

Towards Real-World Applications of Personalized Anesthesia Using Policy Constraint Q Learning for Propofol Infusion Control

Xiuding Cai, Jiao Chen, Yaoyao Zhu, Beimin Wang, Yu Yao

Abstract—Automated anesthesia promises to enable more precise and personalized anesthetic administration and free anesthesiologists from repetitive tasks, allowing them to focus on the most critical aspects of a patient’s surgical care. Current research has typically focused on creating simulated environments from which agents can learn. These approaches have demonstrated good experimental results, but are still far from clinical application. In this paper, Policy Constraint Q-Learning (PCQL), a data-driven reinforcement learning algorithm for solving the problem of learning strategies on real world anesthesia data, is proposed. Conservative Q-Learning was first introduced to alleviate the problem of Q function overestimation in an offline context. A policy constraint term is added to agent training to keep the policy distribution of the agent and the anesthesiologist consistent to ensure safer decisions made by the agent in anesthesia scenarios. The effectiveness of PCQL was validated by extensive experiments on a real clinical anesthesia dataset we collected. Experimental results show that PCQL is predicted to achieve higher gains than the baseline approach while maintaining good agreement with the reference dose given by the anesthesiologist, using less total dose, and being more responsive to the patient’s vital signs. In addition, the confidence intervals of the agent were investigated, which were able to cover most of the clinical decisions of the anesthesiologist. Finally, an interpretable method, SHAP, was used to analyze the contributing components of the model predictions to increase the transparency of the model.

Index Terms—Anesthesia, offline reinforcement learning (ORL), anesthetic administration, propofol.

I. INTRODUCTION

ANESTHESIA is a critical component of the operating room, with millions of patients requiring general anesthesia during surgery each year [1]. Anesthesiologists face

This work was supported in part by the National Natural Science Foundation of China under Grant 82073338, and in part by the Sichuan Provincial Science and Technology Department under Grant 2022YFS0384 and 2022YFQ0108. (Corresponding authors: Jiao Chen.)

Xiuding Cai, Yaoyao Zhu, Beiming Wang and Yu Yao are with Chengdu Institute of Computer Application, Chinese Academy of Sciences, Chengdu, China, and with the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing, China (e-mail: {caixiuding20, zhuyaoyao19, wang-beimin21}@mailsucas.ac.cn, casitmed2022@163.com).

Jiao Chen was with Department of Anesthesiology, West China Hospital, Sichuan University & The Research Units of West China (2018RU012), and with Chinese Academy of Medical Sciences, China (e-mail: chenjiao@wchscu.cn).

tremendous work pressure every day. They must monitor the patient throughout the procedure while maintaining several aspects of the patient simultaneously, including anesthesia status, physiological stability, pain management, and oxygen delivery. Routinely, an anesthesiologist is responsible for the life support of more than one patient at a time. However, the increasing number of operations in recent years has posed a significant challenge to the limited number of anesthesiologists [2]. This serious challenge not only puts immense pressure on anesthesiologists, but also increases the risk of medical malpractice and burnout [3]. Given the critical role of anesthesia, the repetitive nature of anesthesiologists’ tasks, and the shortage of positions, automated anesthesia infusion has emerged as an essential research direction in healthcare, offering a promising solution to alleviate the challenge. Automatic anesthetic administration has been shown to allow for more accurate and responsive anesthetic drug control with a reduced total dose than complete manual control by the anesthesiologist [4]–[7]. This approach also helps reduce the patient’s postoperative recovery time, as high doses are currently known to substantially increase the likelihood of side effects, such as propofol infusion syndrome [8].

Reinforcement learning (RL) is a subfield of machine learning that aims to solve sequential decision problems. RL interacts with the environment by creating an agent that obtains observational states and rewards from the environment to adjust the policy to maximize cumulative rewards and further achieve optimal control. In recent years, RL has been successfully applied in healthcare, including breast cancer screening [9], sepsis treatment [10], glioblastoma treatment [11], diabetes glucose control [12], [13], etc. Anesthesia infusion can be viewed as a decision-making process over time, in which the anesthesiologist selects the optimal combination of drug dosages based on the patient’s clinical information, as well as current physical status (*e.g.*, blood pressure, heart rate, depth of anesthesia, etc.) to maintain the patient’s vital signs within the target interval. Therefore, the anesthesia infusion problem is naturally amenable to modeling using RL, and there is a promise for using RL for more efficient and effective anesthesia administration.

However, traditional RL learns optimal policies by interacting with the environment through trial and error (see Fig. 1). However, anesthesia is a real clinical procedure, and any dose misuse can harm the patient. For instance,

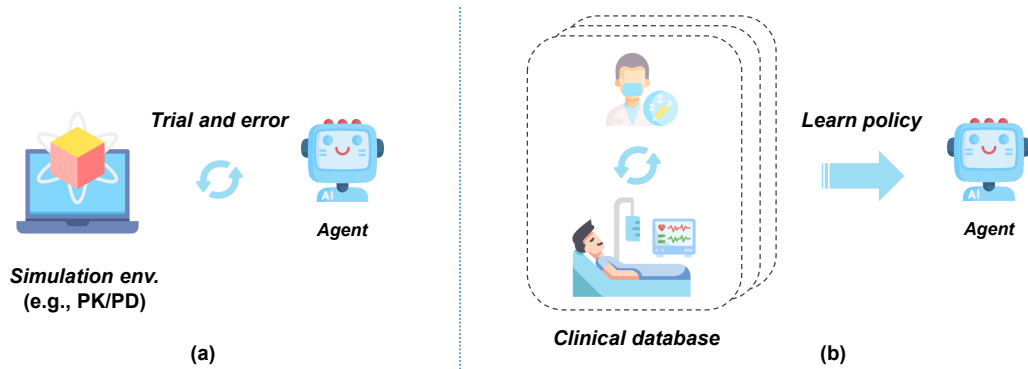


Fig. 1. Comparison of traditional RL (the left) and offline RL (the right) for automatic anesthesia policies learning.

overdosing propofol can damage the patient’s brain, and underdosing can cause unnecessary intraoperative awareness. It is inhumane and high-risk to train RL agents directly with the human body. To this end, Moore et al. [14] proposed pharmacokinetic/pharmacodynamic (PK/PD) models, which simulate the distribution and transfer of drugs in the human body through differential equations. However, such models traditionally involve many hyperparameters, which are usually derived from population statistical data, which means that such models ignore the specificity across patients and therefore may be highly inaccurate when modeling an individual’s environment. By means of personalized multi-task learning, Bird et al. [15] improved the PK/PD model. While such methods have proven successful experimentally, simulation-based learning remains distant from clinical applications. Recently, Shin et al. [16] proposed to learn a function that simulates the human environment based on real-world clinical anesthesia surgery data collected. The agent is then trained in this environment. All these algorithms mentioned above rely on environment modeling, and the reliability of trained agents depends on the simulation credibility.

Recently, offline reinforcement learning (ORL) has achieved impressive success and has received increasingly widespread attention. ORL is an advanced RL paradigm in which agents can learn suboptimal or even optimal policies from a collected offline dataset, while not requiring additional environmental exploration. Given that no interaction with the environment required, this new RL paradigm has attracted considerable interest in healthcare [17]–[20]. However, current ORL methods suffer from the problem of distributional shift [21]. For example, it is easy to learn overestimated Q functions, assigning high Q values to unseen operations, thus causing the model to select uncommon actions in the dataset, leading to undesired results. Moreover, the anesthesia task places higher safety requirements on the learned RL algorithms than common game tasks such as Atari [22]. An abnormal action by the agent may endanger the patient’s life. To this end, we propose Policy Constraint Q-Learning (PCQL), a data-driven RL algorithm for solving the problem of learning anesthesia infusion policies on real clinical datasets. To the best of our knowledge, we are the first to apply ORL in the context of real clinical anesthesia. We first learn a superior anesthesia infusion policy from the collected anesthesia dataset based on Conservative

Q-Learning (CQL; [23]), a class-leading ORL algorithm. To ensure safer decisions by the agent, we add a policy constraint regularization term during the learning of the agent policy. The regularizer is learnable and obtained by supervised training on a collected offline dataset, thus implicitly modeling the distribution of state-action pairs.

We conducted extensive experimental validation using a large clinical anesthesia dataset that we collected. Initially, we evaluated PCQL using the Off-Policy Evaluation method commonly employed in ORL. Our results suggested that PCQL outperformed all baselines, including human anesthesiologists’ policies and other RL methods. Subsequently, we evaluated different RL algorithms trained on retrospective data considering the clinical dose by anesthesiologists as the reference dose. Owing to the policy constraint regularization term, PCQL agrees better with the anesthesiologist’s policy than other RL baselines and has the lowest error in both MAPE and RMSE metrics. We further analyzed PCQL’s dose usage and observed that it required lower total drug doses to maintain patients’ vital signs within the target interval compared to the anesthesiologist’s clinical strategy. Moreover, PCQL’s recommended dose exhibited a stronger correlation with patient vital sign information and a faster adjustment frequency than the anesthesiologist’s policy, indicating the potential for accurate and personalized automated anesthesia. We also analyzed PCQL’s confidence intervals and found that they encompassed most anesthesiologists’ policies, thereby bolstering the model’s credibility and reliability. Finally, we utilized SHAP, an interpretable method, to analyze the model’s decision results and enhance the transparency of the automated anesthesia control system, which is pertinent for anesthesiologists’ utilization of the system.

II. RELATED WORKS

A. Automated anesthesia infusion

Automated anesthesia technology was first pioneered in the 1980s and has since advanced even further, as well as becoming more widely used [24]. This technology frees anesthesiologists from repetitive drug control tasks, allowing them to focus on the most critical aspects of each case, resulting in high-quality care for each patient. Moreover, automated anesthesia has shown great promise in improving

drug infusion regimens and allowing for better precision anesthesia. Pasin *et al.* [8] have shown that Bispectral-guided total intravenous anesthesia can reduce the need for propofol during induction, better maintain the target depth of anesthesia, and reduce recovery time compared to manual control. In a multicenter study, Puri *et al.* [4] showed that automated anesthesia can consistently achieve better performance than manual control. In a retrospective study, Brogi *et al.* [5] found that automated anesthesia can effectively reduce overshooting or undershooting of target physiological indicators and can effectively increase the duration of maintenance at the target interval. The current development of automated anesthesia infusion technology is usually closely related to the evolution of artificial intelligence. Moore *et al.* [14] used a proportional integral derivative (PID) controller for controlling propofol infusion rate, which in turn regulates Bispectral (BIS). Some other control algorithms for automatic anesthesia have also been investigated, such as model-predictive controllers [25], and rule-based controllers [26]. However, such traditional control methods are usually limited by linear assumptions. For this reason, Moore *et al.* [14] first proposed the use of discrete action of RL for anesthesia control and achieved better results than PID controllers. Subsequently, Lowery *et al.* [27] extended the work of [14] to continuous space. Schamberg *et al.* [28] used a more advanced actor-critic algorithm to train an agent for anesthesia control. Yun *et al.* [29] used a hierarchical RL algorithm to learn a high-level and a low-level policy. The high-level policy generates a target BIS trajectory, and the low-level policy uses this information to learn more stable infusion control.

B. Environmental simulation modeling of anesthesia

Environment modeling is fundamental to RL because the agent interacts with the environment through a "trial-and-error" paradigm, from which it learns the target policy. Existing automated anesthesia algorithms typically rely on the development of pharmacokinetic (PK) and pharmacodynamic (PD) models that simulate the response of a patient's BIS level to a specific drug dose. The PK model describes how the drug flows through the various compartments of the body (*e.g.*, brain, slow compartments, etc.) based on a system of equations for a particular drug. The PD model then maps the specific drug concentration at the effect site to the effect level, *i.e.*, BIS. However, simulation-based human environments usually involve a wide range of hyperparameters, which are usually derived from the statistical values of the population. This means that such models ignore the specificity between different patients and may therefore suffer unexpected inaccuracies when simulating the human environment. Bird *et al.* [15] uses multi-task learning techniques to personalize the PK/PD model to an individual level, while retaining statistical power, and the results show improved prediction accuracy. Despite the promising experimental results achieved by such methods, the simulation data are still not convincing for applications in clinical anesthesia. Several studies have attempted to learn on real-world collections of clinical anesthesia procedures, and the goal of such methods is to learn an environmental function that

simulates the body's response after receiving a certain dose of medication. Shin *et al.* [16] learned a transition function that simulates the human environment through supervised learning, and then used proximal policy optimization combined with behavioral cloning algorithms to learn automatic anesthesia policies based on this environment, and achieved promising experimental results. Unlike [16], our goal is to use real clinical data and train a scoring function that evaluates the plausibility of state-action pairs and thus guarantees safer decisions by the agent.

C. Offline Reinforcement Learning

Recently, an advanced RL paradigm, offline reinforcement learning (ORL), also known as batch reinforcement learning, has been proposed. Compared to traditional online reinforcement learning, ORL allows agents to learn superior policies from collected datasets without having to perform additional exploration in the environment. Given the property of "offline", and thus avoiding costly and dangerous, unethical exploratory actions when interacting with the environment, ORL has also gained widespread interest in the medical field. Kondrup *et al.* [17] used deep conservative reinforcement learning to determine the best ventilator settings for ICU patients. Emerson *et al.* [18] proposed using ORL to learn a safer blood glucose control strategy for people with Type 1 diabetes. Shiranthika *et al.* [19] developed the supervised optimal chemotherapy regimen, which can provide cancer patients with an optimal chemotherapy-dosing schedule, thus assisting oncologists in clinical decision-making. Wang *et al.* [20] used ORL to learn the optimal treatment strategy for sepsis patients in ICU. These are good illustrations of the potential of ORL applications in the medical field. However, the current ORL is susceptible to the problem of distributional shift. For example, it is easy to learn overestimated Q functions that assign high Q values to unseen operations, thus causing the model to select uncommon operations in the dataset, putting patients at risk. Researchers have made successive efforts to alleviate the overestimation problem, such as Batch Constrained Q-Learning [21], Advantage Weighted Actor-Critic [30], Conservative Q-Learning [23], etc.

D. Conservative Q-Learning

Conservative Q-Learning (CQL; [23]) is an ORL algorithm based on a Soft Actor-Critic improvement, which solves the problem of overestimating Q values in ORL by learning a conservative estimate of the Q function.

In a standard Q function, the loss function can be written as:

$$\begin{aligned} L_{DQN} &= \mathbb{E}_{s,a \sim D} \left[(Q(s,a) - B^{\pi_k} Q^k(s,a))^2 \right] \\ &= \mathbb{E}_{s,a \sim D} \left[\left(Q(s,a) - \left(r(s,a) + \gamma \max_{a'} Q(s',a') \right) \right)^2 \right], \end{aligned}$$

where B^{π_k} is the Bellman operator [31] on the currently learned policy π_k at iteration k . The true Q-value estimate of the state-action pair is approximated by minimizing $Q(s,a) - (r(s,a) + \gamma \max_{a'} Q(s',a'))$. However, due to the inability to

explore the environment, naive Q-learning tends to get overly optimistic Q values on offline data, which leads to the problem of overestimation. To this end, CQL adds a regularization term:

$$L_{CQL} = L_{DQN} + \alpha \mathbb{E}_{s \sim D} \left[\log \sum_a \exp(Q(s, a)) - \mathbb{E}_{a \sim D} [Q(s, a)] \right], \quad (1)$$

where α is the factor that relaxes the importance of the conservative term in the overall loss. The log-sum-exp term penalizes the action with the largest Q-value. The second term, on the other hand, guarantees to maximize Q-values for state-action pairs in the dataset. Thus, this conservative term allows high Q-values to be assigned only to actions within the distribution.

III. METHODOLOGY

We first formalize the procedure of automatic anesthetic administration, including the state space, action space, and reward design. Subsequently, we introduce a new constraint term that enables to constrain agent-predicted actions in the action distribution of the dataset. Finally, we describe the proposed overall framework and the specific training process.

A. The Problem Setting

In this study, the automatic administration of anesthesia is modeled as a Markov decision process (MDP) at finite time steps. Typically, the MDP is defined as a 5-tuple $\langle S, A, R, P, \gamma \rangle$. At time step t , the agent, in the current state $s_t \in S$, takes an action $a_t \in A$ and moves to the next state $s_{t+1} \in S$ according to the transition probability $P(s'|s, a)$. The agent expects to maximize the accumulated reward as it interacts with the environment.

In the problem setting, the agent is the controller that controls the delivery rate of the drug delivery device, while the environment is the patient in the perioperative phase of anesthesia. The agent is expected to predict the optimal dosage required by the patient at the current moment, based on the state observed from the environment, such as the patient's current vital signs information, and the history of drug administration. A suitable drug delivery controller should be able to maintain the patient's vital signs consistently during the surgery with minimal drug administration costs. Below are detailed definitions of the observed states, environment, reward functions, and agents.

1) *State Space*: The state space defines the information that the agent can observe at each moment, including the patient's clinical information, real-time vital signs, fluids, and other important information. The state space contains a total of 19 variables.

- Clinical information: age, gender, height, weight, BMI, ASA grade.
- Vital Signs: systolic arterial pressure (AP_{sys}), diastolic arterial pressure (AP_{dia}), mean arterial pressure (MAP), $MAP/AP_{sys}/AP_{dia}$ for the previous two moments.
- Analgesics: Sufentanil.
- The others: MAP target, MAP change, MAP target error.

The selection of these 19 observational variables was carefully determined through extensive discussions with clinical anesthesiologists. These variables encompass crucial clinical information, vital signs, and the utilization of Analgesics, which directly influence the dosage administered by the anesthesiologist. Additionally, we have incorporated the concepts introduced in [28], such as MAP target, MAP change, and MAP target error. These variables facilitate a faster perception of the target level, the subject's vital sign changes, and the true error for the agent. In the paper, the MAP target is defined as the average of the MAP of a patient over one operation.

$$MAP_t^* = \frac{1}{T} \sum_{t=1}^T MAP_t,$$

where T is the duration of an operation and MAP_t is the MAP of a patient at time step t . The MAP change is defined as the difference between the MAP at the current moment and the previous moment.

$$MAP(change)_t = MAP_t - MAP_{t-1}.$$

The MAP target error is defined as the distance between the current MAP and the target MAP.

$$MAP(error)_t = MAP_t - MAP_t^*.$$

2) *Action Space*: During automated anaesthetic infusion, more than one drug may be involved. In this study, the agent was only required to control propofol infusion rate. We also included the use of other drugs in the observed state to consider the synergistic effect of propofol with other drugs (e.g., analgesics). We calculated the maximum propofol dose used in the dataset P_{max} and performed a maximum normalization. The action space is thus $a \in \mathcal{A}$, continuous from 0 to 1. The final dose recommended by the agent is $a \cdot P_{max}$, and this allows for safer infusion control.

3) *Reward Function*: During surgery, the primary goal of the agent is to ensure that the patient's vital signs remain within a specific range. In this study, we focused on MAP in particular since it is a readily accessible measurement and serves as an effective indicator of the patient's current vital status. We expected the agent to maintain the patient's MAP around the MAP target. After thorough consultations with clinical anesthesiologists, the reward function is designed as follows,

$$R_{error}(s_t, a_t) = \begin{cases} +1 & \text{if } |MAP_t - MAP_t^*| \leq 15\%MAP_t^*; \\ +0.5 & \text{if } 15\%MAP_t^* \leq |MAP_t - MAP_t^*| \leq 30\%MAP_t^*; \\ -1 & \text{else.} \end{cases}$$

We adopted a segmented reward function design instead of a continuous one to enable the learning of more robust policies. In the surgical environment, there are numerous sources of noise and unpredictable events, such as intubation and sensor detachment. Utilizing a continuous reward function, such as mean squared error, could potentially cause the agent to overfit to these abnormal fluctuations, thereby increasing risks.

In addition, we hoped that the agent would use the lowest possible dosage to maintain the patient's vital signs within

the target interval. To this end, we added a dose penalty as follows.

$$R_{dosage}(s_t, a_t) = -\frac{|\text{MAP}_t - \text{MAP}_t^*|}{\text{MAP}_t^*} a_t,$$

where $\frac{|\text{MAP}_t - \text{MAP}_t^*|}{\text{MAP}_t^*}$ is the adaptive correction factor. We do not want the agent to be overly conservative in the use of propofol, especially when the patient's MAP deviates from the target mean arterial pressure. Instead, the agent is expected to use a lower propofol dose to maintain the patient's vital signs within the target interval only when the patient's MAP is around the target. In summary, the overall reward function is calculated as follows.

$$R_{total}(s_t, a_t) = R_{error}(s_t, a_t) + R_{dosage}(s_t, a_t). \quad (2)$$

B. Extending CQL As An Actor-critic Variant

With the CQL constraint, we are able to learn a more reasonable Q-value estimation function by the Q-Learning method. However, Q-Learning may be too deterministic, especially during offline training, due to the presence of ε -greedy policy, which may limit the exploration of optimal actions by the agent. In this case, a stochastic policy may be a better choice.

To learn an explicitly parameterized stochastic policy, following [23], we first instantiate the CQL as an actor-critic algorithm. In this framework, the CQL constraint term is first added to the training of the Q function and subsequently jointly trained with the actor as a critic. As in common implementations of actor-critic, an entropy term is added to the objective of the policy optimization in order to encourage policy exploration. Hereby, the gradient of the objective function of the actor can be formalized as

$$J(\pi_\theta) = \mathbb{E}_{s \sim D, a \sim \pi_\theta(a|s)} [\log \pi_\theta(a|s) Q(s, a) + \beta H(\pi_\theta(s))], \quad (3)$$

where H is the entropy of the policy and β is the temperature factor used to relax the relative importance between the entropy term and reward.

C. Policy Constraint Q-Learning

In the aforementioned actor-critic framework, CQL serves as a constraint on the Q-function for accurate estimation of Q-values. However, during the iterative process of policy evaluation and improvement, errors in Q-value estimation can persist, potentially leading to an erroneous Q-function that drives the policy towards out-of-distribution (OOD) actions. To safeguard agent's actions, we add a constraint term to the CQL from a policy perspective. Our goal is to learn a function from a collected offline dataset to evaluate the plausibility of state-action pairs.

1) *Policy Constraint*: The scoring function consists of two components, the environment transition model h and the behavior prediction model g . The environmental transition model aims to predict the new state to which the current state s will be transferred after the selection of action a , i.e., $s' = h(s, a)$. And the behavior prediction model predicts the possible actions that may be taken given two consecutive states, i.e., $a = g(s, s')$. For the action $\hat{a} = f(s)$ predicted

Algorithm 1 Policy Constraint Q-Learning (PCQL)

- 1: **Input**: training set \mathcal{D} , number of training epochs E
 - 2: Construct offline replay buffer \mathcal{D}' from \mathcal{D} using (2)
 - 3: Initialize policy constraint network Φ_γ , critic network Q_θ and actor network Π_ϕ with random parameters
 - 4: Initialize target network $Q_{\theta'}$ with weights $\theta' = \theta$
 - 5: **for** epoch=1 **to** E **do**
 - 6: **repeat**
 - 7: Sample a mini-batch from replay buffer \mathcal{D}'
 - 8: Train policy constraint Φ_γ by minimizing Eq. (5) and Eq. (6)
 - 9: Train critic Q_θ by minimizing Eq. (1)
 - 10: Train actor Π_ϕ by maximizing Eq. (4)
 - 11: Every N steps update target network with $\theta' = \theta$
 - 12: **until** all mini-batch are sampled
 - 13: **end for**
 - 14: **Output**: optimal PCQL policy Π_ϕ^*
-

by the agent at the current state s (i.e., the dose), we first imagine its subsequent state using the environmental transition model, and then reason about the possible action \hat{a} using the behavioral prediction model, i.e., $\hat{a} = g(s, h(s, \hat{a}))$. We constrain the action \hat{a} predicted by the agent not to be too far from the action \hat{a} inferred by the behavioral prediction model as likely to occur. A simple approach is to

$$\Phi_{\text{simple}}(s, \hat{a}) = \|g(s, h(s, \hat{a})) - \hat{a}\|_2,$$

where $\|\cdot\|_2$ denotes the Euclidean distance.

However, such a deterministic constraint leads to the fact that the actions predicted by the agent are supervised by the output of the behavior prediction model. This is perhaps too stringent, thus weakening the exploration of other, more optimal actions. For this reason, we propose an alignment constraint in the latent space. The network architectures of the environment transition model h and the behavior prediction model g are similar, and both consist of an encoder, a predictor, and a projector. The predictor aims to predict the target information (e.g. s' for $h(s, a)$), while the projector is responsible for mapping the output of the encoder to the latent space, a process we denote as $\text{Prj}(\cdot)$. We use the cross entropy function for alignment constraints, i.e.

$$\Phi(s, \hat{a}) = \text{softmax} \left(\frac{\text{Prj}(h(s, \hat{a}))}{\tau} \right) \cdot \log \left(\text{softmax} \left(\frac{\text{Prj}(h(s, \hat{a}))}{\tau} \right) \right),$$

where τ is the temperature factor for scaling the learning difficulty. Combining Eq. (3), the final objective function of actor is

$$L_{total} = J(\pi_\theta) + \Phi(s, \hat{a}). \quad (4)$$

2) *Training Strategy*: We next discuss the training of the environment transfer model h and the behavior prediction model g . In fact, both models can be learned from an offline dataset following a supervised learning paradigm. We sample $(s, a, s') \sim \mathcal{D}$ and train these two models using an autoencoder-like approach, i.e.,

$$L_{\text{cycle}} = (g(s, h(s, a)) - a) + (h(s, g(s, s')) - s'). \quad (5)$$

Also, we add the entropy consistency as follow,

$$L_{entropy} = \mathcal{H}(h(s, \hat{a}), h(s, a)) + \mathcal{H}(g(s, s'), g(s, s')), \quad (6)$$

where $\mathcal{H}(a, b) = -\text{softmax}(a/\tau) \cdot \log(\text{softmax}(b/\tau))$ and τ is the temperature factor. We jointly train the scoring function $\Phi(s, a)$ and the agent of the CQL algorithm, while adding consistency constraints to the action loss of CQL. The pseudo-code of the overall training algorithm of the proposed PCQL is shown in Algorithm 1.

IV. EXPERIMENTS AND RESULTS

A. Data Collection and Training Details

We collected 6,303 brain-neurological surgeries involving general anesthesia from West China Hospital. During each surgery, the patient's vital sign signals as well as the infused drugs information was recorded in real time at one-minute intervals. We filtered the surgeries that met any of the following conditions: i) missing corresponding dosing information; ii) surgical time shorter than half an hour; iii) severe missing vital sign records during surgery; and iv) samples using inhaled anesthesia, such as sevoflurane or desflurane instead of propofol. We used the k-NN method to fill in the remaining missing values. Finally, a total of 1,293 surgeries with 284,281 anesthesia records were obtained. Relevant clinical data are summarized in Fig. 2. All experiments are conducted based on d3rlpy [32]. We employed five-fold cross-validation in order to ensure a reliable evaluation of models. We use Adam [33] as the optimizer and train for 200 epochs with batchsize 256. The temperature factor for entropy consistency is set to 0.08. The learning rates of actor, critic, environment transition model, and action prediction model are set to 1×10^{-4} , 3×10^{-4} , 1×10^{-4} and 3×10^{-4} respectively. These hyperparameter settings are obtained from the grid search to ensure the performance and generalization of the agent. Specifically, the learning rate is searched in $\{1 \times 10^{-3}, 3 \times 10^{-3}, 1 \times 10^{-4}, 3 \times 10^{-4}\}$ while the temperature factor is searched in $\{0.03, 0.05, 0.08\}$.

B. Evaluation Protocol

1) *Baselines*: We mainly set the following baselines, including anesthesiologists' clinical policy (ANE), soft actor-critic (SAC, [34]), batch constrained deep Q-learning (BCQ, [35]), conservative Q-learning (CQL, [23]), twin delayed deep deterministic policy gradient with behavior cloning (TD3+BC, [36]) and mildly conservative Q-Learning (MCQ, [37]). ANE represents the actual clinical practice of anesthesiologists in the collected dataset and serves as our primary reference of recommended dose. To demonstrate the advantages of offline reinforcement learning in the absence of an available environment, we include a popular off-policy RL algorithm, SAC, in our comparisons. Furthermore, to highlight the strengths of PCQL, we compare it with several state-of-the-art offline RL techniques, including BCQ, CQL, TD3+BC, and MCQ. BCQ aims to minimize the distance between selected actions and actions in the dataset by employing a conditional VAE on the state. TD3+BC achieves policy constraints by simply adding a behavioral cloning term to the TD3 algorithm [38]. MCQ is an improved algorithm for CQL by assigning pseudo target values to actively train OOD actions.

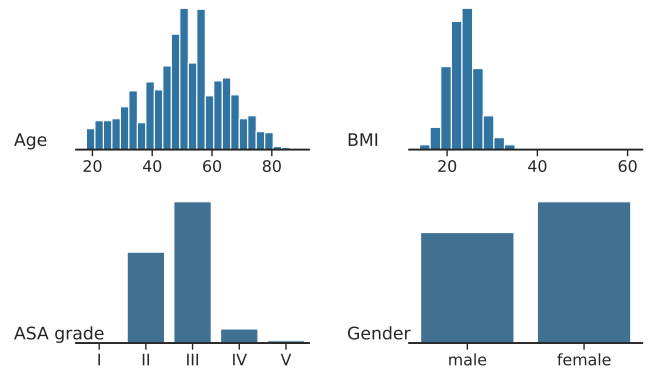


Fig. 2. Clinical data distributions of collected dataset.

2) *Off Policy Evaluation*: In our experimental settings, no interactive environment is available, so trained policies cannot directly compute cumulative rewards by interacting with the environment. We used the off-policy evaluation (OPE) method commonly used in ORL for performance evaluation [39], which was done on a divided test set. In recent experiments in healthcare settings, researchers have found that the Fitted Q Evaluation (FQE; [35]) approach consistently yields accurate policy performance evaluation results [17], [40], [41], and we follow this practice. FQE uses a fixed trained policy, and a re-set Q function, and then re-trains the Q function. The retrained Q function estimates how much return the trained policy can obtain from initial states. The higher the estimated returns, the better the expected performance of the policy. To comprehensively evaluate the performance of policies, we also use Soft Off-Policy Classification (Soft OPC; [42]) as another supplementary metric. Soft OPC is a policy-based evaluation metric that measures algorithm performance by comparing action values between successful and unsuccessful trajectories. If the action values in successful trajectories are significantly higher than those in unsuccessful trajectories, it indicates better policy performance of the algorithm.

3) *Retrospective Evaluation*: Since the data we collect is retrospective from the operating room, we can also evaluate the model in a "consult" mode. In this setting, the intelligence can see the historical state and action information, and then give the recommended dose for the current state. At the same time, we take the actual clinical dose given by the anesthesiologist as the optimal dose and calculate the difference between the recommended dose and the actual dose. Although we could not know how the human body responds to the recommended dose by the agent, we argue that this evaluation can still in part reflect the performance of the agent when faced with the real world. To measure the discrepancy between the recommended doses and the actual ones, we used two common regression metrics for evaluation: mean absolute percentage error (MAPE) and root mean square error (RMSE). In this case, MAPE is defined as

$$\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \frac{|y_{i,t} - y_{i,t}^*|}{\max(\varepsilon, y_{i,t}^*)} 100\%$$

where N is the number of episodes and T is the length of an episode. $y_{i,t}$, $y_{i,t}^*$ denote the recommended dose of the

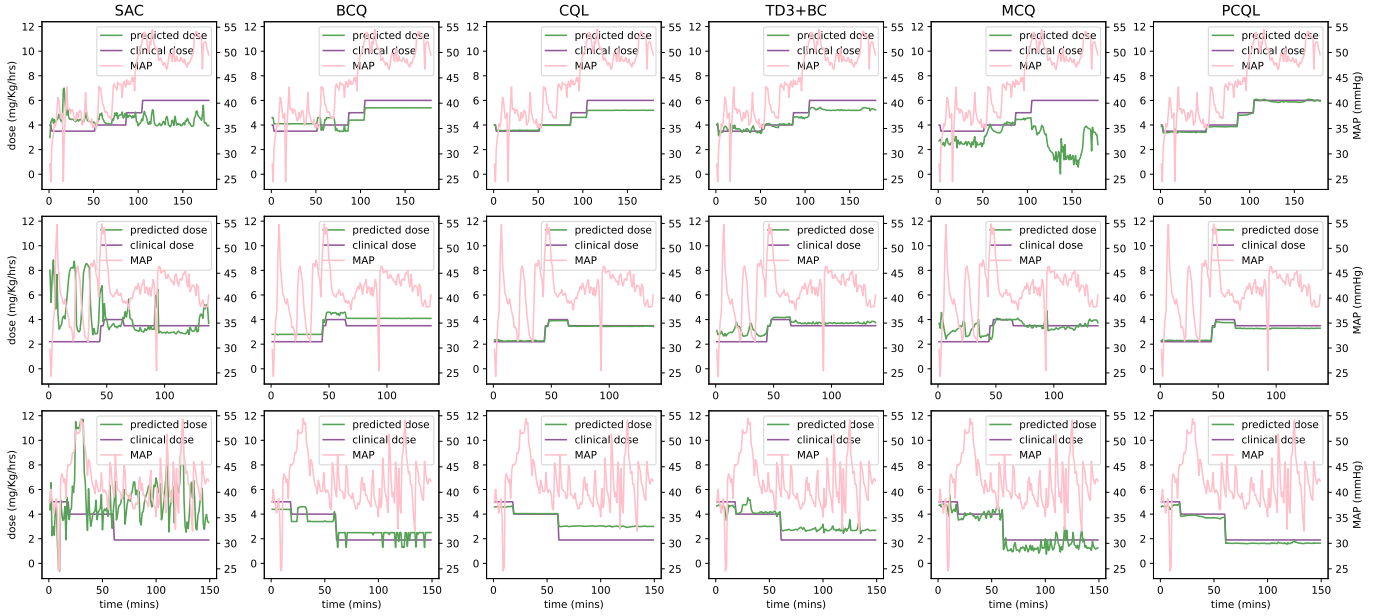


Fig. 3. Comparison of recommended doses in general anesthesia cases among different algorithms and anesthesiologists (each row represents a case).

TABLE I
COMPARISON OF THE PREDICTED RETURNS OF DIFFERENT METHODS

	ANE	SAC [34]	BCQ [35]	CQL [23]	TD3+BC [36]	MCQ [37]	PCQL
Initial state \uparrow	59.071 \pm 0.663	43.325 \pm 12.004	64.919 \pm 2.099	66.598 \pm 1.111	66.693 \pm 1.506	69.043 \pm 1.714	73.069\pm1.302
Soft OPC \uparrow	0.072 \pm 0.002	-1.274 \pm 1.015	-0.184 \pm 2.974	0.633 \pm 0.388	-0.394 \pm 1.814	0.434 \pm 0.879	1.061\pm0.416

TABLE II
PERFORMANCE COMPARISON IN RETROSPECTIVE EVALUATION

	SAC [34]	BCQ [35]	CQL [23]	TD3+BC [36]	MCQ [37]	PCQL
MAPE (%) \downarrow	78.557 \pm 30.034	20.510 \pm 11.236	8.884 \pm 2.435	9.762 \pm 1.822	12.214 \pm 5.486	7.860\pm1.475
RMSE \downarrow	3.522 \pm 1.270	0.800 \pm 0.497	0.323 \pm 0.096	0.390 \pm 0.053	0.437 \pm 0.087	0.279\pm0.068

agent and the actual dose for the i episode at the t moment, respectively. The ε stands for a very small number to avoid division by zero error. RMSE is defined as

$$\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \sqrt{(y_{i,t} - y_{i,t}^*)^2}$$

C. Quantitative Analysis of Performance

In this section, we compare the PCQL with baselines, including ANE, BCQ, CQL, TD3+BC and MCQ. First, we evaluate the policies on the initial state and Soft OPC metrics, and the quantitative results are presented in Table I. We can see that PCQL is expected to yield higher returns than the other policies. The estimated return of PCQL is 1.24 times higher than that of the anesthesiologist's strategy and also superior to MCQ (73.069 compared to 69.043). Moreover, PCQL achieved the highest Soft OPC score, suggesting that the actions recommended by PCQL are more likely to outperform the actions taken by the behavior policy. It is worth noting that all the RL-based policies perform well except SAC. One possible reason for its poor performance is that the off-policy

RL algorithm cannot correct the value estimation errors of the OOD actions in time due to the unavailability of the environment. Consequently, the extrapolation errors continue to propagate during training, ultimately leading to the training failure. In contrast, the ORL algorithms mitigate this issue by constraining the policy distribution either Q-value estimation. PCQL combines their ideas and achieves more advantageous results.

Subsequently, we evaluated the trained models in "consultation" mode, focusing on the discrepancy between the recommended dose from different policies and the actual clinical dose. To compare the performance of the policies, we used two common metrics, namely MAPE and RMSE, and present the results in Table II. As shown in the result, PCQL consistently demonstrates superior and robust performance. While BCQ achieved a less favorable result (20.510%) on the MAPE metric, PCQL, TD3+BC, and CQL were all within 10% of each other. Although MCQ showed promising performance on the initial state metric, it only achieved a 12.214% result on the MAPE metric. Notably, SAC, the traditional RL approach, performed poorly on both metrics, highlighting the importance

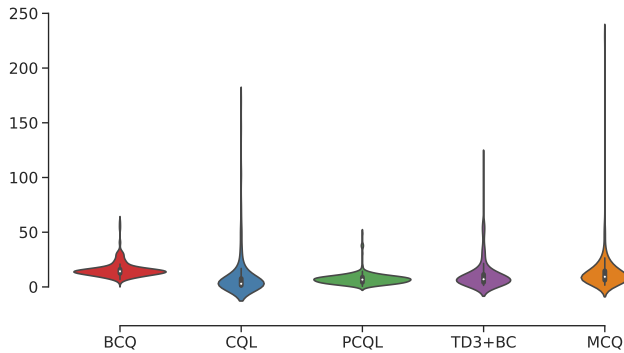


Fig. 4. Comparison of different offline RL algorithms in terms of absolute percentage error (APE) metric.

TABLE III
COMPARISON OF MEAN DOSES USED

	ANE	PCQL
Mean dose (mg/Kg/hrs)	3.921	3.719

of incorporating certain constraints when the environment is not available.

We further analyzed the MAPE performance of different ORL methods in Fig. 4. Our results show that the absolute percentage error (APE) of PCQL had a lower overall distribution, and thanks to the consistency constraint, the APE of PCQL exhibited reduced variance and fewer outliers compared to the baselines. Finally, we visualized the performance of different policies in several randomly selected cases from the test set. Overall, doses recommended by PCQL aligned well with the clinical doses of the anesthesiologists. In contrast, SAC showed poor fluctuations. Other ORL algorithms also aligned with the clinical doses but exhibited different local oscillations (e.g., BCQ in the third row of Fig. 3) and deviations (e.g., MCQ in the first row of Fig. 3).

D. Model Recommended Dosage Analysis

Given that PCQL significantly outperforms baselines in both previous performance evaluations, we next focus exclusively on analyzing PCQL’s recommended dose usage and comparing its advantages and disadvantages with the anesthesiologist’s policy.

The statistical results of mean dose values on the test set are shown in Table III. We found that the mean dose value of 3.719 mg/Kg/hrs recommended by PCQL is slightly lower than the actual clinical dose of 3.921 given by the anesthesiologist by 0.202 mg/Kg/hrs. The difference of 0.202 mg/Kg/hrs reduction is clinically acceptable, because even two experienced anesthesiologists may give a difference of 1 mg/Kg/hrs in the anesthetic dose for the same case of anesthetic surgery. At the same time, it is clinically meaningful if the PCQL policy is able to maintain the patient’s depth of anesthesia within the target range using a lower drug dose. First, it would mean that the problem of over-anesthesia could be alleviated, reducing the incidence of post-anesthesia syndrome and thus improving

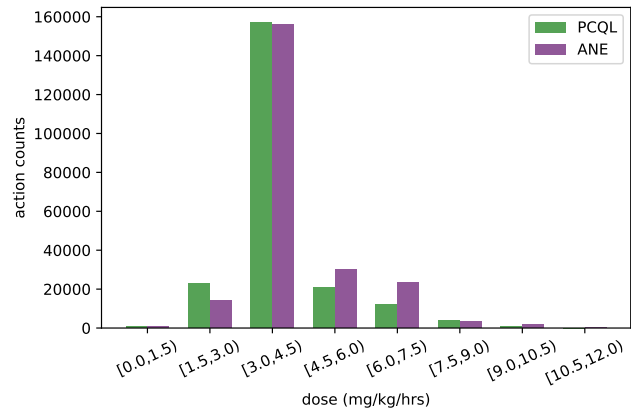


Fig. 5. Comparison of the distribution of doses used by PCQL and anesthesiologists.

TABLE IV
COMPARISON OF THE CORRELATION

	ANE	PCQL
Correlation	0.081	0.245

the quality of patients’ surgery. Second, lower dose use often means less consumption of medical resources, which helps reduce the cost of patient care. As a supplement, we also present the distributional difference of doses recommended by PCQL and human anesthesiologists in the test set (as depicted in Fig. 5)

Intriguingly, we noted that the dose recommended by PCQL fluctuated around the actual dose given by the anesthesiologist and that this fluctuation was positively correlated with the MAP of the patient (for instance, MCQ in the first row of Fig. 3). To verify this, we calculated Pearson correlation coefficients between recommended doses and MAPs for PCQL and anesthesiologist policies, respectively. The Pearson correlation coefficient was calculated by

$$\rho_{X,Y} = \frac{\mathbb{E}((X - \mu_X)(Y - \mu_Y))}{\sigma_X \sigma_Y},$$

where μ_X and σ_X denote the mean and variance of the variable X , respectively; similarly for Y .

As indicated in Table IV, the dose recommended by PCQL showed a stronger correlation with MAP compared to anesthesiologists’ (0.245 compared to 0.081). This is probably because in clinical practice, anesthesiologists are typically responsible for multiple procedures and cannot attend to the same patient continuously. Furthermore, the focus of human anesthesiologists cannot be sustained for extended periods. In contrast, as illustrated in Fig. 3, the dose recommended by PCQL are more sensitive and responsive to changes in the patient’s vital signs. This suggests the potential for computer-assisted precision anesthesia and offers the prospect of delivering higher quality, more personalized anesthesia infusion management.

E. Confidence Interval Analysis of the Model

It is a fact that even experienced anesthesiologists cannot determine the optimal dose—they usually pick a preferred dose

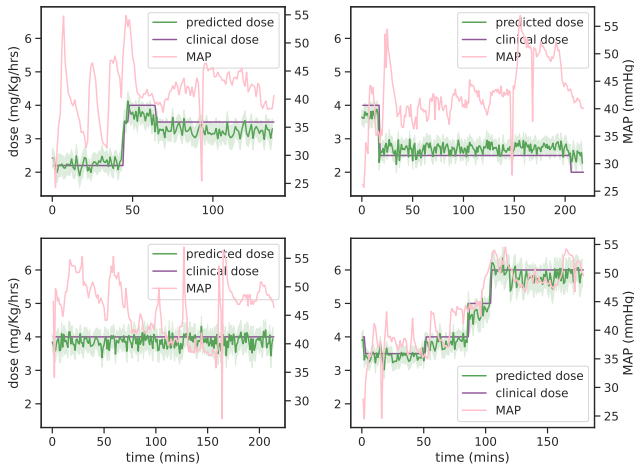


Fig. 6. Confidence interval estimation of PCQL in general anesthesia cases.

within a confidence interval that they believe is reasonable. To estimate the confidence interval for the PCQL, we simply modified the output mechanism of the PCQL so that it probabilistically samples an action when making inferences, rather than the best one. We choose the common Gaussian policy of $\pi_{\theta}(a|s) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{a-f_{\theta}(s)}{2\sigma^2}\right)$. At state s , the actions taken when sampled by this policy obey a normal distribution with a mean of $f_{\theta}(s)$ and a variance of σ^2 . Fig. 6 gives the results of sampling 100 times, and it can be seen that the confidence interval estimated by PCQL is on the safe side without any serious deviation from the actual clinical dose. The confidence interval recommended by PCQL covers the majority of the anesthesiologist’s clinical policy, and this is encouraging. Because it means that under this situation, the anesthesiologist’s policies can be derived from sampling within the confidence interval of the PCQL. Although we could not be informed in retrospective data about the confidence intervals considered by anesthesiologists, we argue that this enhances the credibility and reliability of PCQL.

F. Interpretability of the Model

The deep neural networks used to implement our policy are notorious for their black-box characteristics. However, if an anesthesiologist does not know how the machine model makes decisions, then he will not be able to trust the use of the model. To increase the transparency of the model, we analyzed the prediction results of the model using SHapley Additive exPlanations (SHAP; [43]). SHAP is a game-theoretic interpretation approach that predicts the output of a model by fitting a linear interpretation model, and calculates SHAP values for observed features, which indicate the importance of the features in influencing the model output. SHAP values have been widely applied to analyze the relative importance of features in various machine learning models for medical applications [28], [44], [45]. We used the publicly available SHAP library¹ to calculate the SHAP values for each feature in the observation space, and obtained the Absolute Mean SHAP

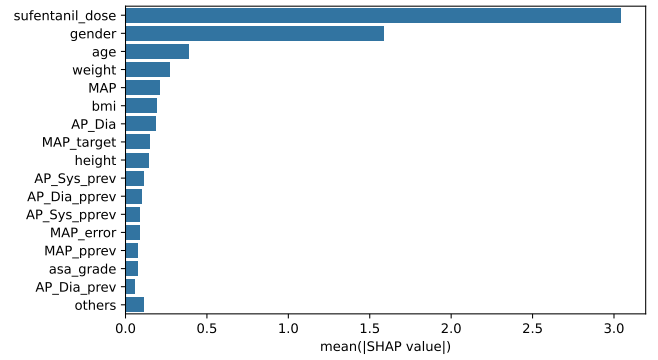


Fig. 7. Results of absolute mean SHAP values for different input components of PCQL.

Score for each feature by taking the average of the absolute values.

As shown in Fig. 7, we found that the output of PCQL (*i.e.*, the infusion rate of propofol) was most affected by sufentanil dose. We attribute this to the fact that sufentanil is an immediate injection analgesic that rapidly depresses the patient’s respiration and, in addition, acts synergistically with propofol, thus greatly influencing the output of the model. Clinical characteristics of patients such as gender, age, and weight also significantly influenced model output, followed by real-time vital sign characteristics such as MAP, AP_{Dia}, and MAP target. We believe this is similar to the clinical decision making of the anesthesiologist—the anesthesiologist usually determines the primary base dose by pre-assessing the patient’s physical condition, and adjust the base dose according to changes in the patient’s vital signs during surgery.

V. CONCLUSION

In this study, we proposed PCQL, an ORL algorithm for automated anesthesia infusion. To ensure the safety of the dosing policy, we have imposed constraints on both the policy distribution and value estimation. Through extensive experiments on a large collected anesthesia dataset, we validated the effectiveness of PCQL. The experimental results demonstrated that our method outperformed baseline methods. The recommended dose generated by PCQL was consistent with the anesthesiologist’s strategy, while utilizing a lower total dose and being more responsive to the patient’s vital signs. This highlights the potential for future personalized and precise anesthesia. Additionally, we employed SHAP to perform interpretability analysis on PCQL’s results, enhancing the transparency of our approach.

However, it is important to acknowledge that further validation through real-world clinical experiments is necessary to fully assess the effectiveness of PCQL. Our study represents a promising step towards achieving true automation in anesthesia, and we recognize the need for continued research and refinement in this area. We hope that our findings will inspire further advancements in the field to ultimately contribute to improved patient care.

¹<https://github.com/slundberg/shap>

VI. LIMITATIONS AND FUTURE WORK

Our study confirms the practicality of using ORL for automated anesthesia, although there are several potential concerns. Firstly, our data is limited to a single center, which does not allow us to assess PCQL's performance in a multicenter context. In future work, we intend to gather clinical anesthesia data from multiple centers and evaluate the model's effectiveness accordingly. For this purpose, an appropriate combination of federated learning [46]–[48] and ORL might be a compelling direction. Secondly, for off-policy evaluation, we employed FQE and Soft OPC as our evaluation protocols. While there is ample literature demonstrating the reliability and credibility of these metrics, and PCQL has shown favorable results based on these metrics, there is still a lack of comprehensive comparative experiments that establish the alignment of these evaluation metrics with real-world performance. It is crucial to ensure that the results of OPE accurately reflect the clinical outcomes observed in real-world scenarios. Once this alignment is established, it may be possible to deploy ORL-trained agents in real clinical settings under controlled conditions. Further research is warranted to explore the design of OPE in the specific context of healthcare.

REFERENCES

- [1] T. G. Weiser, S. E. Regenbogen, K. D. Thompson, A. B. Haynes, S. R. Lipsitz, W. R. Berry, and A. A. Gawande, "An estimation of the global volume of surgery: a modelling strategy based on available data," *The Lancet*, vol. 372, no. 9633, pp. 139–144, 2008.
- [2] P. Kempthorne, W. W. Morriss, J. Mellin-Olsen, and J. Gore-Booth, "The wfsa global anesthesia workforce survey," *Anesthesia & Analgesia*, vol. 125, no. 3, pp. 981–990, 2017.
- [3] A. M. Afonso, J. B. Cadwell, S. J. Staffa, D. Zurakowski, and A. E. Vinson, "Burnout rate and risk factors among anesthesiologists in the united states," *Anesthesiology*, vol. 134, no. 5, pp. 683–696, 2021.
- [4] G. D. Puri, P. J. Mathew, I. Biswas, A. Dutta, J. Sood, S. Gombur, S. Palta, M. Tsering, P. Gautam, A. Jayant *et al.*, "A multicenter evaluation of a closed-loop anesthesia delivery system: a randomized controlled trial," *Anesthesia & Analgesia*, vol. 122, no. 1, pp. 106–114, 2016.
- [5] E. Brogi, S. Cyr, R. Kazan, F. Giunta, and T. M. Hemmerling, "Clinical performance and safety of closed-loop systems: a systematic review and meta-analysis of randomized controlled trials," *Anesthesia & Analgesia*, vol. 124, no. 2, pp. 446–455, 2017.
- [6] C. Zauter, A. Joosten, J. Rinchart, M. M. Struys, and T. M. Hemmerling, "Autonomous systems in anesthesia: where do we stand in 2020? a narrative review," *Anesthesia & Analgesia*, vol. 130, no. 5, pp. 1120–1132, 2020.
- [7] M. Ghita, M. Neckebroek, C. Muresan, and D. Copot, "Closed-loop control of anesthesia: Survey on actual trends, challenges and perspectives," *Ieee Access*, vol. 8, pp. 206 264–206 279, 2020.
- [8] L. Pasin, P. Nardelli, M. Pintaudi, M. Greco, M. Zambon, L. Cabrini, and A. Zangrillo, "Closed-loop delivery systems versus manually controlled administration of total iv anesthesia: a meta-analysis of randomized clinical trials," *Anesthesia & Analgesia*, vol. 124, no. 2, pp. 456–464, 2017.
- [9] A. Yala, P. G. Mikhael, C. Lehman, G. Lin, F. Strand, Y.-L. Wan, K. Hughes, S. Satuluru, T. Kim, I. Banerjee *et al.*, "Optimizing risk-based breast cancer screening policies with reinforcement learning," *Nature medicine*, vol. 28, no. 1, pp. 136–143, 2022.
- [10] A. Raghun, M. Komorowski, L. A. Celi, P. Szolovits, and M. Ghassemi, "Continuous state-space models for optimal sepsis treatment: a deep reinforcement learning approach," in *Machine Learning for Healthcare Conference*. PMLR, 2017, pp. 147–163.
- [11] A. E. Zade, S. S. Haghghi, and M. Soltani, "Reinforcement learning for optimal scheduling of glioblastoma treatment with temozolomide," *Computer Methods and Programs in Biomedicine*, vol. 193, p. 105443, 2020.
- [12] T. Zhu, K. Li, P. Herrero, and P. Georgiou, "Basal glucose control in type 1 diabetes using deep reinforcement learning: An in silico validation," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 4, pp. 1223–1232, 2020.
- [13] G. Noaro, T. Zhu, G. Cappon, A. Facchinetti, and P. Georgiou, "A personalized and adaptive insulin bolus calculator based on double deep q-learning to improve type 1 diabetes management," *IEEE Journal of Biomedical and Health Informatics*, 2023.
- [14] B. L. Moore, T. M. Quasny, and A. G. Doufas, "Reinforcement learning versus proportional–integral–derivative control of hypnosis in a simulated intraoperative patient," *Anesthesia & Analgesia*, vol. 112, no. 2, pp. 350–359, 2011.
- [15] A. Bird, C. Williams, and C. Hawthorne, "Multi-task time series analysis applied to drug response modelling," in *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019, pp. 2174–2183.
- [16] M. Shin and J. Kim, "Joint behavioral cloning and reinforcement learning method for propofol and remifentanyl infusion in anesthesia," in *2021 International Conference on Information Networking (ICOIN)*. IEEE, 2021, pp. 849–852.
- [17] F. Kondrup, T. Jiralerspong, E. Lau, N. de Lara, J. Shkrob, M. D. Tran, D. Precup, and S. Basu, "Towards safe mechanical ventilation treatment using deep offline reinforcement learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 13, 2023, pp. 15 696–15 702.
- [18] H. Emerson, M. Guy, and R. McConville, "Offline reinforcement learning for safer blood glucose control in people with type 1 diabetes," *arXiv preprint arXiv:2204.03376*, 2022.
- [19] C. Shiranthika, K.-W. Chen, C.-Y. Wang, C.-Y. Yang, B. Sudantha, and W.-F. Li, "Supervised optimal chemotherapy regimen based on offline reinforcement learning," *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 9, pp. 4763–4772, 2022.
- [20] Z. Wang, H. Zhao, P. Ren, Y. Zhou, and M. Sheng, "Learning optimal treatment strategies for sepsis using offline reinforcement learning in continuous space," in *International Conference on Health Information Science*. Springer, 2022, pp. 113–124.
- [21] S. Fujimoto, D. Meger, and D. Precup, "Off-policy deep reinforcement learning without exploration," in *International conference on machine learning*. PMLR, 2019, pp. 2052–2062.
- [22] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling, "The arcade learning environment: An evaluation platform for general agents," *Journal of Artificial Intelligence Research*, vol. 47, pp. 253–279, 2013.
- [23] A. Kumar, A. Zhou, G. Tucker, and S. Levine, "Conservative q-learning for offline reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 1179–1191, 2020.
- [24] T. Wingert, C. Lee, and M. Cannesson, "Machine learning, deep learning, and closed loop devices-anesthesia delivery," *Anesthesiology clinics*, vol. 39 3, pp. 565–581, 2021.
- [25] Y. Sawaguchi, E. Furutani, G. Shirakami, M. Araki, and K. Fukuda, "A model-predictive hypnosis control system under total intravenous anesthesia," *IEEE transactions on biomedical engineering*, vol. 55, no. 3, pp. 874–887, 2008.
- [26] N. Liu, T. Chazot, A. Genty, A. Landais, A. Restoux, K. McGee, P.-A. Laloë, B. Trillat, L. Barvais, and M. Fischler, "Titration of propofol for anesthetic induction and maintenance guided by the bispectral index: closed-loop versus manual control: a prospective, randomized, multicenter study," *The Journal of the American Society of Anesthesiologists*, vol. 104, no. 4, pp. 686–695, 2006.
- [27] C. Lowery and A. A. Faisal, "Towards efficient, personalized anesthesia using continuous reinforcement learning for propofol infusion control," in *2013 6th International IEEE/EMBS Conference on Neural Engineering (NER)*. IEEE, 2013, pp. 1414–1417.
- [28] G. Schamberg, M. Badgeley, B. Meschede-Krasa, O. Kwon, and E. N. Brown, "Continuous action deep reinforcement learning for propofol dosing during general anesthesia," *Artificial Intelligence in Medicine*, vol. 123, p. 102227, 2022.
- [29] W. J. Yun, M. Shin, D. Mohaisen, K. Lee, and J. Kim, "Hierarchical deep reinforcement learning-based propofol infusion assistant framework in anesthesia," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [30] A. Nair, A. Gupta, M. Dalal, and S. Levine, "Awac: Accelerating online reinforcement learning with offline datasets," *arXiv preprint arXiv:2006.09359*, 2020.
- [31] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

- [32] T. Seno and M. Imai, "d3rlpy: An offline deep reinforcement learning library," *Journal of Machine Learning Research*, vol. 23, no. 315, pp. 1–20, 2022. [Online]. Available: <http://jmlr.org/papers/v23/22-0017.html>
- [33] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations (ICLR)*, 2014.
- [34] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel *et al.*, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.
- [35] H. Le, C. Voloshin, and Y. Yue, "Batch policy learning under constraints," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3703–3712.
- [36] S. Fujimoto and S. S. Gu, "A minimalist approach to offline reinforcement learning," *Advances in neural information processing systems*, vol. 34, pp. 20 132–20 145, 2021.
- [37] J. Lyu, X. Ma, X. Li, and Z. Lu, "Mildly conservative q-learning for offline reinforcement learning," in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35. Curran Associates, Inc., 2022, pp. 1711–1724.
- [38] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International conference on machine learning*. PMLR, 2018, pp. 1587–1596.
- [39] C. Voloshin, H. M. Le, N. Jiang, and Y. Yue, "Empirical study of off-policy policy evaluation for reinforcement learning," *arXiv preprint arXiv:1911.06854*, 2019.
- [40] S. Tang and J. Wiens, "Model selection for offline reinforcement learning: Practical considerations for healthcare settings," in *Machine Learning for Healthcare Conference*. PMLR, 2021, pp. 2–35.
- [41] T. Zhu, K. Li, and P. Georgiou, "Offline deep reinforcement learning and off-policy evaluation for personalized basal insulin control in type 1 diabetes," *IEEE Journal of Biomedical and Health Informatics*, 2023.
- [42] A. Irpan, K. Rao, K. Bousmalis, C. Harris, J. Ibarz, and S. Levine, "Off-policy evaluation via off-policy classification," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [43] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *Advances in neural information processing systems*, vol. 30, 2017.
- [44] S. M. Lundberg, B. Nair, M. S. Vavilala, M. Horibe, M. J. Eisses, T. Adams, D. E. Liston, D. K.-W. Low, S.-F. Newman, J. Kim *et al.*, "Explainable machine-learning predictions for the prevention of hypoxaemia during surgery," *Nature biomedical engineering*, vol. 2, no. 10, pp. 749–760, 2018.
- [45] E. Tjoa and C. Guan, "A survey on explainable artificial intelligence (xai): Toward medical xai," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 11, pp. 4793–4813, 2020.
- [46] M. Ali, H. Karimipour, and M. Tariq, "Integration of blockchain and federated learning for internet of things: Recent advances and future challenges," *Computers & Security*, vol. 108, p. 102355, 2021.
- [47] M. Ali, F. Naeem, M. Tariq, and G. Kaddoum, "Federated learning for privacy preservation in smart healthcare systems: A comprehensive survey," *IEEE journal of biomedical and health informatics*, vol. 27, no. 2, pp. 778–789, 2022.
- [48] D. Zhou, Y. Zhang, A. Sonabend-W, Z. Wang, J. Lu, and T. Cai, "Federated offline reinforcement learning," *arXiv preprint arXiv:2206.05581*, 2022.