381(9868):774 – 775, 2013. ISSN 0140-6736. doi: [https://doi.org/10.1016/S0140-6736\(12\)61815-7](https://doi.org/10.1016/S0140-6736(12)61815-7). URL <http://www.sciencedirect.com/science/article/pii/S0140673612618157>.
- [30] C Torio and B Moore. National inpatient hospital costs: The most expensive conditions by payer, 2013. Statistical Brief #204. *Healthcare Cost and Utilization Project (HCUP) Statistical Briefs*, May 2016. URL <http://www.hcup-us.ahrq.gov/reports/statbriefs/sb204-Most-Expensive-Hospital-Conditions.pdf>.
- [31] Greg S Martin. Sepsis, severe sepsis and septic shock: changes in incidence, pathogens and outcomes. *Expert review of anti-infective therapy*, 10(6):701–706, 2012.
- [32] P. E. Marik, W. T. Linde-Zwirble, E. A. Bittner, J. Sahatjian, and D. Hansell. Fluid administration in severe sepsis and septic shock, patterns and outcomes: an analysis of a large national database. *Intensive Care Med*, 43(5):625–632, May 2017.
- [33] J. Waechter, A. Kumar, S. E. Lapinsky, J. Marshall, P. Dodek, Y. Arabi, J. E. Parrillo, R. P. Dellinger, A. Garland, S. Dial, P. Ellis, D. Feinstein, D. Gurka, J. Guzman, S. Keenan, A. Kramer, A. Kumar, D. Laporta, K. Laupland, B. Light, D. Maki, G. Martinka, Z. Memish, Y. Mirzanejad, G. Patel, C. Penner, D. Roberts, J. Ronald, D. Simon, S. Sharma, N. A. Shirawi, Y. Skrobik, G. Wood, K. E. Wood, S. Zanotti, M. W. Ahsan, M. Bahrainian, R. Bohmeier, L. Carter, H. Chou, S. Delgra, C. Egbujuo, W. Fu, C. Gonzales, H. Gulati, E. Halmarson, J. Hansen, Z. Haque, J. Harvey, F. Khan, L. Kolesar, L. Kravetsky, R. Kumar, N. Merali, S. Muggaberg, H. Paulin, C. Peters, J. Richards, C. Schorr, H. Serrano, M. Suleman, A. Singh, K. Sullivan, R. Suppes, L. Taiberg, R. Tchokonte, O. A. Torshizi, and K. Wiebe. Interaction between fluids and vasoactive agents on mortality in septic shock: a multicenter, observational study. *Crit. Care Med.*, 42(10):2158–2168, Oct 2014.
- [34] Robert Gauer, Damon Forbes, and Nathan Boyer. Sepsis: diagnosis and management. *American family physician*, 101(7):409–418, 2020.
- [35] Aniruddh Raghu, Matthieu Komorowski, Leo Anthony Celi, Peter Szolovits, and Marzyeh Ghassemi. Continuous state-space models for optimal sepsis treatment: a deep reinforcement learning approach. In *Machine Learning for Healthcare Conference*, pages 147–163. PMLR, 2017.

- [36] Xuefeng Peng, Yi Ding, David Wihl, Omer Gottesman, Matthieu Komorowski, Li-wei H Lehman, Andrew Ross, Aldo Faisal, and Finale Doshi-Velez. Improving sepsis treatment strategies by combining deep and kernel-based reinforcement learning. In *AMIA Annual Symposium Proceedings*, volume 2018, page 887. American Medical Informatics Association, 2018.
- [37] Luchen Li, Matthieu Komorowski, and Aldo A Faisal. Optimizing sequential medical treatments with auto-encoding heuristic search in POMDPs. *arXiv preprint arXiv:1905.07465*, 2019.
- [38] Shengpu Tang, Aditya Modi, Michael Sjoding, and Jenna Wiens. Clinician-in-the-loop decision making: Reinforcement learning with near-optimal set-valued policies. In *International Conference on Machine Learning*, pages 9387–9396. PMLR, 2020.
- [39] Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- [40] Philip S Thomas. *Safe reinforcement learning*. PhD thesis, University of Massachusetts Libraries, 2015.
- [41] Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart Russell. Cooperative inverse reinforcement learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 3916–3924, 2016.
- [42] Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart Russell, and Anca D Dragan. Inverse reward design. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6768–6777, 2017.
- [43] Philip S Thomas, Bruno Castro da Silva, Andrew G Barto, Stephen Giguere, Yuriy Brun, and Emma Brunskill. Preventing undesirable behavior of intelligent machines. *Science*, 366(6468): 999–1004, 2019.
- [44] Romain Laroche, Paul Trichelair, and Remi Tachet des Combes. Safe policy improvement with baseline bootstrapping. In *International Conference on Machine Learning (ICML)*, June 2019.
- [45] Mehdi Fatemi, Shikhar Sharma, Harm Van Seijen, and Samira Ebrahimi Kahou. Dead-ends and secure exploration in reinforcement learning. In *International Conference on Machine Learning*, pages 1873–1881, 2019.
- [46] Alexander Irpan, Kanishka Rao, Konstantinos Bousmalis, Chris Harris, Julian Ibarz, and Sergey Levine. Off-policy evaluation via off-policy classification. *Advances in Neural Information Processing Systems*, 32:5437–5448, 2019.
- [47] Jason H Maley, Kerollos N Wanis, Jessica G Young, and Leo A Celi. Mortality prediction models, causal effects, and end-of-life decision making in the intensive care unit. *BMJ Health & Care Informatics*, 27(3), 2020.
- [48] Alistair Ew Johnson, David J Stone, Leo A Celi, and Tom J Pollard. The MIMIC code repository: enabling reproducibility in critical care research. *J. Am. Med. Inform. Assoc.*, 25(1):32–39, January 2018.
- [49] Taylor W Killian, Haoran Zhang, Jayakumar Subramanian, Mehdi Fatemi, and Marzyeh Ghassemi. An empirical study of representation learning for reinforcement learning in healthcare. In *Machine Learning for Health*, pages 139–160. PMLR, 2020. URL <http://proceedings.mlr.press/v136/killian20a>.
- [50] Jayakumar Subramanian and Aditya Mahajan. Approximate information state for partially observed systems. In *Conference of Decision and Control (CDC), Nice, France*, 2019.
- [51] Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI’16, page 2094–2100. AAAI Press, 2016.
- [52] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

- [53] Jean-Louis Vincent, John C Marshall, Silvio A Ndamendys Silva, Bruno François, Ignacio Martin-Loeches, Jeffrey Lipman, Konrad Reinhart, Massimo Antonelli, Peter Pickkers, Hassane Njimi, Edgar Jimenez, and Yasser Sakr. Assessment of the worldwide burden of critical illness: the intensive care over nations (icon) audit. *The Lancet Respiratory Medicine*, 2(5):380 – 386, 2014. ISSN 2213-2600. doi: [https://doi.org/10.1016/S2213-2600\(14\)70061-X](https://doi.org/10.1016/S2213-2600(14)70061-X). URL <http://www.sciencedirect.com/science/article/pii/S221326001470061X>.
- [54] C. Fleischmann, A. Scherag, N. K. Adhikari, C. S. Hartog, T. Tsaganos, P. Schlattmann, D. C. Angus, and K. Reinhart. Assessment of Global Incidence and Mortality of Hospital-treated Sepsis. Current Estimates and Limitations. *Am. J. Respir. Crit. Care Med.*, 193(3):259–272, Feb 2016.
- [55] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [56] Aaron E Carroll. The high costs of unnecessary care. *Jama*, 318(18):1748–1749, 2017.
- [57] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- [58] Matthieu Komorowski. AI Clinician. https://github.com/matthieukomorowski/AI_Clinician, 2018. Accessed: 2019-08-16.

A1 Formal Results and Proofs

For simplicity, in all the arguments below, we refer to a positive terminal state as *recovery* and a negative terminal state as *death*. As with the main text, we also use *treatment* in place of *action*, which is the common term in RL texts. The rest of terminology follows the definitions presented in the main text.

Lemma 1.

L1.1. $V_D^*(s) = Q_D^*(s, a) = -1$ for all the treatments a if and only if s is a dead-end.

L1.2. $V_R^*(s) = \max_a Q_R^*(s, a) = 1$ if and only if s is a rescue.

Proof. To prove the first part of the lemma, we assume s is a dead-end and prove $Q_D^*(s, a) = -1$ for all the treatments. The definition of return directly implies that the return of all the trajectories from s is precisely -1 since all of them reach a death terminal state and $\gamma = 1$. The expected return is therefore also -1 regardless of stochasticity; hence, $Q_D^*(s, a) = -1$.

Conversely, let for a given state s we have $Q_D^*(s, a) = -1$ for all treatments a . We next prove that s is a dead-end. For a transition (s, a, s') , the next state s' is either of a non-terminal state with $r_D(s, a, s') = 0$, $\max_{a'} Q_D^*(s', a') = -1$; a non-terminal state with $r_D(s, a, s') = 0$, $\max_{a'} Q_D^*(s', a') > -1$; a death terminal state (i.e., $r_D(s, a, s') = -1$, $Q_D^*(s', a') = 0 \forall a'$); or a recovery terminal state (i.e., $r_D(s, a, s') = 0$, $Q_D^*(s', a') = 0 \forall a'$).

Let C_R and C_N denote respectively the sets of “recovery terminal states” and “non-terminal states s' with $\max_{a'} Q_D^*(s', a') > -1$ ”. Note that C_R and C_N are disjoint, and that if a state s' is not in $C_R \cup C_N$ then it is either a death terminal state (hence $r_D(\cdot, \cdot, s') = -1$ and $Q_D^*(s', \cdot) = 0$), or a non-terminal with -1 value (hence $r_D(\cdot, \cdot, s') = 0$ and $Q_D^*(s', \cdot) = -1$). Using Bellman equation, we write

$$\begin{aligned}
 -1 &= Q_D^*(s, a) = \sum_{s'} T(s, a, s') [r_D(s, a, s') + \max_{a'} Q_D^*(s', a')] \\
 &= \sum_{s' \notin C_R \cup C_N} T(s, a, s') \times -1 + \sum_{s' \in C_R} T(s, a, s') \times 0 + \sum_{s' \in C_N} T(s, a, s') \max_{a'} Q_D^*(s', a') \\
 &= - \left[1 - \sum_{s' \in C_R \cup C_N} T(s, a, s') \right] + \sum_{s' \in C_N} T(s, a, s') \max_{a'} Q_D^*(s', a') \\
 &= -1 + \sum_{s' \in C_R} T(s, a, s') + \sum_{s' \in C_N} T(s, a, s') \left[1 + \max_{a'} Q_D^*(s', a') \right] \tag{4}
 \end{aligned}$$

Because $T(s, a, s')$ is non-negative it therefore must be zero for both all $s' \in C_R$ and all $s' \in C_N$ (in the last line $\max_{a'} Q_D^*(s', a') \neq -1$). Hence, the next state is either a death terminal state or a non-terminal state with $Q_D^*(s', \cdot)$ of precisely -1 for all the treatments. Continuing with the same line of argument, it therefore follows that if $Q_D^*(s, a) = -1$ then all possible trajectories after (s, a) will reach a death terminal state and all states on such trajectories assume the value of -1. Finally, if $V_D^*(s) = -1$ then $\max_a Q_D^*(s, a) = -1$, which implies $Q_D^*(s, a) = -1$ for all a . It therefore follows that all trajectories from s will reach a death terminal state, and by definition s is a dead-end.

To prove L1.2, for the sufficiency we cannot use a similar argument as for L1.1 since not all the returns are $+1$; only the maximum needs to be $+1$. If s is a rescue state, then by definition there must exist at least one trajectory w.p.1 to recovery. Starting from the last state before recovery on such a trajectory, we go backward and invoke Bellman equation. For the last state-treatment (s'', a'') that transitions to recovery we have $Q_R^*(s'', a'') = +1$, hence $\max_{a'} Q_R^*(s'', a') = +1$. Similarly, for all other states s' on the deterministic trajectory to recovery, we conclude that $\max_{a'} Q_R^*(s', a') = +1$, which implies $\max_a Q_R^*(s, a) = +1$, as stated in the lemma.

For the necessity, the argument is similar to that of L1.1. In particular, we can show in a similar way as in L1.1 that if $Q_R^*(s, a) = 1$ then $\max_{a'} Q_R^*(s', a') = 1$ for all the immediate next states s' after (s, a) if they are non-terminal. It implies that if s' is non-terminal, then at least one treatment exists whose value at s' is $+1$. Furthermore, for all state-treatment pairs whose values are $+1$, if the treatment causes transitioning to a terminal state it deterministically must be recovery (i.e., it cannot be recovery w.p. p and death w.p. $1 - p$). We therefore conclude that there is at least one trajectory from s to recovery with probability one; hence s is a rescue state. ■

Lemma 2. Let treatment a be administered at state s , and $F_D(s, a)$ and $F_R(s, a)$ denote the probability that the next state will be terminal with death or recovery, respectively. Let further $P_D(s, a)$ and $P_R(s, a)$ denote the probability of transitioning to a dead-end or a rescue state, respectively, i.e. $P_D(s, a) = \sum_{s' \in \mathcal{S}_D} T(s, a, s')$ and $P_R(s, a) = \sum_{s' \in \mathcal{S}_R} T(s, a, s')$. Let $M_D(s, a)$ be the probability that the next state is neither a dead-end nor immediate death, and the patient ultimately expires while all the treatments are selected according to the greedy policy with respect to Q_D^* . Similarly, let $M_R(s, a)$ be the probability that the next state is neither immediate recovery nor a rescue state, but the patient ultimately recovers while all future treatments are selected according to the greedy policy with respect to Q_R^* . We have

$$\text{L2.1. } -Q_D^*(s, a) = P_D(s, a) + M_D(s, a) + F_D(s, a)$$

$$\text{L2.2. } Q_R^*(s, a) = P_R(s, a) + M_R(s, a) + F_R(s, a)$$

Proof. For the first part, Bellman equation reads as the following:

$$Q_D^*(s, a) = \sum_{s'} T(s, a, s') [r_D(s, a, s') + \max_{a'} Q_D^*(s', a')] \quad (5)$$

The next state s' is either of the following:

1. a dead-end state, where $r_D(s, a, s') = 0$; $Q_D(s', a') = -1$, $\forall a'$ (due to Lemma 1); and $\sum_{s'} T(s, a, s') = P_D(s, a)$,
2. a death terminal state, where $r_D(s, a, s') = -1$; $Q_D(s', a') = 0$, $\forall a'$; and $\sum_{s'} T(s, a, s') = F_D(s, a)$,
3. a recovery terminal state where $r_D(s, a, s') = 0$, and $Q_D(s', a') = 0$, $\forall a'$, and
4. a non-terminal, non dead-end state, where $r_D(s, a, s') = 0$.

Item 3 vanishes and items 1 and 2 result in the first and the last terms in L2.1. For the item 4 above, assume any treatment a' at the state s' and consider all the possible roll-outs starting from (s', a') under execution of the greedy policy w.r.t. Q_D^* (which maximally avoids future mortality). At the end of each roll-out, the roll-out trajectory necessarily either reaches death with the \mathcal{M}_D return of -1 for the trajectory, or it reaches recovery with the \mathcal{M}_D return of 0 for the trajectory. Hence, the expected return from (s', a') will be -1 times the sum of probabilities of all the roll-outs that reach death (plus zero times sum of the rest). That is, $Q_D^*(s', a')$ is the *negative total probability of future death* from (s', a') if optimal treatments (w.r.t. Q_D^*) are always known and administered afterwards. Consequently, $\max_{a'} Q_D^*(s', a')$ would be *negative minimum probability of future death* from state s' , again if optimal treatments are known and administered at s' and afterwards. Therefore, $\sum_{s'} T(s, a, s') \max_{a'} Q_D^*(s', a')$ is negative minimum probability of future death from (s, a) under optimal policy, which by definition is $-M_D(s, a)$. This shapes the middle term of L2.1 and concludes the proof.

The second part follows a similar argument. In particular, $Q_R^*(s', a')$ is the probability of reaching recovery under the execution of greedy policy w.r.t. Q_R^* (which itself maximizes reaching a recovery terminal). Therefore, $\max_{a'} Q_R^*(s', a')$ is the maximum probability of reaching recovery under optimal policy from s' , and finally $\sum_{s'} T(s, a, s') \max_{a'} Q_R^*(s', a')$ induces maximum probability of reaching recovery from (s, a) . ■

Lemma 3.

L3.1. State s is a dead-end if and only if $P_D(s, a) + F_D(s, a) = 1$ for *all* treatments a .

L3.2. State s is a rescue if and only if $P_R(s, a) + F_R(s, a) = 1$ for *at least one* treatment a .

Proof. For part one, we note that $P_D(s, a)$, $M_D(s, a)$, and $F_D(s, a)$ are parts of the transition probability to the next state, hence

$$P_D(s, a) + M_D(s, a) + F_D(s, a) \leq 1$$

Therefore, $P_D(s, a) + F_D(s, a) = 1$ deduces $P_D(s, a) + M_D(s, a) + F_D(s, a) = 1$ (i.e., $M_D(s, a) = 0$). Invoking L2.1 induces $Q_D^*(s, a) = -1$ for all a ; hence, s is a dead-end due to L1.1. Conversely, if s is a dead-end, L1.1 induces that $Q_D^*(s, a) = -1$ for *all* treatments a . Invoking (4) again, it follows that the next state cannot be a recovery terminal state or a non-terminal state with $\max_{a'} Q_D^*(s', a') > -1$, which implies the next state is either a dead-end or a death terminal state. Hence, $P_D(s, a) + F_D(s, a) = 1$ for all treatments a .

Similar proof holds for L3.2. Note that the counterpart of (4) for this case is as the following:

$$1 = 1 - \left(\sum_{s' \in C'_D} T(s, a, s') + \sum_{s' \in C'_N} T(s, a, s') \left[1 - \max_{a'} Q_R^*(s', a') \right] \right) \quad (6)$$

with C'_D and C'_N denoting, respectively, the sets of death terminal states and non-terminal states with $\max_{a'} Q_R^*(s', a') < 1$. Similarly, (6) necessitates $T(s, a, \cdot)$ must be zero for all transitions to C'_D and C'_N . ■

Theorem 1. The followings hold:

- T.1 $P_D(s, a) + F_D(s, a) = 1$ if and only if $Q_D^*(s, a) = -1$.
- T.2 $P_R(s, a) + F_R(s, a) = 1$ if and only if $Q_R^*(s, a) = 1$.
- T.3 There exists a threshold $\delta_D \in (-1, 0)$ independent of states and treatments, such that $Q_D^*(s, a) \geq \delta_D$ for all s and a , unless if and only if $P_D(s, a) + F_D(s, a) = 1$.
- T.4 There exists a threshold $\delta_R \in (0, 1)$ independent of states and treatments, such that $Q_R^*(s, a) \leq \delta_R$ for all s and a , unless if and only if $P_R(s, a) + F_R(s, a) = 1$.
- T.5 For any policy π , state s , and treatment a , if $\pi(s, a) \leq 1 + Q_D^*(s, a)$ and $\lambda \in [0, 1]$ exists such that $P_D(s, a) + F_D(s, a) \geq \lambda$, then $\pi(s, a) \leq 1 - \lambda$.
- T.6 For any policy π , state s , and treatment a , if $\pi(s, a) \geq Q_R^*(s, a)$ and $\lambda \in [0, 1]$ exists such that $P_R(s, a) + F_R(s, a) \geq \lambda$, then $\pi(s, a) \geq \lambda$.

Proof. (T.1) and (T.2) are immediate from Lemma 1 and 3. For (T.3), it follows from (L1.1) that for a non-dead-end state s , we have $Q_D^*(s, a) > -1$. We choose $\Delta_D = \max_{s, a} [P_D(s, a) + M_D(s, a) + F_D(s, a)]$ for all non-dead-end and non-terminal states s and all treatments a . If all the transition probabilities are stationary (or more generically, $\exists \lambda < 1 : T(s, a, s') < \lambda$ for all non-dead-end and non-terminal transitions) then Δ_D is a fixed value even though it might be very close to -1 in principle. As a result, it follows from L2.1 that for any threshold $\delta_D \in (-1, -\Delta_D]$ we have $Q_D^*(s, a) \geq -\Delta_D$ unless s is a dead-end for which $Q_D^*(s, a) = -1$ due to L1.1. Furthermore, Δ_D only depends on the transition probabilities $T(s, a, s')$ and not the length of dead-ends. Similar proof concludes (T4).

In order to prove (T.5) and (T.6), we note that both $M_D(\cdot, \cdot)$ and $M_R(\cdot, \cdot)$ are non-negative for all state-treatments. Using the antecedent of (T.5), $P_D(s, a) + F_D(s, a) \geq \lambda$, as well as invoking Lemma 2, it yields:

$$\begin{aligned} Q_D^*(s, a) &\leq Q_D^*(s, a) + M_D(s, a) \\ &= -(P_D(s, a) + F_D(s, a)) \leq -\lambda \end{aligned}$$

which implies $1 + Q_D^*(s, a) \leq 1 - \lambda$. Hence, setting $\pi(s, a) \leq 1 + Q_D^*(s, a)$ deduces $\pi(s, a) \leq 1 - \lambda$.

Similarly, for (T.6) we have $P_R(s, a) + F_R(s, a) \geq \lambda$, therefore

$$\begin{aligned} Q_R^*(s, a) &\geq Q_R^*(s, a) - M_R(s, a) \\ &= P_R(s, a) + F_R(s, a) \geq \lambda \end{aligned}$$

As a result, $\pi(s, a) \geq Q_R^*(s, a)$ deduces $\pi(s, a) \geq \lambda$. ■

Proposition 1. Let $Q_D(s, a)$ be an approximation of $Q_D^*(s, a)$, such that

1. $Q_D(s, a) = Q_D^*(s, a) = -1$ for all $s \in \mathcal{S}_D$.
2. For all other states, the values satisfy monotonicity with respect to the Bellman operator \mathcal{T}^* , i.e. $Q_D(s, a) \leq (\mathcal{T}^*Q_D)(s, a)$ for all (s, a) .
3. All values of $Q_D(s, a)$ remain non-positive.

The security condition still holds if $\pi(s, a) \leq 1 + Q_D(s, a)$.

Proof. Using assumptions 1 and 2 we write

$$Q_D(s, a) \leq (\mathcal{T}^*Q_D)(s, a) \tag{7}$$

$$\begin{aligned} &= \sum_{s'} T(s, a, s') \left[r_D(s, a, s') + \max_{a'} Q_D(s', a') \right] \\ &= - \sum_{s' \in \mathcal{S}_D} T(s, a, s') - \sum_{s' \in \mathcal{C}'_D} T(s, a, s') + \sum_{s' \notin \mathcal{S}_D \cup \mathcal{C}'_D} T(s, a, s') \left[r_D(s, a, s') + \max_{a'} Q_D(s', a') \right] \end{aligned} \tag{8}$$

$$= -P_D(s, a) - F_D(s, a) - \beta_D(s, a) \tag{9}$$

in which, $-\beta_D$ is the last term of (8). The reward of \mathcal{M}_D is always zero unless at death terminal states where $r_D(s, a, s') = -1$. Hence, assumption 3 implies that $\beta_D(s, a)$ is always non-negative, regardless of how much $Q_D(s', a')$ is inaccurate. The rest of argument in Theorem 1 remains valid with Q_D and β_D replacing Q_D^* and M_D . ■

Remark 1. One setting that holds assumption 3 of Proposition 1 is in the tabular case where each $Q_D(s, a)$ is stored separately and under the assumption that all (s, a) pairs are initialized with any non-positive number (naturally in $[-1, 0]$). In the general case involving non-tabular estimators, a practical way to assure that Assumption 3 of Proposition 1 holds is to clip all the values at -1 and 0 .

Remark 2. There are certain cases that formally satisfy assumption 2. For example, the true value of *any* policy (not necessarily optimal) satisfies this inequality [11]. Another example is in the tabular setting when all values are initialized *pessimistically* (e.g., at -1); however, pessimistic initialization may increase false positives because all unseen (s, a) pairs will be inferred as dead-ends. In other cases, since $Q_D(s, a)$ is the convergence point of Bellman error, it is likely that for many state-treatment pairs assumption 2 holds. Nevertheless, one should note that this assumption needs further scrutiny and may not hold in general when function approximation is used. In particular, over-estimation issue (if exists for any state-treatment pair) will forfeit assumption 2.

Remark 3. Proposition 1 implies that under certain assumptions, at each state only the value of treatments that lead to dead-end states w.p.1 has to be fully converged. Importantly, such values are independent of the values of other (non-dead-end) states, since according to Lemma 3 a dead-end's next state is also always either a dead-end or a death terminal state, regardless of the administered treatment. In an abstract way, it leaves out the necessity of learning the value for all the resulting trajectories from other treatments at the initial state as well as in the future resulting trajectories, which grow exponentially. Hence, at least in the tabular case with -1 value-initialization, learning the treatment avoidance method by securing the behavioral policy is an exponentially smaller problem than learning optimal policy (or optimal values), which advises for best treatments.

Remark 4. Full convergence of values of dead-end states \mathcal{S}_D to -1 in Assumption 1 can be relaxed to $-(1 - \epsilon)$ for some $\epsilon \in [0, 1)$. In that case, rewriting (9) induces that the security guarantee will degrade to $\pi(s, a) \leq 1 - (1 - \epsilon)\lambda$. That is, for a risky treatment, abiding by $1 + Q_D$ guarantees less decrease of its probability than what the security conditions requires. This may be addressed by adjusting the thresholds more conservatively.

A2 Further Remarks on Related Work

In light of discussions with reviewers during the rebuttal period, we feel the need to honor similarities and differences between our work and those introduced in Irpan et al. [46] more thoroughly than space constraints allow in the main body of this paper. While there are partial parallels in terms of grounding ideas, our theoretical development vastly diverges from Irpan et al. [46], which relies wholly on empirical exploration and is centered wholly on policy evaluation rather than the assessment of specific decisions an agent may make. We summarize key important differences as follows:

- Their concept of *feasible* is simply being non-catastrophic and is different from *rescue*, which is a state where recovery is reachable w.p.1. (i.e., there is no parallel for rescue states in their work).
- The properties of the Q function and how it formally links to the probabilities of a state being feasible or catastrophic is not derived, discussed, or used in their work.
- Their OPC metric is a proxy for evaluation/ranking learned policies. They do not use the framing to identify problematic or high-risk actions that may lead to catastrophic behavior. More accurately, there is no particular parallel for the concept of (treatment) security, its definition, and the formal guarantees which then shape the foundation of DeD.
- In their work, the classification component is used to identify the value of state-action pairs on a binary $\{0, 1\}$ scale. This makes negative behavior somewhat unidentifiable (they acknowledge this) from intermediate feasible states that do not correspond to terminal conditions.
- Our dead-end construction (reward of -1 for bad outcomes + no-discounting) provides an inherently different value function, which (with a negative sign) formally gives rise to the minimum probability of bad outcomes in the future.
- Side note: in dangerous and stochastic environments and for sufficiently long episodes, their Theorem 1 results in the trivial bound (since the lower-bound becomes a negative value). Their experiments are restricted to robotic tasks and the Atari game of Pong; thus, this core problem has remained hidden in their work.

At a high level both Irpan et al. [46] and our work exploit constructed asymmetries within the state space to identify regions that are undesirable and should be avoided. The notions of *feasible* and *catastrophic* in Irpan et al. [46] are related, in context of an optimal policy π^* , with $P_{\pi^*}(\text{success}|\text{feasible}) > 0$ where $P_{\pi^*}(\text{success}|\text{catastrophic}) = 0$ always. Thus, by being able to classify which states are *catastrophic* evaluation of any trajectory containing such states is made significantly easier when evaluating policies developed from observational data. Irpan, et al. worked to label all state-action pairs as either *feasible* or *catastrophic* using positively-unlabeled classification.

With a similar asymmetry, but generalized to encompass the delicate dynamics often observed in safety-critical domains, we formalize the relationship between the special states (described in Section 3.2) and the terminal conditions of death or recovery as follows: $P(\text{recovery}|\text{rescue}) = 1$ for some policy π (including the optimal policy π^*). In contrast, dead-end states have a more extreme condition where $P(\text{death}|\text{dead-end}) = 1$ for all policies π . This helps to emphasize the importance of identifying treatments that may lead to dead-end states and subsequently influence decision-makers to avoid selecting those treatments. The means by which we infer the risk of a treatment (or action) is through a pair of independent MDPs used to identify the value of a state-treatment pair in accordance to its risk of being a dead-end or the chance it may lead to rescue and being a rescue state. This joint inference problem is used to affix and confirm whether a state should be avoided (and all treatments leading to this state) as discussed in Theorem 1 in Section 3.4.

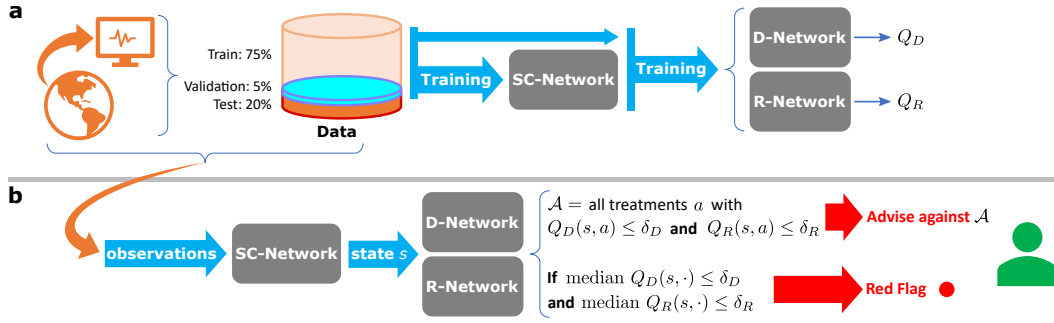


Fig. A2: **Dead-ends Discovery (DeD)**. Our pipeline includes two phases: **(a)** Training phase: using real-world data, we train the three neural networks to set-up i) state construction (SC-Network), ii) dead-end values (D-Network) and iii) rescue values (R-Network). **(b)** Test phase: the trained networks are used to map the immediate history of observations and the last action into Q_D and Q_R to infer risky conditions and dead-end outcomes, which is passed to the human decision-maker.

A4 State Construction Details and Training

This section highlights the construction and development of the state construction network used to embed the observation sequences of a patient’s health condition into a state representation to be used in the reinforcement learning networks used for the detection and avoidance of dead-end states.

A4.1 Notation

Let $\mathcal{D} = \{\tau_j\}_{j=1}^n$ denote the batch data of n trajectories obtained from the database of patients with sepsis in the intensive care unit. We assume that this data is generated from a time-homogeneous partially observable Markov decision process (POMDP). Each trajectory τ_j has a finite number of transitions m_j . Each transition in a trajectory j is a tuple with four entries $(O_{t,j}, A_{t,j}, R_{t,j}, O_{t+1,j})$, where $j \in \{1, \dots, n\}$, $t \in \{1, \dots, m_j\}$. The observation and action (treatment) spaces are defined as in Sec. A6 where:

- $O_{t,j}, O_{t+1,j} \in \mathcal{O}$ are the observations received at times t and $t + 1$ respectively in trajectory j and $\mathcal{O} \subset \mathbb{R}^{d_{\mathcal{O}}}$ is the observation space. In our case for the sepsis treatment problem, the observation space is 44 dimensional.
- $A_{t,j} \in \mathcal{A}$ is the action taken at time t in trajectory j and \mathcal{A} is the action space. In this work we restrict attention to discrete action spaces of finite cardinality, $|\mathcal{A}| = n_a$. In our case for the sepsis treatment problem, $n_a = 25$.
- $R_{t,j} \in \mathbb{R}$ is the per-step reward received at time k in trajectory j . We use an end-of-trajectory binary reward signal of ± 1 (we, however, do not explicitly make use of the reward in the state construction network because we only focus on state representation learning for dynamics prediction).

For clarity we drop the trajectory index j throughout the remainder of this section unless it is necessary to differentiate between trajectories. Let \hat{d}_S denote the dimension of the learned state representation (\hat{S}), which is a hyper-parameter that needs to be chosen. Our objective is to learn a state construction function $\psi : \{O_{0:t}, A_{0:t-1}\} \mapsto \hat{S}_t$, $t \geq 1$, and $\hat{S}_t \in \hat{S} \subset \mathbb{R}^{\hat{d}_S}$. In addition to ψ , the approaches outlined in the next section also involve another function: a dynamics predictor ϕ that involves predicting the next observation \hat{O}_{t+1} . Hence, the function $\phi : \hat{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{O})$, where $\Delta(x)$ denotes a probability distribution of x , estimates the conditional distribution of the next observation given the current state representation and action.

A4.2 State Construction (SC) Network

We construct the state representation of a patient’s condition by training a set of coupled functions, as motivated by the Approximate Information State (AIS) approach [50]. AIS satisfies two key properties: 1) each state is “Markovian” or sufficient for the prediction of the next state, and 2) observations are distinguishable when mapped

to their corresponding states if they result in different future trajectories. The first function, denoted by ψ , encodes the observed sequence patient conditions and the treatments administered into a compressed representation. This representation (corresponding to the state used in the reinforcement learning networks) is then passed, along with the current treatment, to a decoding function ϕ to predict the next patient observation.

The input to ψ is the concatenation of the observation O_t and last selected action A_{t-1} . For the function ψ we use a 3-layer Recurrent Neural Network (RNN), where the first layer is a fully connected layer that maps the current observation and action (69 dimensional input: 44 dimensional observation with a 25 dimensional one-hot encoded action) to 64 units with ReLU activation. This is followed by another (64, 128) fully connected layer with ReLU activation which is followed by a gated recurrent unit [57] layer with hidden state size \hat{d}_S . The output of this recurrent layer is used as the state representation \hat{S}_t . The current action A_t is concatenated to the state representation \hat{S}_t and then fed through the decoder function ϕ to predict the next observation \hat{O}_{t+1} . The function ϕ is comprised of a three layer neural network with sizes $(\hat{d}_S + 25, 64)$, $(64, 128)$ and $(128, 44)$ (with ReLU activation for the first two layers). The last layer outputs a 44-dimensional vector, which forms the mean vector of a unit-variance multivariate Gaussian distribution, samples from which are used to predict the next observation. A schematic of the the state construction network is provided in Fig. A3. The two functions ψ and ϕ that comprise the state construction network are jointly trained by maximizing the negative log likelihood of the predicted next observation \hat{O}_{t+1} .

This is formulated by maximizing the objective:

$$\mathcal{L}(\mathcal{O}_{t+1}, \hat{\mathcal{O}}_{t+1}) = - \sum^{d_{\mathcal{O}_j}} \log \mathcal{N}(\mathcal{O}_{t+1,j}; \mu_j, \sigma_j^2)$$

where $\mu_j = \hat{\mathcal{O}}_{t+1}$, $\sigma_j^2 = 1$, and $\hat{\mathcal{O}}_{t+1} = \psi(\phi(O_t, A_{t-1}), A_t)$.

A4.3 Hyperparameter selection

The dimension of the state representation \hat{d}_S was chosen from among $\{4, 8, 16, 32, 64, 128, 256\}$ dimensions. The choices of the size of neural network layers was chosen proportional to the size of \hat{d}_S , with the final values reported in the prior subsection following the optimal choice of \hat{d}_S being equal to 64. The model construction network was trained for 600 epochs with learning rates of $\{0.0001, 0.0005, 0.001, 0.005, 0.01\}$ with the choice of $lr = 0.0005$ providing the optimal training of the network. We demonstrate the evaluation of the choice of the dimension for the state representation in Fig. A4.

A5 D- and R-Networks Training Details

We use double-DQN algorithm to train both networks. We refer the reader to our code for the implementation details (and we tried to make the code straightforward and relatively easy to understand). In particular, both D- and R-Networks consist of two linear layers with 64 nodes. The first layer is followed by ReLU nonlinearity and the second layer directly outputs 25 nodes corresponding to the 25 treatments. We use learning rate of 0.0001 and minibatch size of 64. In each minibatch, we select 62 transitions uniformly from the train data and append it with two uniformly selected “death” transitions (last transitions of nonsurvivor patients). All other chosen hyper-parameters can be found in the *config.yaml* file in the root directory of our code.

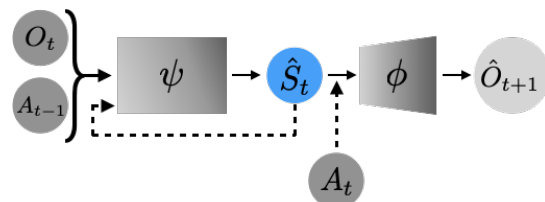


Fig. A3: The state construction network, comprised of the encoding function ψ that provides the state representation \hat{S}_t that is used with the decoding function ϕ to predict the next observation \hat{O}_{t+1} .

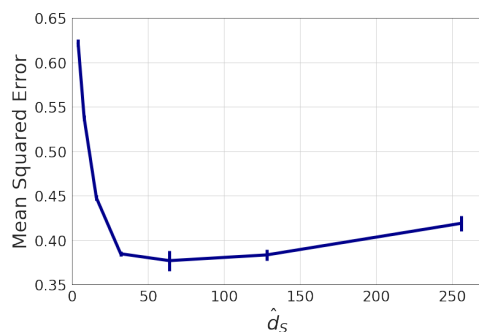


Fig. A4: Analysis of setting the dimension of the learned state representation \hat{d}_S and its effect on the accuracy of predicting the next observation. The bars represent standard deviation. With this, we determine to set $\hat{d}_S = 64$ in the SC-network.

A6 Data Details

We use the MIMIC (Medical Information Mart for Intensive Care) - III dataset (v1.4), which has been sourced from the Beth Israel Deaconess Medical Center in Boston, Massachusetts [22, 48]. This dataset comprises of deidentified patient treatment records of patients admitted to critical care units (CCU, CSRU, MICU, SICU, TSICU). The database includes data collected from 53,423 distinct hospital admissions of patients over 16 years of age for a period of 12 years from 2001 to 2012. The MIMIC dataset has been used in many reinforcement learning for health care projects, including mechanical ventilation and sepsis treatment problems. There are various preprocessing steps that are performed on the MIMIC-III dataset in order to obtain the cohort of patients and their relevant observables for the sepsis treatment study.

To extract and process the data, we follow the approach described in [25] and the associated code repository given in [58]. This includes all ICU patients over 18 years of age who have some presumed onset of sepsis (following the Sepsis 3 criterion) during their initial encounter in the ICU after admission, with a duration of at least 12 hours. These criteria provide a cohort of 19,611 patients, among which there is an observed mortality rate just above 9%, where mortality is determined by patient expiration within 48h of the final observation. Observations are processed and aggregated into 4h windows with treatment decisions (administering fluids, vasopressors, or both) discretized into 5 volumetric categories. All data is normalized to zero-mean and unit variance and missing values are imputed using k-Nearest Neighbor imputation, where possible. In the absence of similar observations any remaining missing values filled with the population mean. We report the 44 features used for the Dead-end approach proposed in this paper in Table A1 with high-level statistics for the extracted cohort in Table A2.

Table A1: Patient features used for learning state representations for predicting future observations

| | | | |
|--------------------|--------------------|-------------------|-------------------|
| Age | Gender | Weight (kg) | Re-admission |
| Glasgow Coma Scale | Heart Rate | Sys. BP | Dia. BP |
| Mean BP | Respiratory Rate | Body Temp (C) | FiO2 |
| Potassium | Sodium | Chloride | Glucose |
| INR | Magnesium | Calcium | Hemoglobin |
| White Blood Cells | Platelets | PTT | PT |
| Arterial pH | Lactate | PaO2 | PaCO2 |
| PaO2 / FiO2 | Bicarbonate (HCO3) | SpO2 | BUN |
| Creatinine | SGOT | SGPT | Total Bilirubin |
| Output (4h) | Output (total) | Cumulated Balance | SOFA |
| SIRS | Shock Index | Base Excess | Mech. Ventilation |

Table A2: MIMIC Sepsis Cohort Statistics

| Variable | MIMIC (<i>n</i> = 19611) | Variable | MIMIC (<i>n</i> = 19611) |
|--|---------------------------|--|---------------------------|
| Demographics | | Outcomes | |
| Age, years | 66.2 (53.8-78.1) | Deceased | 1881 (9.6%) |
| Age range, years | | Vasopressors administered | 5664 (28.9%) |
| 18-29 | 741 (3.8%) | Fluids administered | 17812 (90.8%) |
| 30-39 | 896 (4.6%) | Ventilator used | 9353 (47.7%) |
| 40-49 | 2029 (10.3%) | | |
| 50-59 | 3471 (17.7%) | Severity Scores | |
| 60-69 | 4321 (22.0%) | SOFA | 5 (3.0-8.0) |
| 70-79 | 4086 (20.8%) | SIRS | 2 (1.0-2.0) |
| 80-89 | 3069 (15.6%) | Shock Index | 0.72 (0.6-0.86) |
| ≥90 | 998 (5.1%) | | |
| Gender | | | |
| Male | 10917 (55.6%) | | |
| Female | 8694 (44.3%) | | |
| Re-admissions | 1424 (7.3%) | | |
| Physical exam findings | | | |
| Temperature (°C) | 37.2 (36.6-37.7) | | |
| Weight (kg) | 79.7 (66.7-95.2) | | |
| Heart rate (beats per minute) | 86.0 (75.0-98.0) | | |
| Respiratory rate (breaths per minute) | 19.8 (16.6-23.3) | | |
| Systolic blood pressure (mmHg) | 118.3 (105.8-133.6) | | |
| Diastolic blood pressure (mmHg) | 56.6 (48.6-65.4) | | |
| Mean arterial pressure (mmHg) | 77.0 (69.0-86.7) | | |
| Fraction of inspired oxygen (%) | 40.0 (35.0-50.0) | | |
| P/F ratio | 307.5 (192.0-579.0) | | |
| Glasgow Coma Scale | 14.8 (11.0-15.0) | | |
| Laboratory findings | | | |
| Hematology | | Coagulation | |
| White blood cells (thousands/ μ L) | 10.8 (7.7-14.8) | Prothrombin time (sec) | 14.3 (13.1-16.4) |
| Platelets (thousands/ μ L) | 202.0 (137.0-286.0) | Partial thromboplastin time (sec) | 32.6 (27.6-44.9) |
| Hemoglobin (mg/dL) | 10.2 (9.1-11.4) | INR | 1.3 (1.1-1.5) |
| Base Excess (mmol/L) | 0.5 (0.0-2.6) | | |
| Chemistry | | Blood gas | |
| Sodium (mmol/L) | 138.9 (136.0-141.0) | pH | 7.41 (7.35-7.44) |
| Potassium (mmol/L) | 4.0 (3.7-4.4) | Oxygen saturation (%) | 97.3 (95.5-98.8) |
| Chloride (mmol/L) | 105.0 (101.0-108.5) | Partial pressure of O ₂ (mmHg) | 124.0 (85.0-241.1) |
| Bicarbonate (mmol/L) | 25.0 (22.0-28.0) | Partial pressure of CO ₂ (mmHg) | 40.6 (36.0-46.0) |
| Calcium (mg/L) | 8.3 (7.8-8.8) | | |
| Magnesium (mg/L) | 2.0 (1.8-2.2) | | |
| Blood urea nitrogen (mg/dL) | 22.0 (14.0-36.0) | | |
| Creatinine (mg/dL) | 1.0 (0.7-1.5) | | |
| Glucose (mg/dL) | 127.4 (107.0-156.0) | | |
| SGOT (units/L) | 38.0 (22.0-74.0) | | |
| SGPT (units/L) | 30.0 (17.0-64.0) | | |
| Lactate (mg/L) | 1.5 (1.1-2.2) | | |
| Total bilirubin (mg/L) | 0.7 (0.4-1.5) | | |

A7 Supporting Figures and Tables

a Red flag thresholds

| | D-Network | | | | R-Network | | | | Full | | | |
|--------------|-----------|-------|--------------|-------|-----------|-------|--------------|-------|-----------|------|--------------|-------|
| | Survivors | | Nonsurvivors | | Survivors | | Nonsurvivors | | Survivors | | Nonsurvivors | |
| | Q_D | V_D | Q_D | V_D | Q_R | V_R | Q_R | V_R | Q | V | Q | V |
| -72 h | 0.2% | 0.2% | 0.0% | 0.0% | 0.5% | 0.2% | 2.8% | 0.9% | 0.0% | 0.2% | 0.0% | 0.0% |
| -48 h | 1.2% | 0.4% | 8.1% | 5.4% | 1.5% | 0.5% | 5.9% | 4.3% | 0.7% | 0.2% | 2.7% | 2.7% |
| -24 h | 1.1% | 0.4% | 16.3% | 13.0% | 1.2% | 0.3% | 16.7% | 13.0% | 0.6% | 0.1% | 12.2% | 10.6% |
| -12 h | 0.9% | 0.4% | 20.2% | 18.2% | 0.7% | 0.3% | 20.2% | 17.4% | 0.4% | 0.2% | 12.8% | 14.7% |
| -8 h | 1.0% | 0.4% | 24.5% | 21.9% | 0.9% | 0.3% | 19.3% | 20.4% | 0.6% | 0.2% | 13.4% | 17.8% |
| -4 h | 1.2% | 0.5% | 29.7% | 26.4% | 0.7% | 0.5% | 24.9% | 22.7% | 0.5% | 0.3% | 20.1% | 22.0% |

a Yellow flag thresholds

| | D-Network | | | | R-Network | | | | Full | | | |
|--------------|-----------|-------|--------------|-------|-----------|-------|--------------|-------|-----------|------|--------------|-------|
| | Survivors | | Nonsurvivors | | Survivors | | Nonsurvivors | | Survivors | | Nonsurvivors | |
| | Q_D | V_D | Q_D | V_D | Q_R | V_R | Q_R | V_R | Q | V | Q | V |
| -72 h | 1.6% | 0.5% | 4.6% | 2.8% | 1.8% | 0.2% | 5.6% | 2.8% | 0.5% | 0.0% | 2.8% | 2.8% |
| -48 h | 3.1% | 2.2% | 12.4% | 11.9% | 2.7% | 2.1% | 14.1% | 11.9% | 1.6% | 1.5% | 10.3% | 9.7% |
| -24 h | 2.7% | 1.8% | 17.1% | 20.3% | 2.2% | 1.8% | 13.8% | 15.9% | 1.4% | 1.4% | 10.6% | 13.4% |
| -12 h | 3.3% | 2.5% | 19.4% | 19.4% | 3.4% | 2.4% | 17.4% | 17.8% | 2.0% | 1.7% | 15.1% | 17.1% |
| -8 h | 3.0% | 2.1% | 20.8% | 21.9% | 2.6% | 2.1% | 21.6% | 17.8% | 1.2% | 1.5% | 18.6% | 15.6% |
| -4 h | 3.2% | 2.4% | 20.1% | 20.1% | 3.0% | 2.4% | 16.8% | 21.2% | 1.7% | 1.5% | 16.1% | 17.6% |

Table A3: **Prediction of potentially life-threatening treatments and states (full list)**. Similarly to 2, the results correspond to the part of test data that satisfies having minimum length of the corresponding time step (X hours before terminal). To raise a flag, a patient must concurrently violate the corresponding thresholds, as specified in 2. Q columns correspond to the value of actually selected treatments, while V columns correspond to the median value of patients' state at the corresponding time.

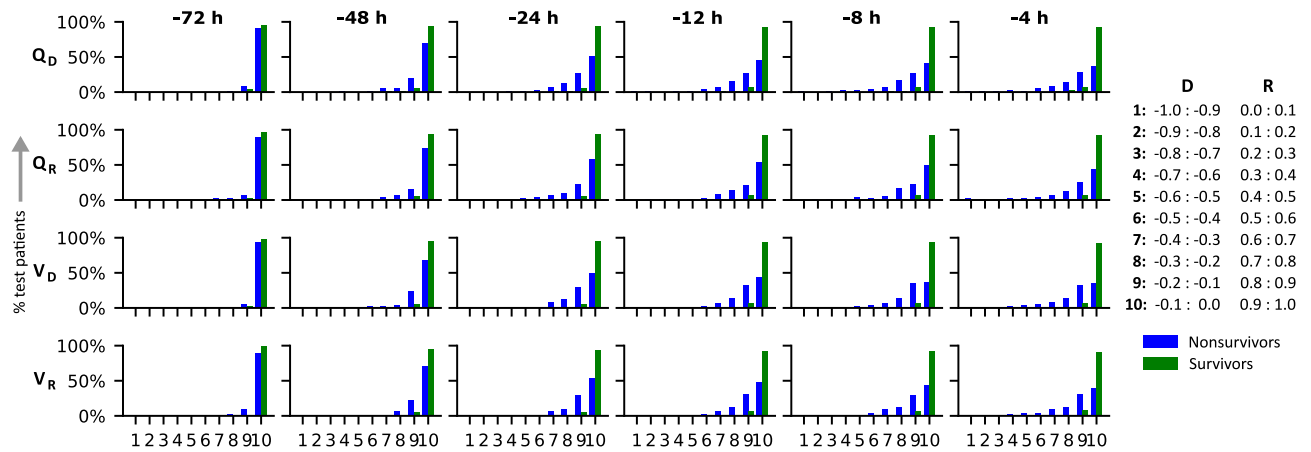


Fig. A5: **Full histogram of values in different time steps.** The histograms are plotted from the part of test data that satisfies having minimum length of the time step. The four rows are corresponding to the following: V_D and V_R : median value of states from D-Network and R-Network, respectively, and Q_D and Q_R : value of the selected treatments at the given time step from D-Network and R-Network, respectively. Note the distinctive difference between the trend of values in survivor (green bars) and nonsurvivor (navy bars) trajectories. In particular, in the course of 72 hours in the ICU, there is not much change in the value of selected or median treatment for the survivor patients, which is completely in contrast with those of nonsurvivor patients.

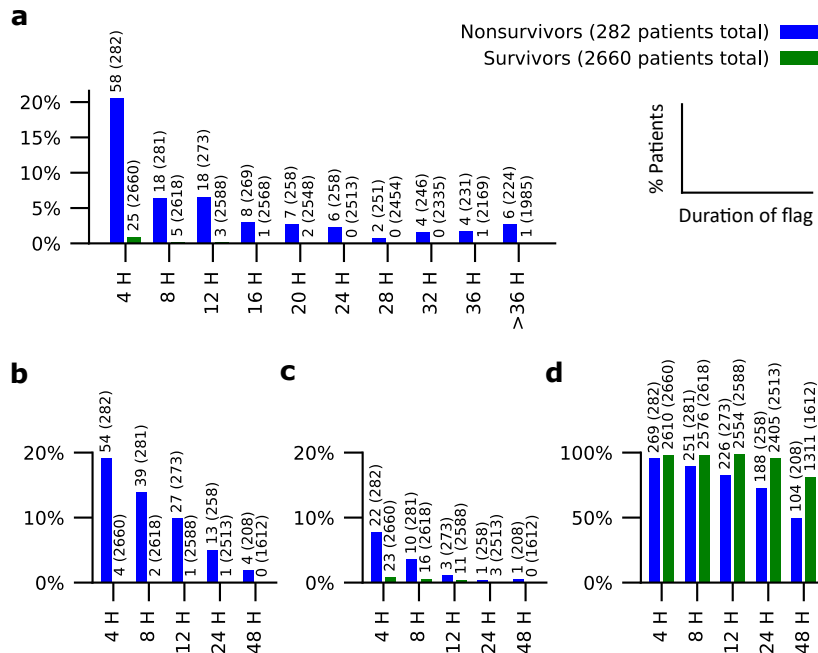


Fig. A6: **Flag duration for ICU patients.** Remaining on *confirmed red-flag* is measured for both survivor and nonsurvivor patients. **a** The bars represent the percentage of patients who experience at least one red-flag with the exact duration on the horizontal axis. Texts depict number of patients (out of total patients with the minimum of specified stay duration). **b** and **c** depict patients who “finish” their ICU stay remaining on red and yellow flags, respectively, at the final X hours before terminal. **d** presents patients who “start” their trajectory with *no flag* at all for the first X hours on the horizontal axis. We found that for the large part, both survivors and nonsurvivors start their trajectory without any flag, suggesting that they do not necessarily start with an unrecoverable situation. Further, nearly zero percent of survivors would raise and remain on red-flag for more than eight hours (even eight hours is quite rare compared to the total number of survivor patients). In contrast, nonsurvivor patients demonstrate a fat tail in the duration distribution **a** and repeatedly remain on the red-flag for eight hours or more. This result suggests that remaining on the red-flag for long periods strongly correlates with mortality, which is inline with our theoretical analysis.

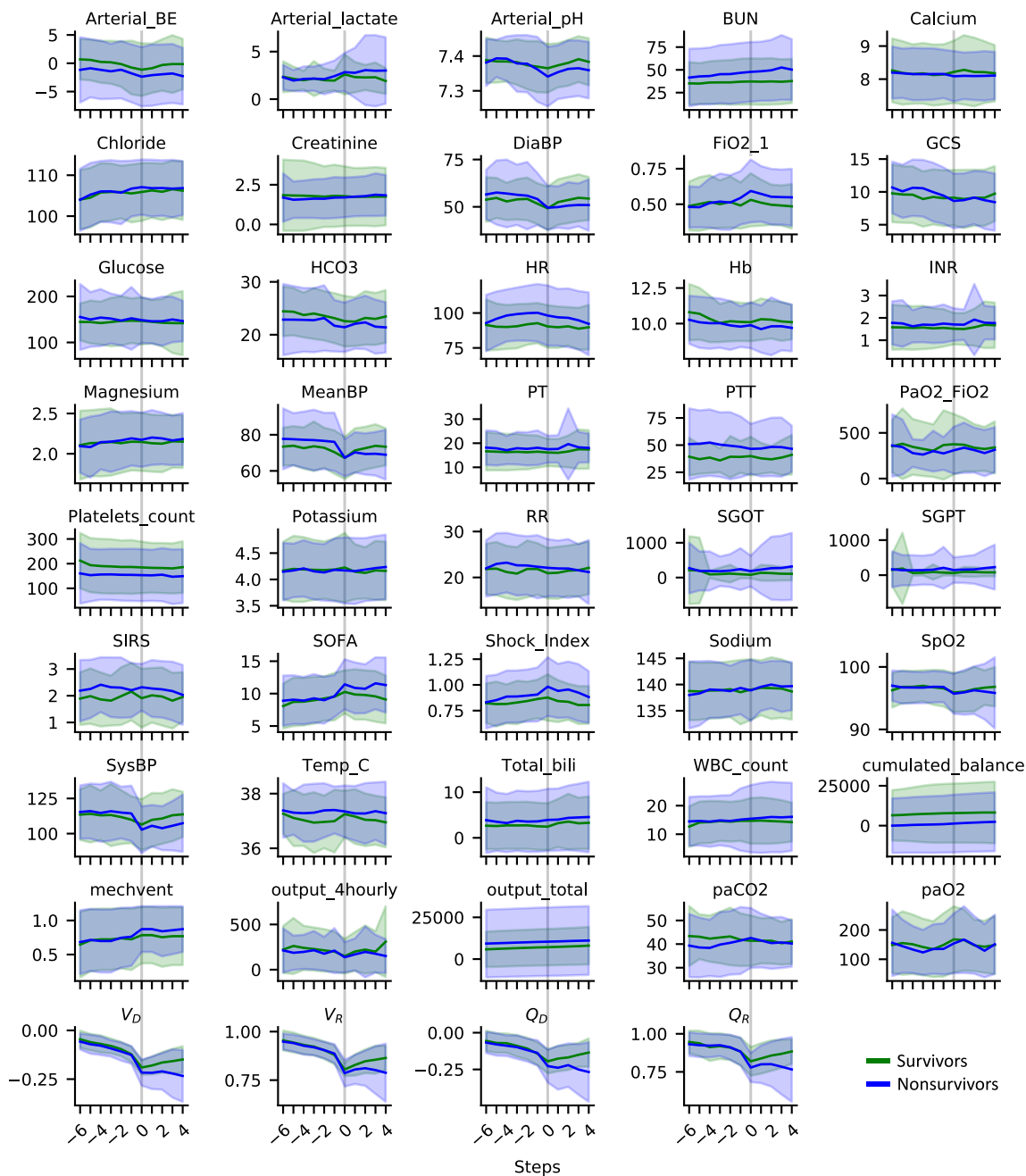


Fig. A7: **Signals prior to the first flag.** Complete list of vitals and standard measures in addition to our dead-end and secure values are shown for both survivor and nonsurvivor patients 24 hours (6 steps, 4 hours each) before and 16 hours (4 steps) after the first raised flag (red or yellow), indicated at point zero. Shaded areas represent standard deviation.

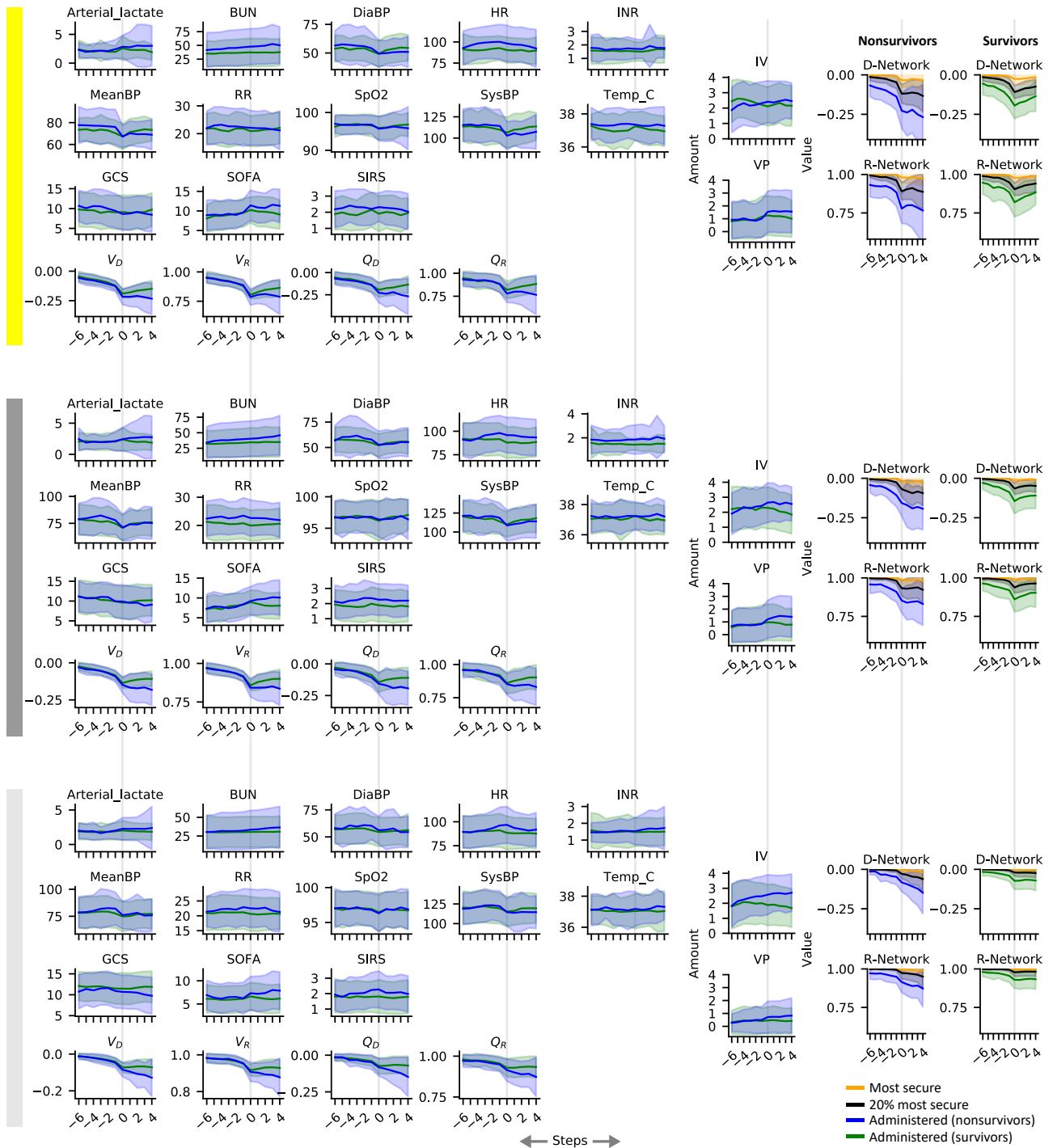
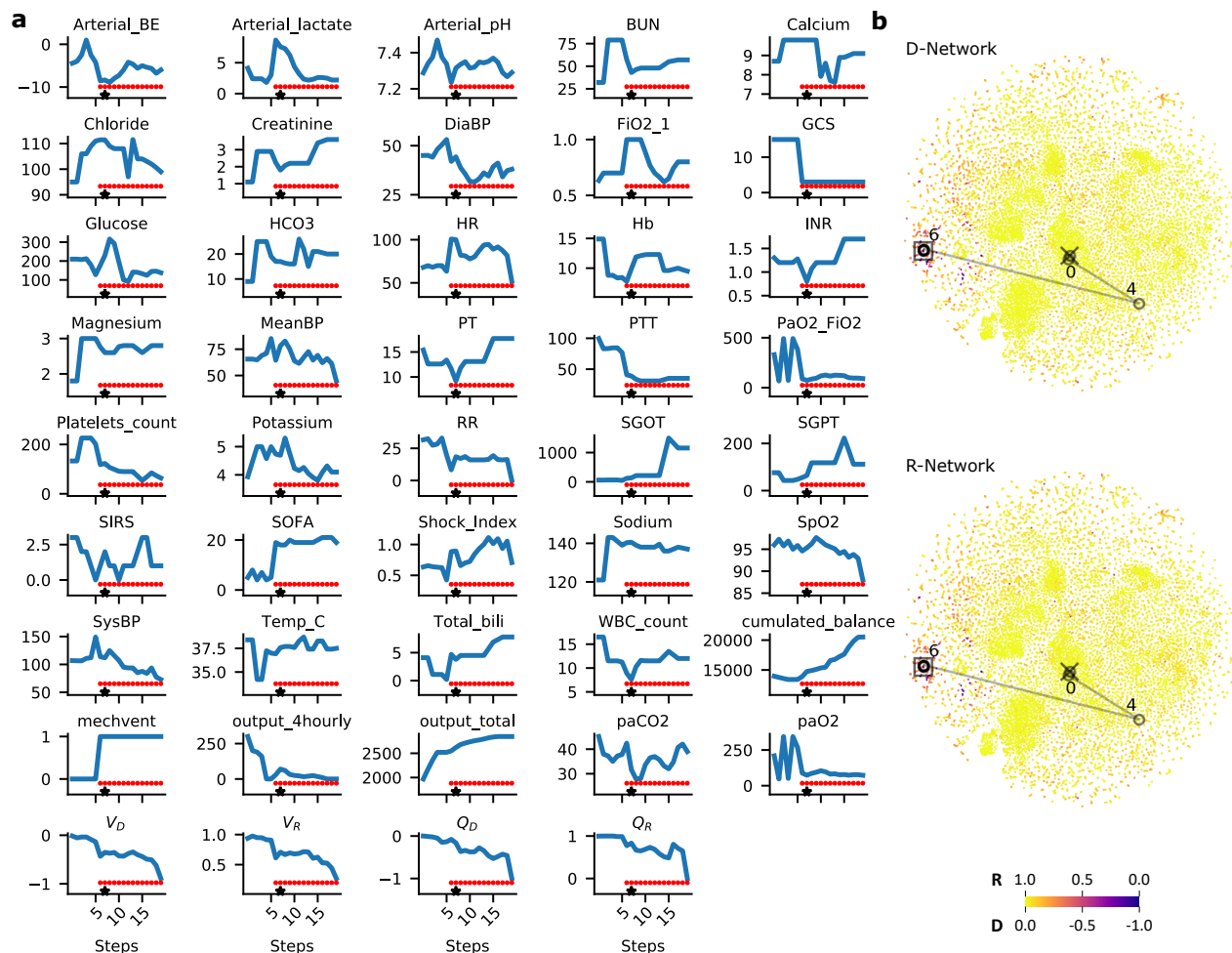


Fig. A8: **Trend of various measures before and after the first raised flags.** Various measures are shown 24 hours (6 steps) before and 16 hours (4 steps) after the first threshold crossing. The colors respectively corresponds to the following thresholds: yellow: $\delta_D = -0.15$, $\delta_R = 0.85$; dark grey: $\delta_D = -0.10$, $\delta_R = 0.90$; light grey: $\delta_D = -0.05$, $\delta_R = 0.95$. Shaded areas represent standard deviation.



C

Step 0: 2181-06-16 18:49:00 (Chest X-Ray Report): "...multifocal pneumonia, asymmetric pulmonary edema, or both... Concern for fluid overload... worsening bilateral opacification..."

Step 2: 2181-06-17 03:24:00 (Nursing Report): "...[family] wish expressed that pt supported fully, including intubation if necessary, ... pt is lethargic, answers questions intermittently w/ unclear speech... continue sepsis protocol..."

Step 3: 2181-06-17 06:04:00 (Chest X-Ray Report): "...Improved aeration of the lungs with features of fluid overload and possible worsening right effusion..."

Step 4: 2181-06-17 12:23:00: Patient intubated

Step 4: 2181-06-17 13:05:00 (Chest X-Ray Report): "Some worsening of airspace findings bilaterally in the lower lung zones -- fluid overload likely -- ..."

Step 6: 2181-06-17 18:07:00 (Nursing Report): "Pt. intubated for impending resp. failure... Became hypotensive shortly after intubation and started on vasopressin... Lactate trending up... Awaiting brother's visit tonight to ? make cmo..."

Step 8: 2181-06-18 03:37:00 (Nursing Report): "Awaiting arrival of brother and continuing w/ full aggressive treatments until his arrival..."

Step 9: 2181-06-18 05:13:00 (Nursing Report): "...plan is to continue with support..."

Step 11: 2181-06-18 16:36:00 (Nursing Report): "Pt remains unresponsive, no longer breathing over vent. Lactate has been trending down... Brother has been bedside, is leaning toward CMO, will consult w/ other family... Continue current care..."

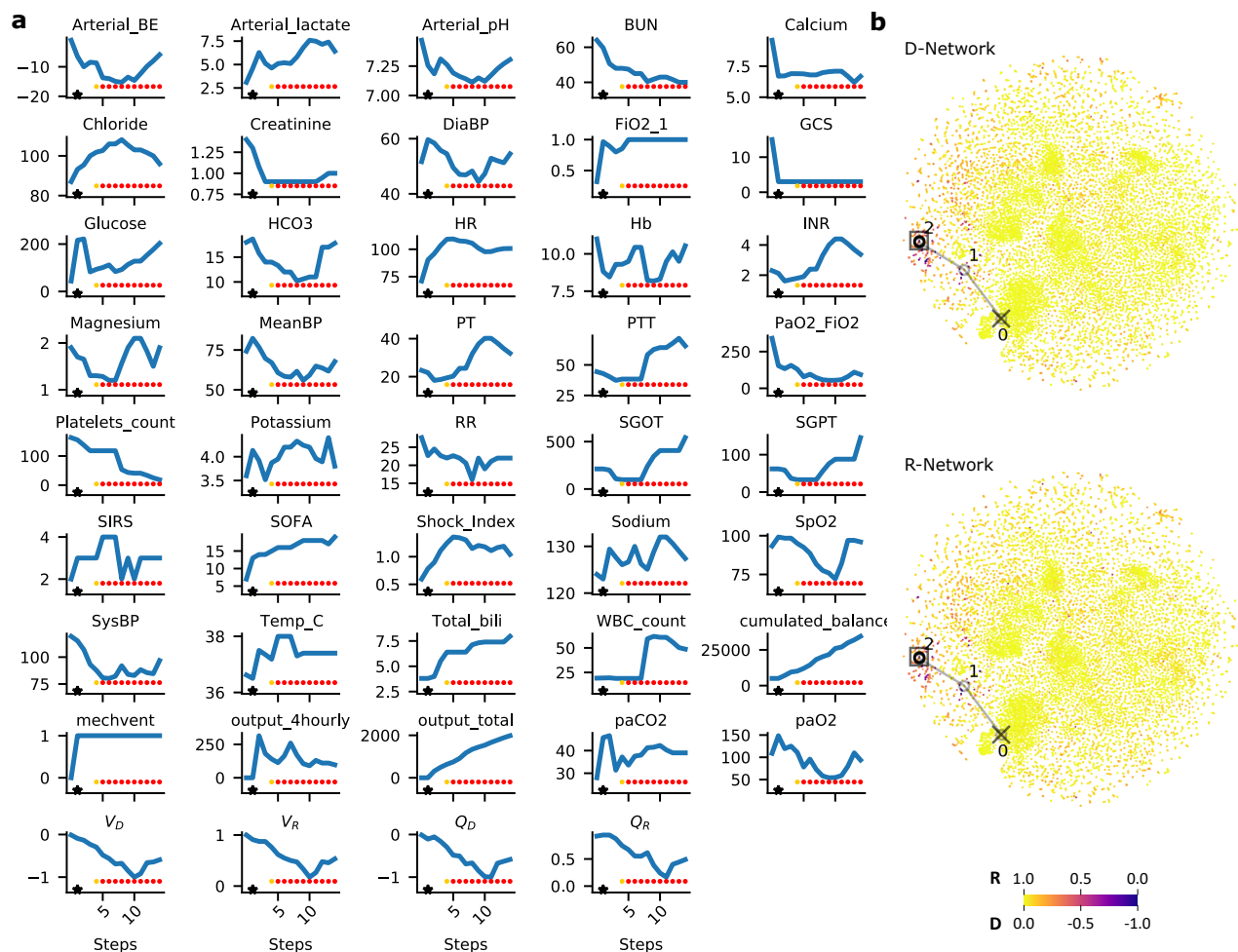
Step 12: 2181-06-18 17:22:00 (Nursing Report): "Brother arrived with sister... verbalizing wishes to withdraw life support, maintain comfort care..."

Step 15: 2181-06-19 05:53:00 (Chest X-Ray Report): "Multifocal infection including nodules in the left lower lung... Distension of the stomach with air and fluid is improved..."

Step 15: 2181-06-19 06:16:00 (Nursing Report): "Pt. continues to be non-responsive... Pt. made DNR... Continue current level of care, ? making pt. CMO if no improvement over next 24 hours..."

Step 17: 2181-06-19 16:13:00 (Nursing Report): "Pt. remains intubated... Dropped RR from 26 to 20... Family still undecided on whether to change pt's code status... Pt. remains unresponsive... Pt. is DNR, family meeting later to discuss status change..."

Fig. A9: Complete analysis of nonsurvivor patient 262011. **a** all the vitals, standard measures, max treatments, and network values for a nonsurvivor patient ICU-Stay-ID 262011. Red dots, yellow dots, and the asterisks show red and yellow flags and the presumed onset of sepsis, respectively. **b** patient's trajectory on the t-SNE plot, and **c** extracted chart notes from different source with their corresponding time stamp and quantized step.



Step 0: 2176-04-09 18:42:00 (Abdomen Report): "Large abd wall hernia...p/w abd wall infection, likely necrotic. Also w/ diffuse TTP over abd... Visualized lung bases are clear"

Step 0: 2176-04-09 20:48:00 (Chest port line placement): "Abd hernia, now w/ hypoxia... w/ minimal linear opacity at the rt lung base"

Step 0: 2176-04-09 20:48:00 (Chest port line placement): "Atelectasis or developing infiltrate at base of rt lower lobe is less conspicuous"

Step 1: 2176-04-09 21:20:00 (Abdomen Report; post intubation): "Collapse of rt middle lobe, rt lower lobe and significant left to right cardiomeastinal shift..."

Step 2: 2176-04-10 01:21:00 (Abdominal CT): "Persistent volume loss in rt lunch w/ some expansion in rt middle lobe... Probably newly developing rt small pleural effusion"

Step 3: 2176-04-10 05:09:00 (Nursing Report): "Pt had hernia w/ necrotic abd wall... Pt has coarse bilateral LS w/ very thick brown secretions...Pt remains sedated w/ no spont. Respirations."

Step 3: 2176-04-10 05:36:00 (Nursing Report): "ABG's becoming progressively more acidotic..."

Step 4: 2176-04-10 09:22:00 (Chest CT): "Substantial clearing of opacification at the right base, consistent with re-expansion of lung following removal of mucous plug or repositioning of endotracheal tube..."

Step 5: 2176-04-10 13:46:00 (Chest CT): "Improved expansion of rt lower lobe... Minimal residual atelectasis is seen in rt middle lobe..."

Step 6: 2176-04-10 17:31:00 (Nursing Report): "Large volume resuscitation for hypotension... Persistently hypotensive... Lungs coarse, decreased at bases... DP/PT pulse present in AM, now absent... Code status changed to DNR..."

Step 8: 2176-04-11 04:28:00 (Nursing Report): "Pt on AC vent... Situation went from bad to worse..."

Step 9: 2176-04-11 06:37:00 (Nursing Report): "Pt sedated and paralyzed, doesn't appear to be in pain... Severe metabolic acidosis on max vaso and vent support... Worsening condition..."

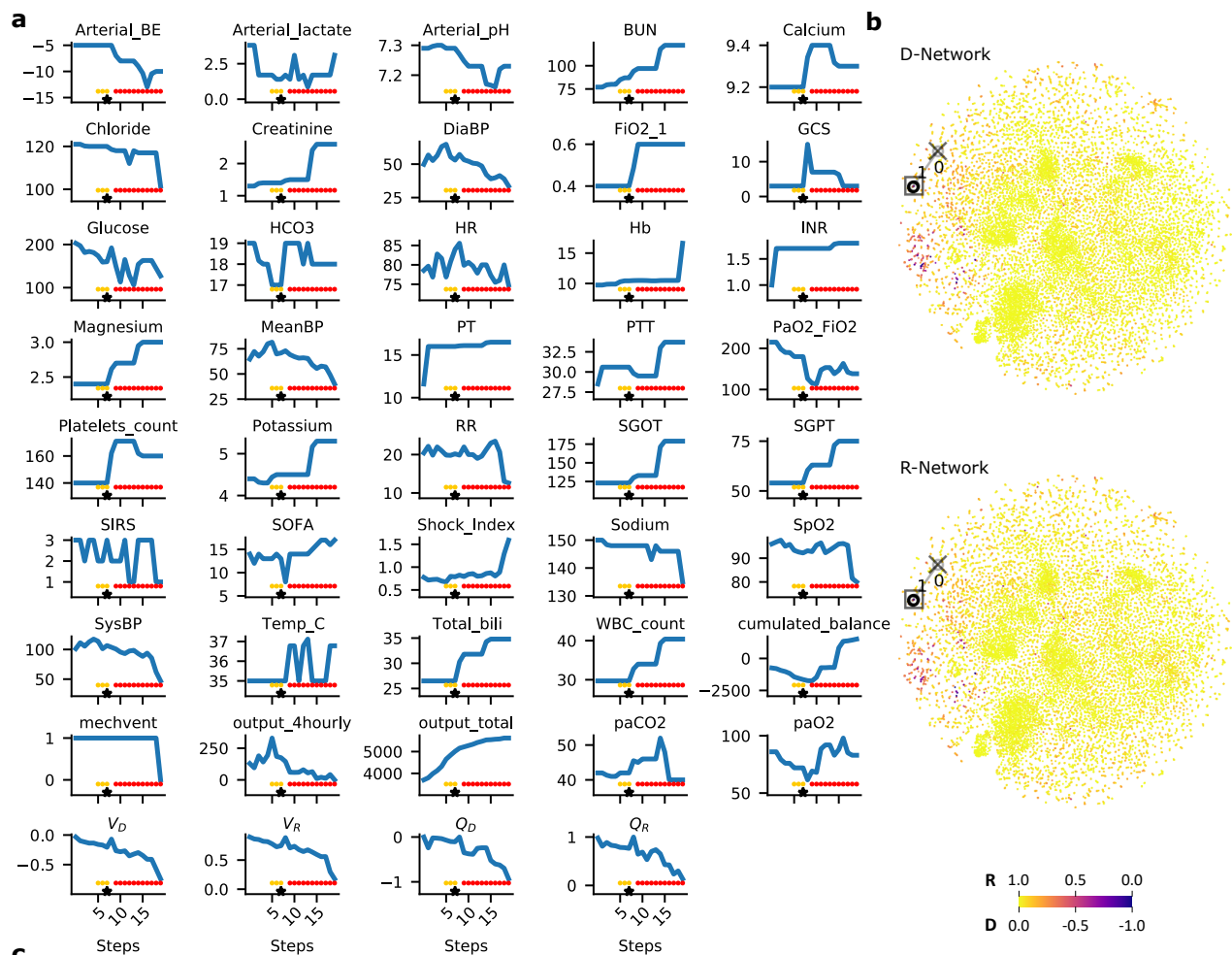
Step 11: 2176-04-11 16:19:00 (Nursing Report): "Pt. remains intubated and currently vented on full support... pt. Remains metabolically acidotic and severely hypoxic..."

Step 11: 2176-04-11 16:21:00 (Nursing Note): "Pt. very fluid positive w/ total body anasarca. Areas of necrosis remaining... Very poor prognosis and condition... Pt remains critically ill w/ profound hypoxia and met. Acidosis as well as sepsis... Continue aggressive ICU care."

Step 14: 2176-04-12 06:12:00 (Nursing report): "Pt requiring multiple fluid boluses to maintain BP... Improved metabolic acidosis... on max vaso and vent support..."

[8 hours after this report, the family requests the patient to be made CMO and expires shortly thereafter]

Fig. A10: Complete analysis of nonsurvivor patient 270174. **a** all the vitals, standard measures, max treatments, and network values for a nonsurvivor patient ICU-Stay-ID 270174. Red dots, yellow dots, and the asterisks show red and yellow flags and the presumed onset of sepsis, respectively. **b** patient's trajectory on the t-SNE plot, and **c** extracted chart notes from different source with their corresponding time stamp and quantized step.



Step 0: 2193-05-15 21:31:00 (CT Head Report): "Liver failure... Uneven but reactive pupils on phys. exam..."

Step 2: 2193-05-16 04:44:00 (Nursing Report): "Pt. unresponsive to painful stimulation, extremities flaccid..."

Step 2: 2193-05-16 05:44:00 (Chest X-Ray Report): "Hepatic failure and GI bleed... Pt. remains intubated... Marked improvement of left sided pleural effusion..."

Step 4: 2193-05-16 14:24:00 (Abdomen CT Report): "Rising amylase, investigate for necrotizing pancreatitis... cirrhotic liver... opacification in the lower left lobe... suggesting a focal infectious or inflammatory process..."

Step 5: 2193-05-16 15:25:00 (Nursing Report): "Oxygenation improved!!! on R+L side w/ sat of 95-98%... In the setting of pancreatitis and worsening LFTs... family contacted by phone, informed of status and DNR status obtained..."

Step 5: 2193-05-16 17:45:00 (Nursing Report): "Pt. essentially unresponsive on fentanyl and ativan drips... He is overbreathing the vent..."

Step 8: 2193-05-17 04:54:00 (Nursing Report): "Plan family meeting..."

Step 11: 2193-05-17 18:21:00 (Nursing Report): "There were a few vent changes made in hopes of forcing a compensation... Pt. w/o any improvements, worsening acidosis..."

Step 14: 2193-05-18 05:59:00 (Nursing Report): "No spontaneous or purposeful movement... Worsening renal failure, worsening overall system failure..."

Step 17: 2193-05-18 17:49:00 (Nursing Report): "Pt. had issues w/ hypotension today... Continues to be acidotic... Family was called this AM and was told the severity of the situation... Pt. will be extubated and made CMO..."

Fig. A11: Complete analysis of nonsurvivor patient 235403. **a** all the vitals, standard measures, max treatments, and network values for a nonsurvivor patient ICU-Stay-ID 235403. Red dots, yellow dots, and the asterisks show red and yellow flags and the presumed onset of sepsis, respectively. **b** patient's trajectory on the t-SNE plot, and **c** extracted chart notes from different source with their corresponding time stamp and quantized step.