

# IMPLICIT BAYESIAN MARKOV DECISION PROCESS FOR RESOURCE-EFFICIENT EXPERIMENTAL DESIGN IN DRUG DISCOVERY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In drug discovery, researchers make sequential decisions to schedule experiments, aiming to maximize probability of success towards drug candidates while simultaneously minimizing expected costs. However, such tasks pose significant challenges due to complex trade-offs between uncertainty reduction and allocation of constrained resources in a high-dimensional state-action space. Traditional methods based on simple rule-based heuristics or domain expertise often result in either inefficient resource utilization due to risk aversion or missed opportunities arising from reckless decisions. To address these challenges, we developed an Implicit Bayesian Markov Decision Process (IB-MDP) algorithm that constructs an implicit MDP model of the environment’s dynamics by integrating historical data through a similarity-based metric, and enables effective planning by simulating future states and actions. To enhance the robustness of the decision-making process, the IB-MDP also incorporates an ensemble approach that recommends maximum likelihood actions to effectively balance the dual objectives of reducing state uncertainty and optimizing expected costs. Our experimental results demonstrate that the IB-MDP algorithm offers significant improvements over traditional rule-based methods by identifying optimal decisions that ensure more efficient use of resources in drug discovery.

## 1 INTRODUCTION

In drug discovery, strategic planning and selection of experiments play a pivotal role in impacting the pace and expenses of R&D activities. The identification of potential drug candidates requires conducting numerous assays at various stages of preclinical studies. The process often begins with limited information, creating significant challenges for achieving optimized outcomes due to time and budget constraints. Optimizing the use of resources to achieve targeted goals within these limitations is among the most demanding tasks in creating effective Research Operation Plans (ROP). Conventional approaches, often relying on simple rule-based heuristics or domain expertise, struggle to adapt as new data emerges and typically fail to address state, model, and parameter uncertainties effectively. Consequently, this results in suboptimal decision-making and inefficient allocation of resources Puterman (2014).

To address these challenges, we propose the **Implicit Bayesian Markov Decision Process (IB-MDP)** algorithm, a *model-based* approach that constructs an implicit model of the environment’s dynamics by integrating historical data through a distance-based similarity metric. Unlike traditional MDP methods that require explicit modeling of transition probabilities, the IB-MDP leverages historical data to build a flexible model of the environment without the need for precise parameterization Rainforth et al. (2024); Bellet et al. (2013). This implicit model captures complex, nonlinear relationships within the data manifold, enabling efficient planning by simulating future states and actions Alagoz et al. (2010).

Moreover, to improve the robustness and reliability of decision-making, we incorporate an ensemble approach into the IB-MDP. Ensemble methods aggregate multiple policies derived from independent algorithm runs, reducing variance and mitigating bias in policy estimation Dietterich (2000); Osband

054 et al. (2016); Zhou (2012). This approach ensures more stable and generalizable policies, particu-  
055 larly in high-dimensional and resource-constrained environments Lakshminarayanan et al. (2017).  
056

057 Our algorithm is demonstrated in the context of assay scheduling and ROP optimization, where it  
058 significantly improves resource utilization and decision quality compared to traditional heuristic-  
059 based approaches. The IB-MDP framework is broadly applicable to various resource-constrained  
060 decision-making tasks in drug discovery, making it a valuable tool for optimizing sequential deci-  
061 sions in preclinical studies.

### 062 **Summary of Contributions:**

- 063 • We introduce the Implicit Bayesian Markov Decision Process (IB-MDP), a model-based  
064 algorithm that integrates historical data using a distance-based similarity metric within the  
065 MDP framework, enabling efficient planning in sequential decision-making tasks.  
066
- 067 • We incorporate an ensemble approach to enhance policy estimation, providing theoretic-  
068 al justification for its effectiveness in variance reduction, bias mitigation, and improved  
069 generalization.
- 070 • We validate our approach through experiments in assay scheduling, demonstrating signif-  
071 icant improvements in resource utilization and decision quality over traditional heuristic  
072 methods. However, our algorithm is broadly applicable across a range of decision-making  
073 problems.

## 074 2 RELATED WORK

075 The optimization of decision-making under uncertainty has been a central focus in various domains,  
076 including drug discovery.

077 **Markov Decision Processes and Model-Based Reinforcement Learning:** Markov Decision Pro-  
078 cesses (MDPs) provide a mathematical framework for modeling sequential decision-making where  
079 outcomes are partly random and partly under the control of a decision-maker (Puterman, 1994;  
080 2014). In drug discovery, MDPs have been applied to tasks such as clinical trial optimization (Ben-  
081 nett & Hauser, 2013; Eghbali-Zarch et al., 2019; Abbas et al., 2007; Fard et al., 2018). However,  
082 applying MDPs to experimental scheduling has remained limited due to the difficulty of accurately  
083 specifying transition probabilities and reward functions. Model-based reinforcement learning (RL)  
084 offers an alternative by learning models of the environment to improve sample efficiency and plan-  
085 ning accuracy (Sutton, 2018; Kaiser et al., 2019). In the drug discovery field, model-based RL  
086 has been used for molecule generation (Wang et al., 2021; Bengio et al., 2021; You et al., 2018;  
087 Zhou et al., 2019), synthesis planning (Segler et al., 2018), and experimental design (Schneider  
088 et al., 2020). These methods typically require accurate environment models, a challenge in high-  
089 dimensional and complex biological systems such as those found in preclinical studies.

090 **Incorporating Historical Data and Similarity Metrics:** Leveraging historical data is crucial for  
091 improving decision-making in contexts with limited experimental data. Bayesian approaches, in-  
092 cluding Bayesian reinforcement learning and optimization, maintain a posterior distribution over pa-  
093 rameters or value functions, updating beliefs based on new data (Ghavamzadeh et al., 2015; Shahriari  
094 et al., 2015). In drug discovery, Bayesian optimization has been applied to optimize molecular prop-  
095 erties (Griffiths & Hernández-Lobato, 2020; Gómez-Bombarelli et al., 2018), but such methods are  
096 often less effective in sequential decision-making scenarios. Using similarity metrics within MDPs  
097 can further enhance the integration of historical data. Kernel-based methods, which use similar-  
098 ity functions to generalize across states, have been explored in reinforcement learning to estimate  
099 transition dynamics more accurately (Ormonéit & Sen, 2002; Kveton & Theodorou, 2012; Xu  
100 et al., 2007). Our approach extends this by incorporating a variance-normalized distance metric to  
101 dynamically integrate historical data into the MDP transition function.  
102

103 **Ensemble Methods in Reinforcement Learning (RL):** Ensemble methods have gained popularity  
104 for improving the robustness and reliability of decision-making. By aggregating multiple models  
105 or policies, ensemble techniques reduce the variance and bias inherent in individual estimates (Di-  
106 etterich, 2000; Osband et al., 2016; Wiering & Van Hasselt, 2008; Zhou, 2012). In RL, ensemble  
107 methods are particularly effective in improving exploration and generalization, as demonstrated by  
their successful application in model-based RL (Lakshminarayanan et al., 2017). Our IB-MDP

framework incorporates ensemble methods to enhance decision robustness, where multiple policies derived from independent algorithm runs are aggregated to produce more reliable decision paths.

**Applications in ADME Studies and Comparison to Existing Methods:** ADME studies focus on the absorption, distribution, metabolism, and excretion (ADME) of drugs to understand their pharmacokinetic properties and impact on effectiveness and safety (Hoffman, 1998; Hoffman et al., 2004; Hughes et al., 2011). Decision-making in ADME studies requires balancing information gain with constrained resources. Although RL has been applied to clinical trial optimization (Coronato et al., 2020; Escandell-Montero et al., 2014; Martín et al., 2020), its application to preclinical experimental scheduling is underexplored. Our IB-MDP framework addresses this gap by providing a flexible and scalable decision-making approach that integrates historical data and real-time experimental results. Unlike traditional methods, the IB-MDP does not require manual specification of transition probabilities, instead leveraging a similarity-based metric to model the environment’s dynamics implicitly. This, combined with the ensemble approach, distinguishes our method from existing techniques, ensuring both adaptability and robustness in decision-making.

### 3 A SEQUENTIAL DECISION-MAKING PROBLEM STATEMENT

In ADME studies, a primary challenge is the optimal scheduling of multiple experimental assays that contribute to evaluating a drug’s ADME profile. Critical ADME assays for central nervous system (CNS) drugs involve assessing whether the drug acts as a substrate for transporters such as P-glycoprotein (PgP) and Breast Cancer Resistance Protein (BCRP). The goal is to plan in vitro PgP and BCRP assays to maximize information gain towards the drug’s brain penetration potential, which can be evaluated through in vivo unbound brain-to-plasma partition coefficient ( $k_{puu}$ ), while minimizing operational costs and adhering to resource limitations.

The focus here is on reducing state uncertainty, which refers to the incomplete knowledge about the final target feature, such as  $k_{puu}$ , rather than model uncertainty. Ensuring these features fall within desirable ranges is key to determining a drug’s efficacy and safety.

The problem can be formulated as finding an optimal policy  $\pi^*$  that minimizes cost and reduces state uncertainty, subject to constraints ensuring the likelihood of achieving experimental outcomes.

This can be expressed as:  $\min_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \gamma^t R(s_t, \pi(s_t)) \right]$  subject to:

1. State uncertainty at the terminal stage  $\mathcal{H}(s_T)$  must be below a threshold  $\epsilon$ :  $\mathcal{H}(s_T) \leq \epsilon$ .
2. The likelihood of achieving desirable outcomes  $\mathcal{L}(s_T)$  must exceed a minimum value  $\tau$ :  $\mathcal{L}(s_T) \geq \tau$ .
3. At each intermediate step  $t$ , the likelihood  $\mathcal{L}(s_t)$  must also exceed the threshold  $\tau$ :  $\mathcal{L}(s_t) \geq \tau, \quad \forall t = 0, \dots, T - 1$ .

Checking the likelihood at each intermediate step ensures that the decision-making process stays aligned with the final goal, maintaining a high probability of achieving desired experimental outcomes. This dynamic constraint helps the policy continuously adapt as new data emerge, enforcing the likelihood requirement throughout the MDP decision process and preventing early decisions from compromising long-term objectives. By consistently ensuring both cost efficiency and target feature accuracy, the policy remains robust and focused until the end of the search, aligning with the principles of constraint optimization in MDP frameworks.

## 4 IMPLICIT BAYESIAN MARKOV DECISION PROCESS (IB-MDP) FOR RESOURCE-EFFICIENT DECISION MAKING

### 4.1 FRAMEWORK DESCRIPTION

The IB-MDP algorithm is designed to optimize experimental scheduling in resource-constrained settings by leveraging historical data through a distance-based similarity metric. This framework aims to strategically select assays to minimize costs and maximize information gain, particularly in high-dimensional decision spaces, such as assays scheduling in preclinical pharmacokinetics and pharmacodynamics (PKPD) space.

The algorithm starts with a partially known initial state and a collection of potential experimental configurations (i.e., action sets in an MDP framework). As it explores the state-action space, the IB-MDP dynamically adjusts its strategy based on emerging evidence, ensuring that the policies remain optimal under given constraints. By constructing an implicit model of the environment’s dynamics, IB-MDP eliminates the need for a parameterized transition probabilities. Instead, the transition dynamics are inferred from historical data using a variance-normalized similarity metric. This method significantly reduces computational complexity while retaining the flexibility to refine decisions as new data become available.

A key feature of IB-MDP is its use of Monte Carlo Tree Search with Double Progressive Widening (MCTS-DPW), which enables efficient navigation through large state spaces without exhaustive data collection. This approach is particularly suited for experimental planning, where the goal is to balance exploration and exploitation in a computationally efficient manner. Furthermore, the framework incorporates a Bayesian sampling method, continuously refining the policy to incorporate new information, thus ensuring that the decision-making process adapts to changes in state uncertainty and target feature values over time.

To further enhance robustness and accuracy, IB-MDP integrates an ensemble method. By aggregating multiple policies generated from independent runs, the ensemble method mitigates inference bias and reduces variance, ensuring that the decision-making process is both reliable and adaptive. This combination of implicit modeling, dynamic policy adjustment, and ensemble learning offers a powerful tool for optimizing resource usage in complex experimental designs.

## 4.2 IB-MDP FORMULATION

The IB-MDP (Implicit Bayesian Markov Decision Process) framework can be defined as a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma \rangle$ , where:

- **States** ( $\mathcal{S}$ ): The state space represents the knowledge about the drug candidate or system at each decision point.
- **Actions** ( $\mathcal{A}$ ): A set of actions, where each action corresponds to selecting assays or experiments to perform.
- **Transition Function** ( $\mathcal{T}$ ): Transition probabilities between states, implicitly derived from historical data  $\mathcal{D}$  using a similarity-based metric.
- **Reward Function** ( $\mathcal{R}$ ): The reward function that penalizes resource costs and rewards uncertainty reduction and goal achievement.
- **Discount Factor** ( $\gamma$ ): A scalar discount factor that determines the present value of future rewards.

### 4.2.1 SIMILARITY WEIGHT FUNCTION

The transition function relies on a similarity weight  $w_i(s)$  for each historical data point  $D_{s_i}$ . This is computed based on the distance between the current state  $s$  and the historical data point  $D_{s_i}$  using a variance-normalized distance metric:  $w_i(s) = \exp(-\lambda_w \cdot d(s, D_{s_i}))$ ,

where  $\lambda_w$  is a scaling factor, and the distance metric  $d(s, D_{s_i})$  is:

$$d(s, D_{s_i}) = \sum_{k=1}^n \lambda_k \cdot \frac{(s_k - (D_{s_i})_k)^2}{\sigma_k^2},$$

with  $\lambda_k$  representing feature-specific scaling factors, and  $\sigma_k^2$  being the variance of the  $k$ -th feature in  $\mathcal{D}$ .

### 4.2.2 IMPLICIT TRANSITION MODELING VIA SAMPLING

In the IB-MDP framework, the transition function  $\mathcal{T}(s, a, s')$  is not explicitly defined through a known analytical function. Instead, it is implicitly modeled using a weighted sampling process that leverages historical data  $\mathcal{D}$  and the similarity weights vector  $W$ . This process allows the transition

to a new state  $s'$  to be based on past data points most similar to the current state  $s$ , according to a variance-normalized distance metric.

The transition probability from state  $s$  to state  $s'$  given action  $a$  is defined as:

$$P(s'|s, a) = \sum_{i=1}^N \frac{w_i(s)}{\sum_{j=1}^N w_j(s)} \cdot \mathbb{I}[s' = s \oplus \Delta s(a, D_{s_i})],$$

where:

- $w_i(s)$  is the similarity weight between the current state  $s$  and the historical data point  $D_{s_i}$ .
- $\Delta s(a, D_{s_i})$  is the change in state resulting from action  $a$  applied to the historical data point  $D_{s_i}$ .
- $\oplus$  denotes the state update operation, combining the current state  $s$  with the effect of action  $a$  based on the sampled historical data.
- $\mathbb{I}[\cdot]$  is an indicator function that ensures the state update conforms to the sampled transition.
- $N$  denotes the total number of historical data points.

This transition is realized through a weighted sampling process. Specifically, a historical state  $D_{s_{\text{sampled}}}$  is selected with probability proportional to its similarity weight  $w_i(s)$ . The sampling function  $\delta(W, s)$  selects a data point  $D_{s_{\text{sampled}}}$  from the historical dataset  $\mathcal{D}$ :  $D_{s_{\text{sampled}}} = \delta(W, s) \cdot \mathcal{D}$ , where  $\delta(W, s)$  uses the similarity weights  $W$  to sample a historical state  $D_{s_i}$  from  $\mathcal{D}$ . Once the historical state is sampled, the new state  $s'$  is updated as:  $s' = \beta(s, \mathcal{D}, a) = s \oplus \Delta s(a, D_{s_{\text{sampled}}})$ , where  $\Delta s(a, D_{s_{\text{sampled}}})$  represents the change in state resulting from action  $a$  applied to the sampled data point  $D_{s_{\text{sampled}}}$ . The  $\beta$  function is considered as the implicit Bayesian update for the state via sampling. This update mechanism allows the system to dynamically evolve by incorporating the effects of historical actions, without needing an explicit transition function.

#### 4.2.3 BAYESIAN UPDATE MECHANISM

After the transition to a new state  $s'$ , the similarity weights  $W$  are updated based on the new state and historical data. The updated weights  $W'$  are computed as:  $W' = \text{update}(W, s', \mathcal{D})$  where the update mechanism reflects how the similarity weights are adjusted based on the new state  $s'$ , the action  $a$ , and the historical data  $\mathcal{D}$ .

This process dynamically adjusts the state transition based on the updated belief about the system, providing a probabilistic framework for modeling uncertainties.

#### 4.2.4 REWARD FUNCTION

The reward function is defined to balance cost minimization and uncertainty reduction:

$$R(s, a) = \begin{cases} -\mathbf{c}(s, a) \cdot \boldsymbol{\lambda}, & \text{if } a \neq \text{eox}, \\ -M, & \text{if } a = \text{eox}, \end{cases}$$

where  $\mathbf{c}(s, a)$  represents the cost vector for action  $a$ ,  $\boldsymbol{\lambda}$  is a vector of trade-off parameters, and  $M$  is a large penalty for premature termination.

#### 4.2.5 STATE UNCERTAINTY AND LIKELIHOOD

Uncertainty in state  $s$ , denoted as  $\mathcal{H}(s)$ , is computed as:  $\mathcal{H}(s) = \frac{\sum_{i=1}^N w_i(s) (k_i - \bar{k}_w(s))^2}{\sum_{i=1}^N w_i(s)}$  where  $k_i$  is the value of the target feature in  $D_{s_i}$  and  $\bar{k}_w(s)$  is the weighted mean of the feature:  $\bar{k}_w(s) = \frac{\sum_{i=1}^N w_i(s) \cdot k_i}{\sum_{i=1}^N w_i(s)}$ . The likelihood of achieving the desired outcome is computed as:  $\mathcal{L}(s) = \frac{\sum_{i=1}^N w_i(s) \cdot \mathbb{I}[k_i \in [k_{\min}, k_{\max}]]}{\sum_{i=1}^N w_i(s)}$  where  $k_i$  is the value of the target feature  $k$  in the  $i$ -th historical data point,  $w_i(s)$  is the similarity weight for  $D_{s_i}$ , and  $\mathbb{I}$  is the indicator function.  $k_{\min}$  and  $k_{\max}$  are the user-defined lower and upper bound values of the target features.

#### 270 4.2.6 TERMINAL CONDITION

271  
272 The state  $s$  is considered terminal when:  $\mathcal{H}(s) \leq \epsilon$  and  $\mathcal{L}(s) \geq \tau$  where  $\epsilon$  is the state uncertainty  
273 threshold, and  $\tau$  is the likelihood threshold.

### 275 4.3 THE IB-MDP ALGORITHM

276  
277 The IB-MDP framework models sequential decision-making within resource-constrained environ-  
278 nments, such as drug discovery, where an optimal sequence of experiments needs to be chosen under  
279 uncertainty. The action set at each state, represented as the power set of available assays  $\mathcal{P}(A)$ ,  
280 may be constrained by the maximum number of assays  $m$  (i.e.,  $\mathcal{A}_m$ ). Transitions between states are  
281 modeled using historical data and similarity weights, as outlined in the formulation section.

#### 282 4.3.1 SOLVING IB-MDP WITH MCTS-DPW

283  
284 To solve the IB-MDP problem, we employ Monte Carlo Tree Search (MCTS) with Double Pro-  
285 gressive Widening (DPW). MCTS is a powerful search algorithm used to explore large state-action  
286 spaces by building a search tree through iterative simulations Browne et al. (2012). DPW is utilized  
287 to handle large and continuous action spaces by initially restricting the number of explored actions  
288 at each state and progressively widening the action set as more iterations are performed Couëtoux  
289 et al. (2011).

290 MCTS operates in four main steps: Selection, Expansion, Simulation, and Backpropagation. The  
291 Upper Confidence Bound (UCB) policy is used in the selection phase to balance exploration and  
292 exploitation:  $a = \arg \max_{a' \in \mathcal{A}(s)} \left( Q(s, a') + c \sqrt{\frac{\ln N(s)}{N(s, a')}} \right)$  where  $Q(s, a')$  is the estimated value of  
293 action  $a'$  in state  $s$ , and  $N(s)$  and  $N(s, a')$  represent the number of visits to state  $s$  and action  $a'$  in  
294 state  $s$ , respectively.

295  
296 The state transitions during the simulation phase are modeled implicitly via weighted sampling  
297 based on historical data, using the Bayesian update mechanism described earlier. This allows the  
298 IB-MDP to adapt dynamically to new information while maintaining computational efficiency.

#### 299 4.3.2 PARETO FRONT GENERATION

300  
301 After solving the IB-MDP problem and obtaining a set of optimal policies  $\pi_j^*$  from the MCTS-DPW  
302 runs, the next step is to generate the Pareto front. This front helps to discern optimal trade-offs  
303 between competing objectives, such as minimizing cost and reducing state uncertainty.

304  
305 The Pareto front consists of non-dominated points in the objective space, where each point repre-  
306 sents a policy that is optimal under certain constraints. Mathematically, this can be formulated as:  
307 minimize  $\{(\mathcal{C}(s), \mathcal{H}(s)) \mid s \in \mathcal{S}\}$  where  $\mathcal{C}(s)$  is the cost associated with state  $s$  and  $\mathcal{H}(s)$  represents  
308 state uncertainty. A state  $s'$  dominates state  $s$  if:  $\mathcal{H}(s') < \mathcal{H}(s)$  and  $\mathcal{C}(s') < \mathcal{C}(s)$  indicating that  
309  $s'$  has lower uncertainty and lower cost. The Pareto front is the set of states that are not dominated  
310 by any other state, ensuring that each point on the front represents a trade-off between cost and  
311 uncertainty.

### 312 4.4 ENSEMBLE METHOD FOR IB-MDP

313  
314 The ensemble method addresses the inherent variability in single runs of the IB-MDP algorithm,  
315 especially when using Monte Carlo Tree Search with Double Progressive Widening (MCTS-DPW).  
316 Stochasticity and sensitivity to initial conditions can lead to different Pareto fronts and optimal  
317 policies. By executing the IB-MDP algorithm  $N$  times, each run generates an optimal policy  $\pi_j^*$  and  
318 a corresponding Pareto front  $\mathcal{P}_j$ . Aggregating the results across multiple runs improves robustness,  
319 reduces bias, and enhances the reliability of decision-making.

320 **Advantages of Ensemble IB-MDP:** The ensemble IB-MDP methodology provides several advan-  
321 tages: **Improved Robustness:** By aggregating results from multiple simulations, the ensemble ap-  
322 proach reduces the impact of variability and randomness in individual runs, enhancing the stability  
323 of decision outcomes. **Bias Reduction:** Exploring diverse decision trajectories across runs mini-  
mizes inference bias and yields more accurate estimates of optimal policies. **Predictive Power:** The

ensemble method helps identify patterns across multiple runs. The aggregation of these results informs the construction of a Maximum Likelihood Action Sets Path (MLASP), providing a candidate metric that can guide decision-making by assessing decisions under different likelihood thresholds  $\tau$ .

#### 4.4.1 MAXIMUM LIKELIHOOD ACTION SETS PATH (MLASP)

The MLASP is a key outcome of the ensemble approach. It is constructed by identifying the most frequently occurring optimal action set at each uncertainty level  $u$  across multiple runs. Specifically, for each uncertainty level, the action set  $A_u^*$  that occurs most frequently is chosen as:  $A_u^* = \arg \max_A \sum_{j=1}^N \mathbb{I}(A \in \mathcal{P}_j(u))$  where  $\mathbb{I}$  is the indicator function that counts whether the action set  $A$  is part of the Pareto front  $\mathcal{P}_j(u)$  for the  $j$ -th run. By connecting all  $A_u^*$  across varying uncertainty levels, we form the MLASP, ensuring robust decision-making across different scenarios.

Figure 2 illustrates the construction of the MLASP by showing a histogram of actions proposed across an ensemble of 50 independent runs of the IB-MDP algorithm. For a given state at the uncertainty threshold  $\tau = 0.2$ , the height of each bar represents the frequency with which a particular action was chosen by the ensemble. The action that appears most frequently across the runs for that threshold is considered the majority-voted action and is included in the MLASP.

The histogram provides a clear visual of how often each action was selected, and the bar with the highest frequency corresponds to the optimal action under the given uncertainty. The key insight from this figure is how the ensemble method ensures more robust and stable decision-making by leveraging majority voting across multiple runs. It demonstrates that, even with variability across different runs, the MLASP method converges on a reliable decision, enhancing both predictive power and robustness in the decision-making process.

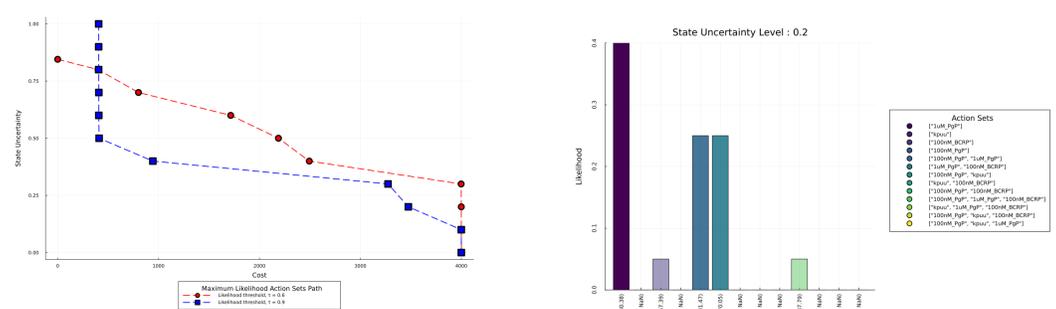


Figure 1: Exemplary Monetary-prioritized MLASP for two  $\tau$  thresholds for the data point with 50 ensembles,  $QSAR_{mrt}$  value of 1.56,  $QSAR_{100mL\_BCRP}$  of 0.87, and  $QSAR_{1uM\_Pgp}$  of 0.513. Different  $L$  likelihood result in different MLASP leading to distinct decision paths. A general trend is that the higher likelihood threshold  $\tau$  value, the lower left MLASP will be.

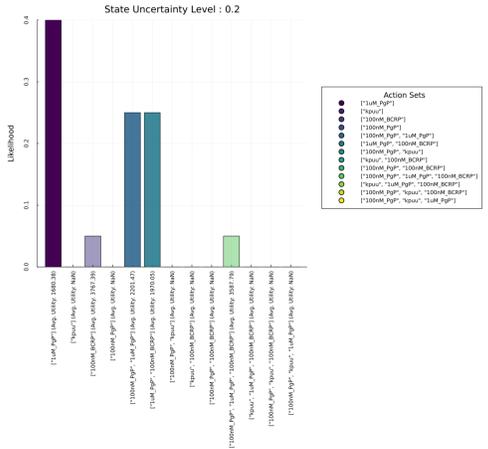


Figure 2: Example of 50 ensemble IB-MDP proposed action in a histogram plot for the state uncertainty level = 0.2, and  $\tau = 0.9$

#### 4.5 ALGORITHM

For the detailed and complete description of the algorithm, see Algorithm 1.

### 5 EXPERIMENTS

#### 5.1 EXPERIMENTAL SETUP

Our experimental setup utilizes a dataset of 220 compounds, each characterized by both *in silico* predictions and physical properties. The *in silico* features include Quantitative Structure-Activity

**Algorithm 1** Ensemble IB-MDP Algorithm

---

**Require:** Initial state  $s_0$ , historical data  $\mathcal{D}$ , similarity function  $W$ , Bayesian update function  $\beta$ , horizon  $H$ , number of iterations  $n_{\text{itr}}$ , number of ensemble runs  $N$

**Ensure:** Pareto front of state uncertainty vs. expected utility costs

- 1: Initialize an array  $\mathcal{P}$  to store Pareto fronts
- 2: **for**  $j = 1$  to  $N$  **do**
- 3:     Initialize MCTS-DPW tree with root node representing  $s_0$
- 4:     **for**  $i = 1$  to  $n_{\text{itr}}$  **do**
- 5:          $s \leftarrow s_0$
- 6:         **while** not terminal and within horizon  $H$  **do**
- 7:             Select action  $a$  using UCB policy:  $a = \arg \max_{a' \in \mathcal{A}(s)} Q(s, a') + c \sqrt{\frac{\ln N(s)}{N(s, a' )}}$
- 8:             Simulate next state  $s'$  using Bayesian update via sampling:  $s' = \beta(s, \mathcal{D}, a)$
- 9:             Update similarity weights  $W$  based on new state  $s'$
- 10:            Update tree with  $s'$  and reward  $R(s, a)$
- 11:             $s \leftarrow s'$
- 12:         **end while**
- 13:         Backpropagate rewards and update  $Q$  values along the path
- 14:     **end for**
- 15:      $\pi_j^* \leftarrow$  Extract optimal policy from tree
- 16:      $\mathcal{P}_j \leftarrow$  Compute Pareto front from  $\pi_j^*$
- 17:     Append  $\mathcal{P}_j$  to  $\mathcal{P}$
- 18: **end for**
- 19: **for** each uncertainty level  $u$  **do**
- 20:      $A_u^* = \arg \max_A \sum_{j=1}^N \mathbb{I}(A \in \mathcal{P}_j(u))$
- 21: **end for**
- 22: Construct Maximum Likelihood Action Sets Path (MLASP) from  $A_u^*$
- 23: **return** MLASP

---

Relationship (QSAR) predictions, such as QSAR<sub>1uM\_PgP</sub>, QSAR<sub>100nM\_BCRP</sub>, and QSAR<sub>mrt</sub>. In addition to these predictions, transporter activity data such as 100nM PgP, 1uM PgP, and 100nM BCRP are also considered. The financial and time costs associated with these transporter activities are estimated at \$400 per assay with a turnaround of 7 days, while  $k_{puu}$  measurements incur a higher cost of \$4000 and take 21 days. These values highlight the substantial resource investment required for these tests.

To generate the Maximum Likelihood Action Sets Path (MLASP), we conduct up to three parallel assays, allowing simultaneous experimental operations. This setup helps reduce state uncertainty more efficiently while maximizing information gain, both of which are essential for effective decision-making. A computational threshold of 10 is applied to assess state uncertainty, ensuring that the algorithm captures meaningful differences in uncertainty levels.

We employ the IB-MDP algorithm, integrated with a Monte Carlo Tree Search (MCTS) solver using Double Progressive Widening (DPW). This solver runs for 20,000 iterations, with an exploration constant of 5.0, 50 ensembles, providing a balance between exploring new actions and exploiting known outcomes.

The primary goal is to identify the actions that achieve the greatest reduction in state uncertainty, comparable to performing the final target assay ( $k_{puu}$ ), while minimizing both costs and resource use throughout the decision-making process.

### Experimental Computing Resources:

We performed the IB-MDP simulations on an Apple M1 Pro chip with 16GB of memory. For each ensemble run, with 100 iterations of the IB-MDP per  $\tau$  value, the estimated completion time was approximately 1 hour.

## 5.2 TRADITIONAL HEURISTIC DECISION RULES

The decision-making process for brain penetration assays typically relies on heuristic rules, primarily using QSAR (Quantitative Structure-Activity Relationship) predictions and the unbound brain-to-plasma partition coefficient ( $k_{puu}$ ). These rules can be summarized as follows:

A compound is considered **promising** if:  $QSAR_{1\mu M\_PgP} < 2$ ,  $QSAR_{100nM\_BCRP} < 2$ , and  $0.5 \leq kpuu \leq 1$ . A compound is considered **non-promising** if either:  $QSAR_{1\mu M\_PgP}$  or  $QSAR_{100nM\_BCRP}$  exceeds 4, regardless of the  $kpuu$  value.

### 5.3 SELECTIVE CASE STUDY FOR COMPOUND SELECTION DECISION-MAKING

We tested the framework using three scenarios, designed to reflect different QSAR conditions. These case studies demonstrate the flexibility and robustness of our decision-making framework, showing its potential to improve the drug discovery process by identifying promising compounds that might be missed by traditional methods.:

**Baseline Confirmation:** This scenario tests compounds where both  $QSAR_{1\mu M\_PgP}$  and  $QSAR_{100nM\_BCRP}$  are below 2, and  $kpuu$  values fall within normal ranges. It serves to validate traditional decision-making processes.

**Heuristic Challenge:** In this scenario, compounds present borderline or conflicting QSAR data, with at least one QSAR value exceeding 4. This scenario tests the framework’s ability to interpret complex signals and identify viable candidates.

**Opportunity Discovery:** This scenario evaluates compounds with high  $QSAR_{1\mu M\_PgP}$  and  $QSAR_{100nM\_BCRP}$  values, but acceptable  $kpuu$ . It aims to discover overlooked compounds that could be promising despite failing traditional heuristics.

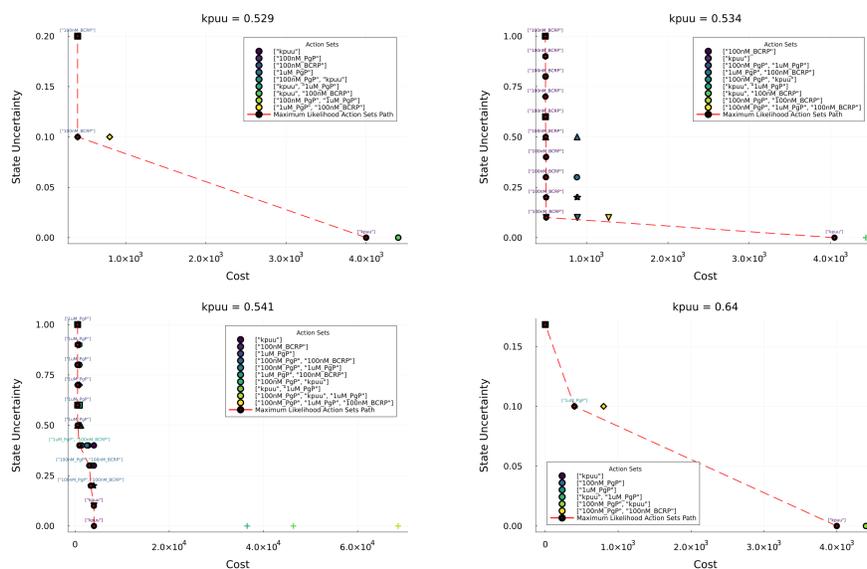


Figure 3: Monetary-prioritized IB-MDP results with MLASPs for four representative compounds, ordered by  $kpuu$  values to illustrate variations in QSAR metrics and corresponding recommended actions. For  $kpuu = 0.529$ ,  $QSAR_{1\mu M\_PgP} = 5.0$ ,  $QSAR_{100nM\_BCRP} = 9.6$ , and  $QSAR_{mrt} = 0.99$ . The IB-MDP recommends action is [100nM\_BCRP]. For  $kpuu = 0.534$ ,  $QSAR_{1\mu M\_PgP} = 0.903$ ,  $QSAR_{100nM\_BCRP} = 8.5$ , and  $QSAR_{mrt} = 2.64$ . The recommended action is [100nM\_BCRP]. For  $kpuu = 0.5407$ ,  $QSAR_{1\mu M\_PgP} = 1.68$ ,  $QSAR_{100nM\_BCRP} = 1.3$ , and  $QSAR_{mrt} = 1.82$ . The IB-MDP suggests actions are either [100nM\_PgP, 100nM\_BCRP] or [1uM\_PgP, 100nM\_BCRP]. For  $kpuu = 0.6400$ ,  $QSAR_{1\mu M\_PgP} = 21.4$ ,  $QSAR_{100nM\_BCRP} = 0.73$ , and  $QSAR_{mrt} = 1.2$ . Recommended actions include [1uM\_PgP], indicating a high probability of effectiveness under the given experimental conditions.

### 5.4 EXPERIMENTAL RESULTS : COST COMPARISON BETWEEN CONVENTIONAL AND IB-MDP DECISIONS

The results of the IB-MDP exploration for the representative cases in Table 1 are shown in Figure 3. In the baseline scenario, the IB-MDP recommends actions involving [1uM\_PgP, 100nM\_BCRP],

[100nM\_PgP, 100nM\_BCRP], or 1uM\_PgP, resulting in monetary costs ranging from \$400 to \$800, compared to the traditional cost of \$5200.

In the heuristic challenge scenario, the IB-MDP still proposes a single action along the MLASP with a \$400 cost, whereas traditional heuristic rules completely miss the opportunity to identify this promising compound. For the opportunity discovery scenario where all QSAR values are greater than 4, the IB-MDP successfully identifies a unique set of actions [100nM\_PgP, 1uM\_PgP] that significantly reduce state uncertainty. In contrast, the traditional rules fail to recognize this specific compound as a promising candidate.

Table 1: Comparison of Traditional Approach and IB-MDP Generated Costs for Selected Compounds

QSAR <sub>1uM_PgP</sub>	QSAR <sub>100nM_BCRP</sub>	QSAR <sub>mt</sub>	kpuu	100nM_PgP	1uM_PgP	100nM_BCRP	Traditional Cost	IB-MDP Cost
1.68	1.3	1.82	0.5407	1.06	0.79	1.32	\$5200	\$400 - \$800
0.903	8.5	2.64	0.5343	2.16	1.14	14.16	\$5200	\$400
21.4	0.73	1.2	0.6400	17.42	19.69	0.83	\$5200	\$400 - \$800
5.0	9.6	0.99	0.5289	15.92	12.86	8.23	\$5200	\$800

## 6 CONCLUSIONS

In this study, we present the Implicit Bayesian Markov Decision Process (IB-MDP), a framework designed to improve decision making under uncertainty in resource-constrained environments. By dynamically integrating historical data using a similarity-based metric, the IB-MDP refines beliefs about the current state in relation to target features. This refinement is achieved through implicit Bayesian updates with a sampling approach, ensuring the policy search maximizes information gain, minimizes costs, and achieves key experimental objectives within an optimal range.

A key advantage of the IB-MDP is its ability to significantly reduce state uncertainty without the need to perform the most expensive and resource-intensive assays, such as the target assay (kpuu). This not only enhances cost-efficiency but also accelerates decision-making by bypassing less critical, high-cost experiments.

Furthermore, IB-MDP benefits from an ensemble approach that aggregates policies across multiple runs, reducing variance, and improving robustness. By aligning maximum likelihood action sets with predefined probability bounds for target features, the methodology ensures consistent decision quality. Through its comprehensive and data-driven approach, the IB-MDP outperforms traditional heuristic methods, offering enhanced adaptability, precision, and resource optimization in experimental planning. It demonstrates significant potential in streamlining decision-making tasks in drug discovery and other fields requiring strategic resource management under uncertainty.

## 7 BROADER IMPACTS

The IB-MDP framework provides a versatile approach to adaptive decision making, offering benefits beyond preclinical assay scheduling. By integrating historical data, dynamic Bayesian updates, and ensemble methods, the framework enhances decision-making efficiency in various fields such as healthcare, logistics, and financial risk management. Its ability to handle uncertainty and resource constraints ensures robust, data-driven decisions, making it a valuable tool to improve outcomes in diverse industries.

## 8 LIMITATIONS

Increasing the number of runs  $N$  enhances the accuracy and robustness of the optimal action set estimation but also increases computational costs. While the ensemble method helps reduce variance and bias, its efficiency decreases as  $N$  grows, with diminishing returns typical in ensemble-based decision-making frameworks. The optimal value of  $N$  depends on the problem’s complexity and the available computational resources, as larger ensembles may be necessary to fully explore complex state-action spaces. In particularly intricate environments, more iterations may be required to ensure

540 reliable convergence, although the ensemble approach generally stabilizes after a sufficient number  
541 of runs.

542 Another important limitation lies in the framework’s reliance on historical data. Although leveraging  
543 historical data helps to integrate similarity-based metrics for decision-making, it may not sufficiently  
544 account for novel scenarios in real-world experiments. Future extensions of the IB-MDP framework  
545 could incorporate more flexible strategies, such as adaptive kernel-based methods or deep learning  
546 approaches, to extrapolate to states not represented in the existing dataset.

547 Moreover, while thresholds for decision-making tend to converge toward a maximum likelihood ac-  
548 tion path, further exploration of how state uncertainty reduction influences the likelihood of achiev-  
549 ing desired outcomes could offer additional insights. This might enable more effective ensemble  
550 adjustments, ultimately improving policy reliability and performance in dynamic environments.

#### 552 ACKNOWLEDGMENTS

553 We thank Karsten Menzel, Karin Otte, Kevin Bateman and Antong Chen for their helpful feedback  
554 and suggestions.

#### 557 REFERENCES

558  
559 Ismail Abbas, Joan Rovira, and Josep Casanovas. Clinical trial optimization: Monte carlo simulation  
560 markov model for planning clinical trials recruitment. *Contemporary clinical trials*, 28(3):220–  
561 231, 2007.

562 Oguzhan Alagoz, Heather Hsu, Andrew J Schaefer, and Mark S Roberts. Markov decision processes:  
563 a tool for sequential decision making under uncertainty. *Medical Decision Making*, 30(4):474–  
564 483, 2010.

565  
566 Aurélien Bellet, Amaury Habrard, and Marc Sebban. A survey on metric learning for feature vectors  
567 and structured data. *arXiv preprint arXiv:1306.6709*, 2013.

568 Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow  
569 network based generative models for non-iterative diverse candidate generation. *Advances in*  
570 *Neural Information Processing Systems*, 34:27381–27394, 2021.

571  
572 Casey C Bennett and Kris Hauser. Artificial intelligence framework for simulating clinical decision-  
573 making: A markov decision process approach. *Artificial intelligence in medicine*, 57(1):9–19,  
574 2013.

575  
576 Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp  
577 Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey  
578 of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in*  
579 *games*, 4(1):1–43, 2012.

580 Antonio Coronato, Muddasar Naeem, Giuseppe De Pietro, and Giovanni Paragliola. Reinforcement  
581 learning for intelligent healthcare applications: A survey. *Artificial intelligence in medicine*, 109:  
582 101964, 2020.

583  
584 Adrien Couëtoux, Jean-Baptiste Hooek, Nataliya Sokolovska, Olivier Teytaud, and Nicolas Bon-  
585 nard. Continuous upper confidence trees. In *Learning and Intelligent Optimization: 5th Interna-*  
586 *tional Conference, LION 5, Rome, Italy, January 17-21, 2011. Selected Papers 5*, pp. 433–445.  
587 Springer, 2011.

588 Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on multi-*  
589 *ple classifier systems*, pp. 1–15. Springer, 2000.

590  
591 Maryam Eghbali-Zarch, Reza Tavakkoli-Moghaddam, Fatemeh Esfahanian, Amir Azaron, and Mo-  
592 hammad Mehdi Sepehri. A markov decision process for modeling adverse drug reactions in  
593 medication treatment of type 2 diabetes. *Proceedings of the Institution of Mechanical Engineers,*  
*Part H: Journal of Engineering in Medicine*, 233(8):793–811, 2019.

- 594 Pablo Escandell-Montero, Milena Chermisi, Jose M Martinez-Martinez, Juan Gomez-Sanchis, Carlo  
595 Barbieri, Emilio Soria-Olivas, Flavio Mari, Joan Vila-Francés, Andrea Stopper, Emanuele Gatti,  
596 et al. Optimization of anemia treatment in hemodialysis patients via reinforcement learning.  
597 *Artificial intelligence in medicine*, 62(1):47–60, 2014.
- 598  
599 Mahdi M Fard, Sandor Szalma, Shashikant Vattikuti, and Gyan Bhanot. A bayesian markov decision  
600 process framework for optimal decision making in clinical trials. *IEEE Journal of Biomedical and*  
601 *Health Informatics*, 22(6):2061–2068, 2018.
- 602  
603 Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, and Aviv Tamar. Bayesian reinforcement  
604 learning: A survey. *Foundations and Trends® in Machine Learning*, 8(5-6):359–483, 2015.
- 605  
606 Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato,  
607 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel,  
608 Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven contin-  
609 uous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- 610  
611 Ryan-Rhys Griffiths and José Miguel Hernández-Lobato. Constrained bayesian optimization for  
612 automatic chemical design using variational autoencoders. *Chemical science*, 11(2):577–586,  
613 2020.
- 614  
615 Amnon Hoffman. Pharmacodynamic aspects of sustained release preparations. *Advanced drug*  
616 *delivery reviews*, 33(3):185–199, 1998.
- 617  
618 Amnon Hoffman, David Stepensky, Eran Lavy, Sara Eyal, Eytan Klausner, and Michael Friedman.  
619 Pharmacokinetic and pharmacodynamic aspects of gastroretentive dosage forms. *International*  
620 *journal of pharmaceutics*, 277(1-2):141–153, 2004.
- 621  
622 John P Hughes, Simon Rees, Sonya B Kalindjian, and Roger K Philpott. Principles of early drug  
623 discovery. *British journal of pharmacology*, 162(6):1239–1249, 2011.
- 624  
625 Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell, Krzysztof  
626 Czechowski, Dumitru Erhan, Chelsea Finn, Patryk Kozakowski, Sergey Levine, et al. Model-  
627 based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.
- 628  
629 Branislav Kveton and Georgios Theodorou. Kernel-based reinforcement learning on represen-  
630 tative states. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pp.  
631 977–983, 2012.
- 632  
633 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
634 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*,  
635 30, 2017.
- 636  
637 Mercedes F Martín, Carmen Sánchez, and Eva Gómez. Markov decision processes for modeling  
638 and optimization of drug discovery. *Journal of Chemical Information and Modeling*, 60(5):2494–  
639 2506, 2020.
- 640  
641 Dirk Ormoneit and Śaunak Sen. Kernel-based reinforcement learning. *Machine learning*, 49:161–  
642 178, 2002.
- 643  
644 Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via  
645 bootstrapped dqn. *Advances in neural information processing systems*, 29, 2016.
- 646  
647 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
Wiley & Sons, 1994.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
Wiley & Sons, 2014.
- Tom Rainforth, Adam Foster, Desi R Ivanova, and Freddie Bickford Smith. Modern bayesian ex-  
perimental design. *Statistical Science*, 39(1):100–114, 2024.

- 648 Petra Schneider, W Patrick Walters, Alleyn T Plowright, Norman Sieroka, Jennifer Listgarten,  
649 Robert A Goodnow Jr, Jasmin Fisher, Johanna M Jansen, José S Duca, Thomas S Rush, et al.  
650 Rethinking drug design in the artificial intelligence era. *Nature reviews drug discovery*, 19(5):  
651 353–364, 2020.
- 652 Marwin HS Segler, Mike Preuss, and Mark P Waller. Planning chemical syntheses with deep neural  
653 networks and symbolic ai. *Nature*, 555(7698):604–610, 2018.
- 654 Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the  
655 human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):  
656 148–175, 2015.
- 657 Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 658 Jike Wang, Chang-Yu Hsieh, Mingyang Wang, Xiaorui Wang, Zhenxing Wu, Dejun Jiang, Benben  
659 Liao, Xujun Zhang, Bo Yang, Qiaojun He, et al. Multi-constraint molecular generation based  
660 on conditional transformer, knowledge distillation and reinforcement learning. *Nature Machine  
661 Intelligence*, 3(10):914–922, 2021.
- 662 Marco A Wiering and Hado Van Hasselt. Ensemble algorithms in reinforcement learning. *IEEE  
663 Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 38(4):930–936, 2008.
- 664 Xin Xu, Dewen Hu, and Xicheng Lu. Kernel-based least squares policy iteration for reinforcement  
665 learning. *IEEE transactions on neural networks*, 18(4):973–992, 2007.
- 666 Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy  
667 network for goal-directed molecular graph generation. *Advances in neural information processing  
668 systems*, 31, 2018.
- 669 Zhenpeng Zhou, Steven Kearnes, Li Li, Richard N Zare, and Patrick Riley. Optimization of  
670 molecules via deep reinforcement learning. *Scientific reports*, 9(1):10752, 2019.
- 671 Zhi-Hua Zhou. *Ensemble methods: foundations and algorithms*. CRC press, 2012.
- 672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701