# A Markov Decision Process Approach to Glucose Monitoring Allocation for ICU Patients

Anthony Maltsev
*Department of Computer Science*
*Stanford University*
Stanford, California
amaltsev@stanford.edu

Justin Jasper
*Department of Computer Science*
*Stanford University*
Stanford, California
jtjasper@stanford.edu

Jiayu Sui
*Department of Statistics*
*Stanford University*
Stanford, California
echosui@stanford.edu

*Abstract*—Critically ill patients in the intensive care unit (ICU) require frequent blood glucose monitoring, yet resource and operational constraints limit how often clinicians can obtain point-of-care measurements. This creates a complex sequential decision problem where each decision to invoke a blood glucose measurement must weigh the operational cost of the test against the clinical risk and uncertainty that arises from deferring the test. In this project, we model glucose-monitoring allocation as a Markov Decision Process (MDP) and develop offline reinforcement learning (RL) and behavioral cloning approaches to learn monitoring policies from historical ICU data. Using a curated MIMIC-III dataset of glucose measurements and insulin events, we construct a feature-based state representation that approximates the latent glycemic state under partial observability and we define a clinically informed reward structure that captures measurement costs and proxies for hyperglycemic risk [1]. We evaluate heuristic, supervised, and Deep Q-Learning policies on historical trajectories and show that learned policies produce clinically sensible measurement strategies despite sparse and noisy observations. Our results demonstrate that MDP-based methods can support data-driven glucose-monitoring decisions in ICU settings, though they fall short of clinically informed policies.

## I. INTRODUCTION

All critically ill patients admitted to the intensive care unit (ICU) require regular blood glucose measurements, irrespective of a history of diabetes [2]. The physiological stress of critical illness can cause significant fluctuations in blood glucose levels that need to be monitored and managed. The current standard of care for glucose testing in the ICU is point-of-care blood glucose monitoring using a fingerprick or arterial/venous blood sample [3]. Due to time and resource constraints, clinicians must make strategic decisions on when to administer glucose checks to their patients, prioritizing the allocation of limited monitoring capacity to patients with the greatest clinical need. These decisions have direct implications for both patient safety and overall ICU workload efficiency. In this work, we develop a framework for modeling this problem as a sequential decision-making task using a cost-sensitive Markov Decision Process (MDP) and evaluate algorithmic policies that can support glucose-monitoring allocation under realistic ICU constraints.

### A. Reinforcement Learning and Its Application to Blood Glucose Management

Reinforcement learning (RL) provides a framework for addressing sequential decision-making problems by learning policies that optimize for long-term patient outcomes rather than isolated or myopic decisions. Blood glucose management is inherently a sequential decision-making problem: patient glucose levels evolve over time, respond stochastically to interventions, and must be managed under operational constraints and uncertainty. Hence, the paradigm of RL has been increasingly applied to this problem space ( [2], [5], [6]). An advantage of RL is that it allows clinicians to deliver personalized care based on individual clinical needs rather than relying on generalized guidelines.

### B. Formulating Glucose Monitoring as a Markov Decision Process

Markov Decision Processes (MDPs) provide the mathematical structure underlying reinforcement learning and consist of four components: a state space, an action space, a transition model, and a reward function. To capture the sequential nature of ICU glucose-monitoring decisions, we formulate the problem as an MDP. This framework allows us to model how patient physiology evolves over time, how clinicians allocate limited monitoring capacity, and how these choices influence future clinical outcomes. Furthermore, using an MDP allows for the optimization of patient outcomes under clinical uncertainty.

The use of Markov Decision Processes (MDPs) to inform clinical decision-making in glucose management for patients equipped with continuous glucose monitors (CGMs) has gained increasing traction in recent years( [7]–[10]). For our project, we generalize and adapt these methodologies to a distinct clinical context: determining the optimal schedule for intermittent blood glucose testing for critically ill patients in the ICU. In contrast to previous studies, CGMs are not commonly used in our chosen clinical setting due to accuracy concerns. Instead, ICU clinicians must rely on point-of-care glucose testing, which results in a sparsity of ground-truth data. Furthermore, ICU clinicians have operational constraints on how frequently they can monitor each patient, forcing

clinicians to strategically allocate their glucose monitoring capacity. The scarcity of ground-truth data plus operational limits on patient monitoring presents unique modeling challenges not addressed in prior work.

## II. RELATED WORK

Decision-making under partial observability has been widely studied in healthcare, where clinicians must act under uncertain patient states and costly or infrequent measurements. Partially Observable Markov Decision Processes (POMDPs) have been used in chronic disease management, medical screening, and personalized follow-up scheduling, highlighting the central value-of-information (VoI) tradeoff between measurement costs and improved belief accuracy [11], [12].

More closely related to our application, several works model glucose management as a sequential decision problem. Prior POMDP-based approaches have optimized glucose monitoring frequency by treating each measurement as an information-gathering action that reduces uncertainty in latent glucose dynamics [13]. Other studies apply POMDP or reinforcement learning (RL) methods to insulin dosing and closed-loop glucose control under partial observability [14], [15]. Recent work also highlights challenges arising from delayed and prolonged insulin effects, motivating uncertainty-aware formalisms such as PAE-POMDPs [16].

Methodologically, exact POMDP planning is computationally intractable in large real-world domains, motivating approximate dynamic programming, online tree search, and belief-state compression techniques [17], [18]. Practical software toolkits further demonstrate that POMDPs can be deployed for real clinical time-series data despite limited observability [19].

In contrast to full POMDP planning, our work adopts a belief-approximation MDP approach: although the underlying glucose dynamics are only partially observable, we construct a compact feature-based state representation that summarizes recent measurements, trends, and elapsed time, and treat this representation as approximately Markov. Within this MDP framework, we study the problem of *when* to measure glucose under uncertainty and operational cost constraints. Unlike prior work that focuses on optimizing treatment actions such as insulin dosing, we focus specifically on the *measurement scheduling* problem and the associated value-of-information tradeoff in safety-critical ICU glucose management.

## III. DATASET & FEATURES

We use the *Curated Data for Describing Blood Glucose Management in the Intensive Care Unit* dataset from PhysioNet, which aggregates over 500,000 blood glucose readings and more than 140,000 insulin administration events for nearly 9,600 ICU patients from the MIMIC-III database [1]. Each entry corresponds to either a glucose reading or an insulin intervention, with accompanying timestamps and contextual metadata.

The curated data includes patient identifiers, ICU admission information, glucose test types and results, insulin administration types and dosage. For our study, we focus on the following key variables:

1) **TIMER**: Timestamp of the event—either the START-TIME of an insulin entry or the GLCTIMER of a glucose measurement. We use TIMER to chronologically order events within an admission.
2) **EVENT**: The insulin administration type (subcutaneous bolus, intravenous bolus, or infusion). In our environment, insulin events are treated as proxies for potential hyperglycemic episodes.
3) **GLC**: Blood glucose value in mg/dL, used to update the latent glucose state when a measurement is observed.
4) **GLCSOURCE**: The method of glucose measurement (fingerstick vs. laboratory analyzer), which influences the reliability of observations.

To construct a sequential decision-making environment, we convert each patient's event stream into a discretized time series with 15-minute time steps. At each step, we aggregate all events occurring within the corresponding window. Glucose measurement events directly update the observed glucose value and associated statistics in the state, while insulin entries serve as signals that alarm a glucose measurement may be needed. This assumption is motivated by the fact that insulin entries in the ICU are typically administered in response to elevated glucose levels, making them a reasonable proxy for unobserved hyperglycemic states.

## IV. METHODS

### A. State Representation and Reward Function

We formulate the problem of blood glucose monitoring for ICU patients as a partially observable, sequential decision making process, where every fixed timeframe (15 minutes) a choice must be made to administer some blood glucose test or not. The true patient glycemic state is only observable through noisy, sparse measurements taken via finger pricks or laboratory analyses at irregular intervals. We instead approximate our belief state using information about the last observed measurement, summary statistics of all previous measurements, and information about the time since the last measurement. We treat this belief state feature vector as though it was the Markov state for our decision making policies. The state space is comprised of an 8-dimensional feature vector with the following features:

1) Most recent glucose measurement
2) Mean glucose over observation window
3) Standard deviation of glucose
4) Minimum glucose in window
5) Maximum glucose in window
6) Rate of glucose change per hour
7) Fraction of time spent in ideal glucose range (70-180 mg/dL)
8) Time since last measurement

We finalized the feature vector by experimenting with how adding and removing state features impacted the performance of the learned policies relative to the base heuristics. Incorporating additional features derived from the *Curated Data for Describing Blood Glucose Management in the Intensive Care Unit* dataset (including "time since last insulin dosage", "last insulin dosage amount", and "insulin type") worsened the performance of the learned policies. Moreover, removing state feature vectors also resulted in worsened performance, giving us justification for our current state feature vector formulation.

The action space has 3 choices: do nothing, take a fingerprick, or take a laboratory analysis (each with different associated costs). Transitions are derived retrospectively purely from historical data based on how monitoring decisions were made by clinicians. We evaluate our policies also based on historical trajectories, simply capturing the reward for different policies when evaluated on historic trajectories since we cannot collect online rollouts for our policies.

We design the reward function to balance clinical tradeoffs between measurement costs, information value, and patient health. The reward function contains several distinct terms which are summed to provide a total reward. First, there are default costs $c(a)$ for different actions that represent the actual clinical resource costs of performing any measurements. Second, we introduce a reward term $r_t^{\text{treat}}$ that penalizes missing measurements before historic interventions. The reasoning for this decision is that medical interventions/bolus injections are a proxy for hyperglycemic state in a patient. These hyperglycemic states should be observed via some measurement in a short window by performing some measurements close to the historically observed intervention. Specifically, we let the reward at time $t$ be given by $r_t$:

$$r_t = c(a_t) + r_t^{\text{treat}}$$

$$c(a) = \begin{cases} 0 & a = \text{wait} \\ -1 & a = \text{fingerprick} \\ -3 & a = \text{lab analysis} \end{cases}$$

$$r_t^{\text{treat}} = \begin{cases} -20, & L_t > 90, \\ 0, & \text{otherwise.} \end{cases}$$

Here, $L_t$ is the time since the last measurement in minutes. The -20 penalty is applied only when an insulin intervention occurs more than 90 minutes after the last measurement, reflecting a missed opportunity to detect a worsening glycemic state before treatment. We experimented with other terms that intuitively would reward measurements, because without terms that reward penalties, the obviously optimal policy would be the simple heuristic of always waiting. We tried adding an uncertainty penalty term, which would penalize waiting too long between measurements. However, we found that our learned policies performed better when we used only the two reward components above.

## B. Policies

Given this representation of the problem, we constructed several different policies: some hard-coded heuristic baselines, as well as a policy learned using Q-Learning (DQN), and a simple policy derived via behavioral cloning of the historical actions. First, we implement a policy $\pi_{\text{wait}}$ which deterministically performs no measurement in all states. Second, we implement a policy which deterministically measures every 3 hours, $\pi_{\text{threshold}}$. We selected a 3-hour interval because it aligns with a realistic monitoring frequency for ICU workflows and provides a simple, interpretable heuristic baseline. Third, we implement a heuristic one-step look ahead policy, which estimates the value of information gained by any measurement by adding a term that scales linearly with the time since the last measurement as well as a term that scales with the standard deviation of previous glucose measurements in the patient and choosing the action with the highest $-cost + VOI$.

Then, to learn a policy directly from historic patient trajectories, we implement a Deep Q-Learning (DQN) approach that treats the belief-state feature vector described above as the Markov state $s_t$. Because the data is fully offline, no environment interaction occurs during training; instead, we sample transitions $(s_t, a_t, r_t, s_{t+1})$ from the historical dataset. We parameterize the action-value function with a neural network $Q_\theta(s, a)$ and train it to satisfy the Bellman equation for the optimal Q-function:

$$Q^*(s_t, a_t) = r_t + \gamma \max_{a' \in \mathcal{A}} Q^*(s_{t+1}, a').$$

The neural network is 2 layer MLP with hidden dimensions 128. Following the standard DQN training recipe, we maintain a separate target network $Q_\theta^{\text{target}}$ whose parameters are periodically updated from $\theta$. For each sampled transition, we construct the target

$$y_t = r_t + \gamma \max_{a' \in \mathcal{A}} Q_\theta^{\text{target}}(s_{t+1}, a'),$$

and minimize the loss

$$\mathcal{L}(\theta) = \left( Q_\theta(s_t, a_t) - y_t \right)^2.$$

We also implement a simple behavioral cloning algorithm, which bypasses our occasionally finicky reward function and directly trains against the historical actions taken. This treats the historical measurements as a ground truth or expert demonstrations, which is a reasonable assumption since those measurements were performed by trained clinical staff. We implement this as a simple supervised learning classification problem, to which we apply a simple 2 layer MLP with hidden dimensions 64 and 32. This produces logits for each action in the action space, which we convert into a probability distribution over actions using the softmax function. During evaluation, we sample from this resulting probability distribution.

## V. RESULTS

We measure the performance of our policies by checking how they perform when the actions selected by the policy

TABLE I
POLICY EVALUATION STATISTICS

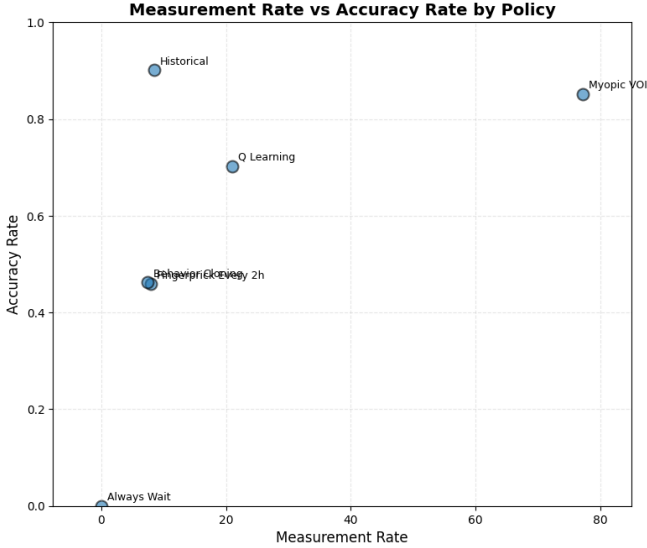| Policy | Measurement Rate ↓ | Interventions w/ Associated Measurement ↑ |
|---|---|---|
| Always Wait | 0.00 | 0.000 |
| Fingerprick Every 3h | 7.97 | 0.458 |
| Myopic VOI | 77.25 | 0.851 |
| QLearner | 21.75 | 0.703 |
| BCPolicy | 7.48 | 0.462 |
| Historical | 8.47 | 0.901 |



Fig. 1. Tradeoff between policy accuracy and measurement rate.

are evaluated on the stream of historic events for patients. We include two metrics: the average number of measurements taken per day by a policy and the number of treatment interventions that have an associated measurement before the intervention (as defined above). We report this as an accuracy rate out of all total treatments. These two metrics capture the clinical tradeoff between measurement cost and observation of hyperglycemic states for patients. The ideal policy would minimize the measurement rate (conserving clinical resources) while maximizing the number of treatments that have an associated measurement (encouraging knowledge of patient glycemic levels during critical moments). The performance of each policy we tested is reported in table I.

We see that our learned policies are only marginally better than even a simple logical policy which measures once every 3 hours (roughly matching the measurement rate of historical measurements). This reflects the difficulty of the problem and the sparsity of the signal in the dataset. It is difficult to recover high accuracy without measuring too much. The Q Learning policy measured more than twice as much as the historical rate while recovering only 80% of the accuracy of the historical policy. This highlights the difficulty of predicting the patients'

glycemic states, especially given a very sparse signal of previous measurements. The policy that directly learns from the historical measurement actions performs very similarly to the simple policy of regularly scheduled measurements. It is likely that more clinical information is needed to inform measurement actions beyond simple heuristics which are not available in the dataset we are using.

All of the measurement actions taken by the learned and heuristic policies are finger pricks actions because they have a lower cost. Our state space models does not capture the differing levels of uncertainty between measurements via finger prick versus more sophisticated lab analyses. Different methods of modeling uncertainty in measurements would help differentiate measurement actions.

It is also worth noting that the historical policy has a significant correlation between measurements and interventions as interventions are usually informed by a measurement. Therefore, there might be other periods of hyperglycemia in patients that are not captured in the data.

## VI. DISCUSSION

### A. Implications of Our MDP Formulation

In this paper, we developed a framework for sequential glucose monitoring allocation in the ICU and evaluated several policies, including heuristics, behavioral cloning, and Deep Q-Learning. Our work highlights both the potential and limitations of applying reinforcement learning and MDP-based methods in clinical decision-making under partial observability.

For this project, a notable challenge arose from the sparsity of ground-truth glucose measurements in real-world ICU datasets. Negative consequences from a lack of glucose measurements only occur when the patient's true glycemic state is poor. However, the majority of the true glucose trajectory is unobserved due to infrequent point-of-care testing. Consequently, penalties computed using the last-known measurements may not accurately reflect the patient's actual physiological state, particularly when the time elapsed since the last measurement is long. This limitation can lead to learned policies that underestimate the risk of hyperglycemia or overestimate the safety of waiting.

The underlying clinical process is fundamentally a POMDP: the true glucose state is only intermittently observed through sparse measurements. In this work, however, we approximate it as a belief-state MDP by constructing a feature-based state representation that encodes the history of measurements, recent trends, and time elapsed since the last observation, and then treating this feature vector as Markov. Within this MDP formulation, we handcraft the reward signal. The action of waiting is assigned a reward of 0, since deferring a glucose check incurs no direct operational cost. Point-of-care finger-prick testing is modeled with a small negative reward of $(-1)$ to represent its modest but non-negligible time and resource burden. In contrast, sending a sample for a full laboratory analysis is assigned a larger penalty of $(-3)$, reflecting its

substantially higher cost. These numerical values are hand-specified and would benefit from tuning by a clinician with more domain expertise in ICU treatment.

Using the defined reward function, the Q-network implicitly learns to trade off between measurement costs and glycemic safety through trial-and-error on historical trajectories. This approach favors practical scalability and data efficiency over the exact optimality provided by POMDP solvers, enabling application to large retrospective datasets. In effect, the network learns policies that capture temporal patterns in glucose dynamics without requiring an explicit model of the underlying physiological processes. However, this method does not recover Pareto optimal performance. This aligns with other lines of work in glucose monitoring which suggest that clinically informed policies outperform purely data driven policies [20].

### B. Adjustments to Model Complexity

Contrary to our initial beliefs, we discovered that our relatively low-dimensional features were sufficient to recover meaningful policies. When we attempted to add additional features beyond the 8 we currently include (e.g. time since last insulin dosage, most recent insulin dose amount, insulin type), we noticed that the performance of the learned policies actually worsened, likely due to the model overfitting sparse insulin data — the vast majority of time steps have no insulin dosing.

## VII. Conclusion

Our findings suggest that feature-based approximations and offline RL are able to produce clinically sensible policies for glucose monitoring allocation, even under sparse and noisy observations. The results highlight that compact state representations and relatively simple function approximations can still capture the essential clinical trade-offs involved in weighing risk, uncertainty, and measurement burden. Overall, this study demonstrates the feasibility of applying RL-based approaches to the problem of ICU glucose monitoring allocation and provides insight into the trade-offs between model complexity, observability, and real-world operational constraints. However, our results also suggest that a more clinically informed model would likely yield better results as our learned policies are only marginally better than simple heuristic approaches.

## VIII. Future Work

Based on the results of our project, there are many potential avenues for future development:

- The project would benefit from a clinically informed transition model, which would enable collecting policy rollouts. This environment would enable us to use a variety of different methods that improve Q-Learning. For example, we could use DAgger. There is support in the literature for using clinically informed models, which are more interpretable than purely data-driven models, especially in regimes with low-data or sparse signal [20].

- Going forward, we could instead represent the problem of glucose measurement allocation as a POMDP by better estimating the temporal dynamics of blood glucose changes using a transformer-based sequence model. This would enable the agent to maintain a more informative belief state over the patient's underlying physiology. Despite sparse observations, the agent would be empowered to select measurement times based on predicted future risk rather than solely on currently observed values. In order to accomplish this goal, we could train our new transformer-based sequence model on data from continuous glucose monitors, as shown in [21].

- The framework that we described in this paper could be extended beyond glucose monitoring to other analogous resource-sensitive measurement allocation tasks, like administering labs, vitals, and clinical imaging.

## IX. Contributions

The contribution by each teammate to the project were as follows:

- Anthony helped formulate the MDP model; wrote the dataset preprocessing and model representation; wrote the Q-Learning code; wrote the Behavior Cloning code; experimented with different reward functions; ran experiments with Q Learning and behavioral cloning including the final training runs; wrote the methods and results sections of the report.

- Justin helped formulate the MDP model and reward function; wrote baseline heuristics policies and simulation environment; wrote the abstract, introduction, discussion, conclusion, and references for the final paper; researched various articles/methods and datasets that could be used for modeling the allocation of point-of-care glucose monitoring devices.

- Jiayu helped formulate the problem statement and MDP model; conducted a literature review; ran some experimentation with Q Learning policies.

Our team spent an additional 30 hours on the project conducting a more thorough literature review, reading about additional methods for offline learning like imitation learning and implementing behavioral cloning and running more experiments. When finalizing our state feature vector, we ran multiple ablation studies to test how removing state features altered model performance. After completing the main project, we explored partial observability by brainstorming the design of transformer-based predictors of glucose dynamics.

## References

[1] Robles Arévalo, A., Mateo-Collado, R., & Celi, L. A. (2021). Curated Data for Describing Blood Glucose Management in the Intensive Care Unit (version 1.0.1). PhysioNet. $RRID : SCR_007345$. https://doi.org/10.13026/517s-2q57

[2] D. Juneja, D. Deepak, and P. Nasa, "What, why and how to monitor blood glucose in critically ill patients," World J. Diabetes, vol. 14, no. 5, pp. 528–538, May 2023, doi: 10.4239/wjd.v14.i5.528.

[3] R. Sreedharan, A. Martini, G. Das, N. Aftab, S. Khanna, and K. Ruetzler, "Clinical challenges of glycemic control in the intensive care unit: A narrative review," World J. Clin. Cases, vol. 10, no. 31, pp. 11260–11272, Nov. 2022, doi: 10.12998/wjcc.v10.i31.11260.

[4] I. S. Jacobs and C. P. Bean, "Fine particles, thin films and exchange anisotropy," in *Magnetism*, vol. III, G. T. Rado and H. Suhl, Eds., New York: Academic, 1963, pp. 271–350.

[5] K.-L. Yau et al., "Reinforcement learning models and algorithms for diabetes management," IEEE Access, pp. 1–1, 2023, doi: 10.1109/ACCESS.2023.3259425.

[6] H. Emerson, M. Guy, and R. McConville, "Offline reinforcement learning for safer blood glucose control in people with type 1 diabetes," J. Biomed. Inform., vol. 142, p. 104376, 2023, doi: 10.1016/j.jbi.2023.104376.

[7] L. H. Dicker et al., "Continuous blood glucose monitoring: A Bayes-hidden Markov approach," Statistica Sinica, vol. 23, no. 4, pp. 1595–1627, 2013. [Online]. Available: http://www.jstor.org/stable/24310814

[8] S. Wang and W. Gu, "An improved strategy for blood glucose control using multi-step deep reinforcement learning," in *Proc. 16th Int. Conf. Bioinformatics and Biomedical Technology (ICBBT)*, New York, NY, USA: ACM, 2024, pp. 196–203, doi: 10.1145/3674658.3674689.

[9] Y. Zhang, H. Wu, B. T. Denton et al., "Probabilistic sensitivity analysis on Markov models with uncertain transition probabilities: An application in evaluating treatment decisions for type 2 diabetes," Health Care Manag. Sci., vol. 22, pp. 34–52, 2019, doi: 10.1007/s10729-017-9420-8.

[10] J. O. Ferstad, E. B. Fox, D. Scheinker, and R. Johari, "Learning explainable treatment policies with clinician-informed representations: A practical approach," arXiv preprint arXiv:2411.17570, 2024. [Online]. Available: https://arxiv.org/abs/2411.17570

[11] O. Alagoz, H. Hsu, A. J. Schaefer, and M. S. Roberts, "A review of decision models in chronic disease screening and treatment," Operations Research, vol. 58, no. 5, pp. 949–964, 2010.

[12] J. Yu and D. Bertsekas, "Partially observable Markov decision processes for disease progression and treatment decisions," in *Proc. IEEE EMBS*, 2004, pp. 300–303.

[13] I. Thapa, E. Rao, and W. Cai, "Dynamic glucose monitoring with partially observable Markov decision processes," Stanford Univ., Tech. Rep., 2019, AA228/CS238 Course Project.

[14] C. Liu and A. Nanduri, "Optimization of insulin dosing in diabetic patients with POMDPs," Stanford Univ., Tech. Rep., 2019, AA228/CS238 Course Project.

[15] Z. Yang et al., "An improved strategy for blood glucose control using multi-step deep reinforcement learning," arXiv preprint arXiv:2403.07566, 2024. [Online]. Available: https://arxiv.org/abs/2403.07566

[16] S. Jaszczur, F. Johansson, and M. Wattenberg, "On the challenges of using reinforcement learning in precision drug dosing: Delay and prolongedness of action effects," arXiv preprint arXiv:2301.00512, 2023. [Online]. Available: https://arxiv.org/abs/2301.00512

[17] D. Silver and J. Veness, "Monte-Carlo planning in large POMDPs," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2010.

[18] W. S. Lovejoy, "A survey of algorithmic methods for partially observed Markov decision processes," *Annals of Operations Research*, vol. 28, no. 1, pp. 47–66, 1991.

[19] M. Kochenderfer et al., "The `pomdp` package in R," *The R Journal*, 2024.

[20] J. Ferstad et al., "Learning Explainable Treatment Policies with Clinician-Informed Representations: A Practical Approach," *Proceedings of Machine Learning for Health (ML4H)*, 2024

[21] H. Emerson, M.Guy, R. McConville, "Offline reinforcement learning for safer blood glucose control in people with type 1 diabetes," *Journal of Biomedical Informatics*, Volume 142, 104376, ISSN 1532-0464, 2023. [Online]. Available: https://doi.org/10.1016/j.jbi.2023.104376.