

Adaptive Information Gathering via Imitation Learning

Sanjiban Choudhury*, Ashish Kapoor†, Gireeja Ranade†, Sebastian Scherer* and Debadeepta Dey†

*The Robotics Institute

Carnegie Mellon University, Pittsburgh, PA 15232

Email: {sanjibac,basti}@andrew.cmu.edu

†Microsoft Research, Redmond USA

Email: {akapoor,giranade,dedey}@microsoft.com

Abstract—In the adaptive information gathering problem, a policy is required to select an informative sensing location using the history of measurements acquired thus far. While there is an extensive amount of prior work investigating effective practical approximations using variants of Shannon’s entropy, the efficacy of such policies heavily depends on the geometric distribution of objects in the world. On the other hand, the principled approach of employing online POMDP solvers is rendered impractical by the need to explicitly sample online from a posterior distribution of world maps.

We present a novel data-driven imitation learning framework to efficiently train information gathering policies. The policy imitates a *clairvoyant oracle* - an oracle that at train time has full knowledge about the world map and can compute maximally informative sensing locations. We analyze the learnt policy by showing that offline imitation of a clairvoyant oracle is implicitly equivalent to online oracle execution in conjunction with posterior sampling. This observation allows us to obtain powerful near-optimality guarantees for information gathering problems possessing an adaptive sub-modularity property. As demonstrated on a spectrum of 2D and 3D exploration problems, the trained policies enjoy the best of both worlds - they adapt to different world map distributions while being computationally inexpensive to evaluate.

I. INTRODUCTION

This paper¹ examines the following information gathering problem - given a hidden world map, sampled from a prior distribution, the goal is to successively visit sensing locations such that the amount of relevant information uncovered is maximized while not exceeding a specified fuel budget. This problem fundamentally recurs in mobile robot applications such as autonomous mapping of environments using ground and aerial robots [1, 10], monitoring of water bodies [11] and inspecting models for 3D reconstruction [14, 12].

The nature of “interesting” objects in an environment and their spatial distribution influence the optimal trajectory a robot might take to explore the environment. As a result, it is important that a robot learns about the type of environment it is exploring as it acquires more information and adapts its exploration trajectories accordingly. This adaptation must be done online, and we provide such algorithms in this paper.

Consider a robot equipped with a sensor (RGBD camera) that needs to generate a map of an unknown environment. It is given a prior distribution about the geometry of the world - such as a distribution over narrow corridors, or bridges or

power-lines. At every time step, the robot visits a sensing location and receives a sensor measurement (e.g. depth image) that has some amount of information utility (e.g. surface coverage of objects with point cloud). If the robot employs a non-adaptive lawnmower-coverage pattern, it can potentially waste time and fuel visiting ineffective locations. On the other hand, if it uses the history of measurements to infer the geometry of unexplored space, it can plan efficient information-rich trajectories. This incentivizes training policies to optimally gather information on the distribution of worlds the robot expects to encounter.

Even though it is natural to think of this problem setting as a POMDP, we frame this problem as a novel data-driven imitation learning problem [26]. We propose a framework that trains a policy on a dataset of worlds by imitating a *clairvoyant oracle*. During the training process, the oracle has full knowledge about the world map (hence clairvoyant) and visits sensing locations that would maximize information. The policy is then trained to imitate these movements as best as it can using partial knowledge from the current history of movements and measurements. As a result of our novel formulation, we are able to sidestep a number of challenging issues in POMDPs like explicitly computing posterior distribution over worlds and planning in belief space.

Our contributions are as follows:

- 1) We map the information gathering problem to a POMDP and present an approach to solve it using imitation learning of a clairvoyant oracle.
- 2) We present a framework to train such a policy on the non i.i.d distribution of states induced by the policy itself.
- 3) We analyze the learnt policy by showing that offline imitation of a clairvoyant oracle is equivalent to online oracle execution in conjunction with posterior sampling.
- 4) We present results on a variety of datasets to demonstrate the efficacy of the approach.

The remainder of this paper is organized as follows. Section II presents the formal problem and Section III maps it to the imitation learning framework. The algorithms and analysis is presented in Section IV. Section V presents experimental results with conclusions in Section VI.

II. BACKGROUND

A. Notation

Let \mathcal{V} be a set of nodes corresponding to all sensing locations. The robot starts at node v_s . Let $\xi = (v_1, v_2, \dots, v_p)$

¹This work was conducted by Sanjiban Choudhury as part of a summer internship at Microsoft Research, Redmond, USA.

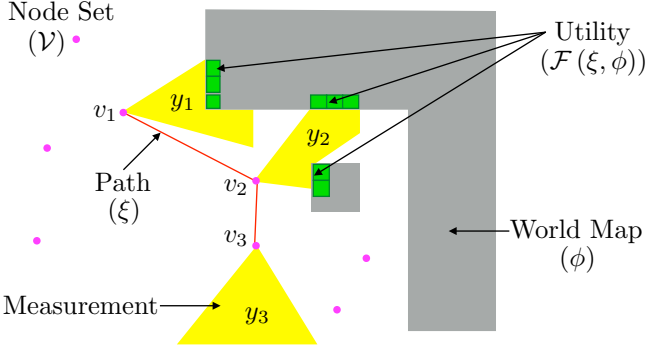


Fig. 1: The adaptive information gathering problem. Given a world map ϕ , the robot plans a path ξ which visits a node $v_i \in \mathcal{V}$ and receives measurement y_i , such that information gathered (utility) $\mathcal{F}(\xi, \phi)$ is maximized.

be a sequence of nodes (a path) such that $v_1 = v_s$. Let Ξ be the set of all such paths. Let $\phi \in \mathcal{M}$ be the world map. Let $y \in \mathcal{Y}$ be a measurement received by the robot. Let $\mathcal{H} : \mathcal{V} \times \mathcal{M} \rightarrow \mathcal{Y}$ be a measurement function. When the robot is at node v in a world map ϕ , the measurement y received by the robot is $y = \mathcal{H}(v, \phi)$. Let $\mathcal{F} : 2^{\mathcal{V}} \times \mathcal{M} \rightarrow \mathbb{R}_{\geq 0}$ be a utility function. For a path ξ and a world map ϕ , $\mathcal{F}(\xi, \phi)$ assigns a utility to executing the path on the world. Given a node $v \in \mathcal{V}$, a set of nodes $V \subseteq \mathcal{V}$ and world ϕ , the discrete derivative of the utility function \mathcal{F} is $\Delta_{\mathcal{F}}(v | V, \phi) = \mathcal{F}(V \cup \{v\}, \phi) - \mathcal{F}(V, \phi)$. Let $\mathcal{T} : \Xi \times \mathcal{M} \rightarrow \mathbb{R}_{\geq 0}$ be a travel cost function. For a path ξ and a world map ϕ , $\mathcal{T}(\xi, \phi)$ assigns a travel cost for executing the path on the world. Fig. 1 shows an illustration.

B. Problems with Known World Maps

We define four variants of the information gathering problem. For the first two variants, the world map ϕ is known and can be evaluated while computing a path ξ .

Problem 1 (KNOWN-UNC: Known World Map; Unconstrained Travel Cost). *Given a world map ϕ and a time horizon T , find a path ξ that maximizes utility*

$$\begin{aligned} \arg \max_{\xi \in \Xi} \quad & \mathcal{F}(\xi, \phi) \\ \text{s.t.} \quad & |\xi| \leq T + 1 \end{aligned} \quad (1)$$

Problem 2 (KNOWN-CON: Known World Map; Constrained Travel Cost). *Problem 1 with a travel cost budget B*

$$\begin{aligned} \arg \max_{\xi \in \Xi} \quad & \mathcal{F}(\xi, \phi) \\ \text{s.t.} \quad & \mathcal{T}(\xi, \phi) \leq B \\ & |\xi| \leq T + 1 \end{aligned} \quad (2)$$

Problem 1 is a set function maximization problem which in general can be NP-Hard (Krause and Golovin [20]). However, the utility function \mathcal{F} is a *monotone submodular* function. For such functions, it has been shown that greedy strategies achieve near-optimality (Krause et al. [22], Krause and Guestrin [21]).

Problem 2 introduces a routing constraint (due to \mathcal{T}) for which greedy approaches can perform arbitrarily poorly. Chekuri and Pal [2], Singh et al. [30] propose a quasi-polynomial time recursive greedy approach to solving this problem. Iyer and Bilmes [15] solve a related problem (submodular knapsack constraints) using an iterative greedy approach which is generalized by Zhang and Vorobeychik [35]. Yu et al. [34] propose a mixed integer approach to solve a related correlated orienteering problem. Hollinger and Sukhatme [11] propose a sampling based approach.

C. Problems with Hidden World Maps

We now consider the setting where the world map ϕ is hidden. Given a prior distribution $P(\phi)$, it can be inferred only via the measurements y_i received as the robot visits nodes v_i . Hence, instead of solving for a fixed path, we compute a policy that maps history of measurements received and nodes visited to decide which node to visit.

Problem 3 (HIDDEN-UNC: Hidden World Map; Unconstrained Travel Cost). *Given a distribution of world maps, $P(\phi)$, a time horizon T , find a policy that at time t , maps the history of nodes visited $\{v_i\}_{i=1}^{t-1}$ and measurements received $\{y_i\}_{i=1}^{t-1}$ to compute node v_t to visit at time t , such that the expected utility is maximized.*

Problem 4 (HIDDEN-CON: Hidden World Map; Constrained Travel Cost). *Problem 3 with a travel cost budget B*

Due to the hidden world map ϕ , it is not straight forward to apply the approaches discussed in Section II-B - methods have to reason about how $P(\phi | \{v_i\}_{i=1}^{t-1}, \{y_i\}_{i=1}^{t-1})$ will evolve. However, the utility function \mathcal{F} has an additional property of *adaptive submodularity* [7]. Hence, applying greedy strategies to Problem 3 has near-optimality guarantees (Golovin et al. [8], Javdani et al. [16, 17], Chen et al. [4, 5]).

Problem 4 does not enjoy the adaptive submodularity property. Hollinger et al. [13, 12] propose a heuristic based approach to select a subset of informative nodes and perform minimum cost tours. Singh et al. [31] replan every step using a non-adaptive information path planning algorithm. Inspired by adaptive TSP approaches by Gupta et al. [9], Lim et al. [24, 23] propose recursive coverage algorithms to learn policy trees. However such methods cannot scale well to large state and observation spaces. Heng et al. [10] make a modular approximation of the objective function. Isler et al. [14] survey a broad number of myopic information gain based heuristics that work well in practice but have no formal guarantees.

III. POMDPs AND IMITATION LEARNING

A. Mapping Problems to a POMDP

We now map Problems HIDDEN-UNC and HIDDEN-CON to a Partially Observable Markov Decision Process (POMDP). The POMDP is a tuple $(\mathcal{S}, \mathcal{M}, \mathcal{A}, \Omega, R, \mathcal{O}, Z, T)$ defined upto a fixed finite horizon T . It is defined over an augmented state space comprising of the ego-motion state space \mathcal{S} (which we will refer to as simply the state space) and the space of world maps \mathcal{M} . The first component, \mathcal{S} , is fully observable while

the second component, \mathcal{M} , is partially observable through observations received.

Let the state, $s_t \in \mathcal{S}$, be the set of nodes visited, $s_t = (v_1, v_2, \dots, v_t)$. Let the action, $a_t \in \mathcal{A}$ be the node visited $a_t = v_{t+1}$. Given a world map ϕ , at state s , the utility of a is $\mathcal{F}(s \cup a, \phi)$. For Problem HIDDEN-CON, let $\mathcal{A}_{\text{feas}}(s, \phi) \subset \mathcal{A}$ be the set of feasible actions defined as

$$\mathcal{A}_{\text{feas}}(s, \phi) = \{a \mid a \in \mathcal{A}, \mathcal{T}(s \cup a, \phi) \leq B\} \quad (3)$$

The state transition function, $\Omega(s, a, s') = P(s'|s, a)$, is the deterministic function $s' = s \cup a$. The one-step-reward function, $R(s, \phi, a) \in [0, 1]$, is defined as the normalized marginal gain of the utility function, $R(s, \phi, a) = \frac{\Delta_{\mathcal{F}}(a|s, \phi)}{\mathcal{F}(\mathcal{A}, \phi)}$. Let the observation, $o_t \in \mathcal{O}$ be the measurement $o_t = y_t$. The observation model, $Z(s, a, \phi, o) = P(o|s, a, \phi)$ is the deterministic function $o = \mathcal{H}(s \cup a, \phi)$.

We define the belief, ψ_t , to be the history of state, action, observation tuple received so far, i.e. $\{(s_i, a_i, o_i)\}_{i=1}^t$. The belief transition function $P(\psi'| \psi, a)$ can be computed from Ω and Z . Let $\tilde{\pi}(s, \psi) \in \tilde{\Pi}$ be a policy that maps state s and belief ψ to a feasible action $a \in \mathcal{A}_{\text{feas}}(s, \phi)$. The value of executing a policy $\tilde{\pi}$ for t time steps starting at state s and belief ψ is the expected cumulative reward obtained by $\tilde{\pi}$:

$$\tilde{V}_t^{\tilde{\pi}}(s, \psi) = \sum_{i=1}^t \mathbb{E}_{\substack{\psi_i \sim P(\psi' | \psi, \tilde{\pi}, i) \\ \phi \sim P(\phi | \psi_i) \\ s_i \sim P(s' | s, \tilde{\pi}, i)}} [R(s_i, \phi, \tilde{\pi}(s_i, \psi_i))] \quad (4)$$

where $P(\psi' | \psi, \tilde{\pi}, i)$ is the distribution of beliefs at time i starting from ψ and following policy $\tilde{\pi}$. Similarly $P(s' | s, \tilde{\pi}, i)$ is the distribution of states. $P(\phi | \psi_i)$ is the posterior distribution on worlds given the belief ψ_i .

The state-action value function $\tilde{Q}_t^{\tilde{\pi}}$ is defined as the expected sum of one-step-reward and value-to-go

$$\tilde{Q}_t^{\tilde{\pi}}(s, \psi, a) = \mathbb{E}_{\phi \sim P(\phi | \psi)} [R(s, \phi, a)] + \mathbb{E}_{\psi' \sim P(\psi' | \psi, a), s' \sim P(s' | s, a)} [\tilde{V}_{t-1}^{\tilde{\pi}}(s', \psi')] \quad (5)$$

The optimal POMDP policy is obtained by minimization of the expected state-action value function

$$\pi^* = \arg \max_{\tilde{\pi} \in \tilde{\Pi}} \mathbb{E}_{\substack{t \sim U(1:T) \\ s \sim P(s | \tilde{\pi}, t) \\ \psi \sim P(\psi | \tilde{\pi}, t)}} [\tilde{Q}_{T-t+1}^{\tilde{\pi}}(s, \psi, \tilde{\pi}(s, \psi))] \quad (6)$$

where $U(1 : T)$ is a uniform distribution over the discrete interval $\{1, \dots, T\}$, $P(s | \tilde{\pi}, t)$ and $P(\psi | \tilde{\pi}, t)$ are the posterior distribution over states and beliefs following policy $\tilde{\pi}$ for t steps. The value of a policy $\tilde{\pi} \in \tilde{\Pi}$ for T steps on a distribution of worlds $P(\phi)$, starting states $P(s)$ and starting belief $P(\psi)$ is the expected value function for the full horizon

$$J(\tilde{\pi}) = \mathbb{E}_{s \sim P(s), \psi \sim P(\psi)} [\tilde{V}_T^{\tilde{\pi}}(s, \psi)] \quad (7)$$

Just as Problems HIDDEN-UNC and HIDDEN-CON map to a POMDP, Problems KNOWN-UNC and KNOWN-CON map to the corresponding MDP. While we omit the details for brevity, the MDP defines a corresponding policy $\pi(s, \phi) \in \Pi$, value function $V_t^\pi(s, \phi)$, state-action value function $Q_t^\pi(s, \phi, a)$ and optimal policy π_{MDP} .

Online POMDP planning also has a large body of work (see Ross et al. [28]). Although there exists fast solvers such as POMCP (Silver and Veness [29]) and DESPOT (Somani et al. [32]), the space of world maps is too large for online planning. An alternative class of approaches is model-free policy improvement ([25]). While these methods make very little assumptions about the problem, they are local and require careful initialization.

B. Imitation Learning

An alternative to policy improvement approaches is to train policies to imitate reference policies (or oracle policies). This is a suitable choice for scenarios where the problem objective is to imitate a user-defined policy. This is also a useful approach in scenarios where there exist good oracle policies for the original problem, however these policies cannot be executed online (e.g due to computational complexity) hence requiring imitation via an offline training phase.

We now formally define imitation learning as applied to our setting. Given a policy $\tilde{\pi}$, we define the distribution of states $P(s | \tilde{\pi})$ and beliefs $P(\psi | \tilde{\pi})$ induced by it (termed as *roll-in*). Let $\mathcal{L}(s, \psi, \tilde{\pi})$ be a loss function that captures how well policy $\tilde{\pi}$ imitates an oracle. Our goal is to find a policy $\hat{\pi}$ which minimizes the expected loss as follows.

$$\hat{\pi} = \arg \min_{\tilde{\pi} \in \tilde{\Pi}} \mathbb{E}_{s \sim P(s | \tilde{\pi}), \psi \sim P(\psi | \tilde{\pi})} [\mathcal{L}(s, \psi, \tilde{\pi})] \quad (8)$$

This is a non-i.i.d supervised learning problem. Ross and Bagnell [26] propose FORWARDTRAINING to train a non-stationary policy (one policy $\hat{\pi}_t$ for each timestep), where each policy $\hat{\pi}_t$ can be trained on distributions induced by previous policies $(\hat{\pi}_1, \dots, \hat{\pi}_{t-1})$. While this has guarantees, it is impractical given a different policy is needed for each timestep. For training a single policy, Ross and Bagnell [26] show how such problems can be reduced to no-regret online learning using dataset aggregation (DAGGER). The loss function they consider \mathcal{L} is a mis-classification loss with respect to what the expert demonstrated. Ross and Bagnell [27] extend the approach to the reinforcement learning setting where \mathcal{L} is the reward-to-go of an oracle reference policy by aggregating *values* to imitate (AGGREVATE).

C. Solving POMDP via Imitation of a Clairvoyant Oracle

To examine the suitability of imitation learning in the POMDP framework, we compare the training rules (6) and (8). We see that a good candidate loss function $\mathcal{L}(s, \psi, \tilde{\pi})$ should incentivize maximization of $\tilde{Q}_{T-t+1}^{\tilde{\pi}}(s, \psi, \tilde{\pi}(s, \psi))$. A suitable approximation of the optimal value function $\tilde{Q}_{T-t+1}^{\pi^*}$ that can be computed at train time would suffice. In this work we define cumulative reward gathered by a *clairvoyant oracle*

as the value function to imitate².

Definition 1 (Clairvoyant Oracle). *Given a distribution of world map $P(\phi)$, a clairvoyant oracle $\pi_{\text{OR}}(s, \phi)$ is a policy that maps state s and world map ϕ to a feasible action $a \in \mathcal{A}_{\text{feas}}(s, \phi)$ such that it approximates the optimal MDP policy, $\pi_{\text{OR}} \approx \pi_{\text{MDP}} = \arg \max_{\pi \in \Pi} J(\pi)$.*

The term *clairvoyant* is used because the oracle has full access to the world map ϕ at train time. The oracle can be used to compute state-action value as follows

$$Q_t^{\pi_{\text{OR}}}(s, \phi, a) = R(s, \phi, a) + \mathbb{E}_{s' \sim P(s'|s, a)} [V_{t-1}^{\pi_{\text{OR}}}(s', \phi)] \quad (9)$$

Our approach is to imitate the oracle during training. This implies that we train a policy $\hat{\pi}$ by solving the following optimization problem

$$\hat{\pi} = \arg \max_{\tilde{\pi} \in \tilde{\Pi}} \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s|\tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi, \tilde{\pi}(s, \psi))] \quad (10)$$

While we will define training procedures to concretely realize (12) later in Section IV, we offer some intuition behind this approach. Since the oracle π_{OR} knows the world map ϕ , it has appropriate information to assign a value to an action a . The policy $\hat{\pi}$ attempts to imitate this action from the partial information content present in its belief ψ . Due to this realization error, the policy $\hat{\pi}$ visits a different state, updates its belief, gains more information, and queries the oracle for the best action. Hence while the learnt policy can make mistakes in the beginning of an episode, with time it gets better at imitating the oracle.

D. Analysis using a Hallucinating Oracle

The learnt policy imitates a clairvoyant oracle that has access to more information (world map ϕ compared to belief ψ). Hence, the realizability error of the policy is due to two terms - firstly the information mismatch and secondly the expressiveness of feature space. This realizability error can be hard to bound making it difficult to bound the performance of the learnt policy. This motivates us to introduce a hypothetical construct, a *hallucinating oracle*, to alleviate the information mismatch.

Definition 2 (Hallucinating Oracle). *Given a prior distribution of world map $P(\phi)$ and roll-outs by a policy $\tilde{\pi}$, a hallucinating oracle $\tilde{\pi}_{\text{OR}}$ computes the instantaneous posterior distribution over world maps and takes the action with the highest expected state-action value as computed by the clairvoyant oracle.*

$$\tilde{\pi}_{\text{OR}} = \arg \max_{a \in \mathcal{A}} \mathbb{E}_{\phi' \sim P(\phi'|\psi, \tilde{\pi}, t)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a)] \quad (11)$$

²Imitation of a clairvoyant oracle in information gathering has been explored by Choudhury et al. [6]. We subsume the presented algorithm, EXPLORE, in our framework (as Algorithm QVALAGG) and instead focus on the theoretical insight on what it means to imitate a clairvoyant oracle. Kahn et al. [18] has explored a similar idea in the context of reactive obstacle avoidance where a clairvoyant MPC is imitated. This too can be subsumed in our framework (as Algorithm REWARDAGG). Additionally, we provide analysis highlighting when such an approach is effective (in Section III-D)

We will now show that by imitating a clairvoyant oracle we effectively imitate the corresponding hallucinating oracle

Lemma 1. *The offline imitation of clairvoyant oracle (12) is equivalent to sampling online a world from the posterior distribution and executing a hallucinating oracle as shown³*

$$\hat{\pi} = \arg \max_{\tilde{\pi} \in \tilde{\Pi}} \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s|\tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\tilde{\pi}_{\text{OR}}}(s, \phi, \tilde{\pi}(s, \psi))]$$

Note that a hallucinating oracle uses the same information content as the learnt policy. Hence the realization error is purely due to the expressiveness of the feature space. However, we now have to analyze the performance of the hallucinating oracle which chooses the best action at a time step given its current belief. We will see in Section IV-C that this policy is near-optimal for Problem HIDDEN-UNC. For Problem HIDDEN-CON, while this does not have any such guarantees, there is evidence to show this is an effective policy as alluded to in Koval et al. [19].

IV. APPROACH

A. Overview

We introduced imitation learning and its applicability to POMDPs in Section III. We now present a set of algorithms to concretely realize the process. The overall idea is as follows - we are training a policy $\hat{\pi}(s, \psi)$ that maps features extracted from state s and belief ψ to an action a . The training objective is to imitate a clairvoyant oracle that has access to the corresponding world map ϕ . Fig. 2 shows an abstract overview. In order to define concrete algorithms, there are two degrees of freedom that need to be specified

- 1) *Clairvoyant oracle to imitate:* The choice of the clairvoyant oracle defines $Q_t^{\pi_{\text{OR}}}(s, \phi, a)$ in (9). Depending on whether we are solving Problem HIDDEN-UNC or HIDDEN-CON, we explore two different kinds of oracles
 - (a) *Clairvoyant one-step-reward:* For Problem HIDDEN-UNC, we use the one-step-reward $R(s, \phi, a)$ in place of $Q_t^{\pi_{\text{OR}}}(s, \phi, a)$. This corresponds to greedily imitating the immediate reward. We will see in Section IV-C that this has powerful near-optimality guarantees.
 - (b) *Clairvoyant reward-to-go:* For Problem HIDDEN-CON, we define an oracle policy π_{OR} that approximately solves the underlying MDP and use its cumulative reward-to-go to compute $Q_t^{\pi_{\text{OR}}}(s, \phi, a)$ (refer (9)).
- 2) *Stationary / Non-stationary policy:* As mentioned in Section III-B, imitation learning deals with non i.i.d distributions induced by the policy themselves. Depending on the choice of policy type (non-stationary / stationary), a corresponding training algorithm exists that offers guarantees
 - (a) *Non-stationary policy:* For the non-stationary case, we have a policy for each timestep $\hat{\pi}^1, \dots, \hat{\pi}^T$. While

³Refer to supplementary for all proofs

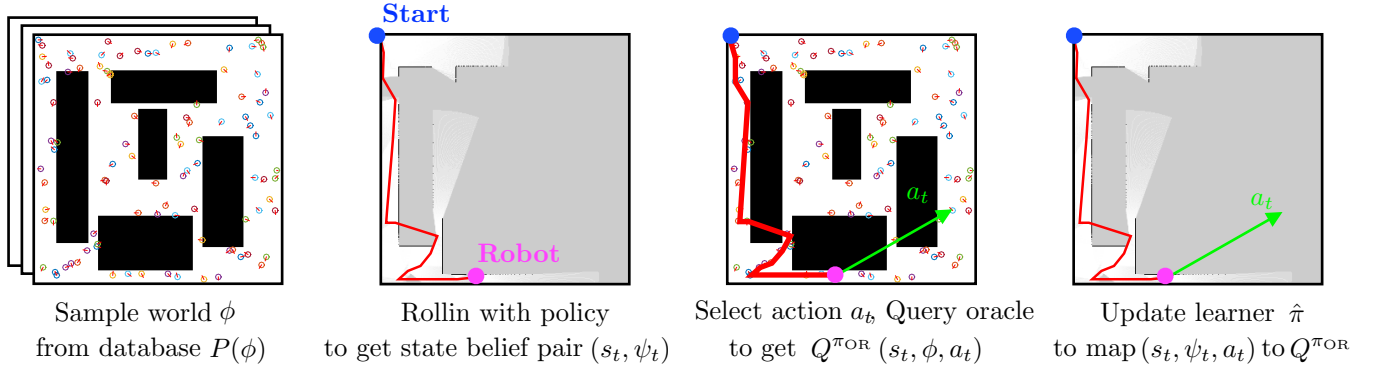


Fig. 2: An abstract overview of the imitation learning architecture where a learner $\hat{\pi}$ is trained to imitate a clairvoyant oracle π_{OR} . There are 4 key steps. Step 1: A world map ϕ is sampled from database representing $P(\phi)$. Step 2: A policy is used to roll-in on ϕ to a timestep t to get state s_t and belief ψ_t . Step 3: A random action a_t is chosen and a clairvoyant oracle π_{OR} is given full access to world map ϕ to compute the cumulative reward to go $Q^{\pi_{\text{OR}}}_{T-t+1}(s_t, \phi, a_t)$. Step 4: The learnt policy $\hat{\pi}$ is updated to map (s_t, ψ_t, a_t) to $Q^{\pi_{\text{OR}}}$.

this can be trained using the FORWARDTRAINING algorithm [26], there are several drawbacks. Firstly, it is impractical to have a different policy for each time-step as it scales with T . Secondly, the training has to proceed sequentially. Thirdly, each policy operates on data for only that time-step, thus preventing generalizations across timesteps. However, the procedure is presented for completeness.

- (b) *Stationary policy*: A single stationary policy $\hat{\pi}$ can be trained using the AGGREGATE algorithm [27] that reduces the problem to no-regret online learning setting. The training procedure is an interactive process where data is aggregated to refine the policy. The advantages are that the policy uses data across all timesteps, only one policy needs to be tracked and the training process can be stopped arbitrarily.

B. Algorithm

We now present concrete algorithms to realize the training procedure. Given the two axes of variation - problem and policy type - we have four possible algorithms

- 1) REWARDFT: Imitate one-step-reward using non-stationary policy by FORWARDTRAINING (Alg. 1)
- 2) QVALFT: Imitate reward-to-go using non-stationary policy by FORWARDTRAINING (Alg. 1)
- 3) REWARDAGG: Imitate one-step-reward using stationary policy by DAGGER (Alg. 2)
- 4) QVALAGG: Imitate reward-to-go using non-stationary policy by AGGREGATE (Alg. 2)

Table. I shows the algorithm mapping.

TABLE I: Mapping from Problem and Policy type to Algorithm

Policy	Problem	
	HIDDEN-UNC	HIDDEN-CON
Non-stationary policy	REWARDFT	QVALFT
Stationary policy	REWARDAGG	QVALAGG

Algorithm 1 Non-stationary policy (REWARDFT, QVALFT)

```

1: for  $t = 1$  to  $T$  do
2:   Initialize  $\mathcal{D} \leftarrow \emptyset$ .
3:   for  $j = 1$  to  $m$  do
4:     Sample world map  $\phi$  from dataset  $P(\phi)$ 
5:     Execute policy  $\hat{\pi}^1, \dots, \hat{\pi}^{t-1}$  to reach  $(s_t, \psi_t)$ .
6:     Execute any action  $a_t \in \mathcal{A}_{\text{feas}}(s_t, \phi)$ .
7:     Execute oracle  $\pi_{\text{OR}}$  from  $t+1$  to  $T$  on  $\phi$ 
8:     Collect value to go  $Q_i^{\pi_{\text{OR}}} = Q^{\pi_{\text{OR}}}_{T-t+1}(s_t, \phi, a_t)$ 
9:      $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, \psi_t, a_t, t, Q_i^{\pi_{\text{OR}}}\}$ 
10:   Train cost-sensitive classifier  $\hat{\pi}^t$  on  $\mathcal{D}$ 
11: Return Set of policies for each time step  $\hat{\pi}^1, \dots, \hat{\pi}^T$ .

```

Algorithm 2 Stationary policy (REWARDAGG, QVALAGG)

```

1: Initialize  $\mathcal{D} \leftarrow \emptyset$ ,  $\hat{\pi}_1$  to any policy in  $\tilde{\Pi}$ 
2: for  $i = 1$  to  $N$  do
3:   Initialize sub dataset  $\mathcal{D}_i \leftarrow \emptyset$ 
4:   Let roll-in policy be  $\pi_{\text{mix}} = \beta_i \pi_{\text{OR}} + (1 - \beta_i) \hat{\pi}_i$ 
5:   Collect  $m$  data points as follows:
6:   for  $j = 1$  to  $m$  do
7:     Sample world map  $\phi$  from dataset  $P(\phi)$ 
8:     Sample uniformly  $t \in \{1, 2, \dots, T\}$ 
9:     Assign initial state  $s_1 = v_s$ 
10:    Execute  $\pi_{\text{mix}}$  up to time  $t-1$  to reach  $(s_t, \psi_t)$ 
11:    Execute any action  $a_t \in \mathcal{A}_{\text{feas}}(s_t, \phi)$ 
12:    Execute oracle  $\pi_{\text{OR}}$  from  $t+1$  to  $T$  on  $\phi$ 
13:    Collect value-to-go  $Q_i^{\pi_{\text{OR}}} = Q^{\pi_{\text{OR}}}_{T-t+1}(s_t, \phi, a_t)$ 
14:     $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \{s_t, \psi_t, a_t, t, Q_i^{\pi_{\text{OR}}}\}$ 
15:   Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ 
16:   Train cost-sensitive classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ 
17:   (Alternately: use any online learner  $\hat{\pi}_{i+1}$  on  $\mathcal{D}_i$ )
18: Return best  $\hat{\pi}_i$  on validation

```

Alg. 1 describes the FORWARDTRAINING procedure to train the non-stationary policy. Previous policies $\hat{\pi}^1, \dots, \hat{\pi}^{t-1}$ are used to create a dataset (Lines 1–5). The oracle provides the value to imitate (Line 7) - if the one-step-reward is used the algorithm is referred to REWARDFT else if an oracle is used QVALFT.

Alg. 2 describes the AGGREGATE procedure to train the stationary policy. The algorithm iteratively trains a sequence of learnt policies $(\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_N)$ by aggregating data for an online cost-sensitive classification problem. At any given iteration, data is collected by roll-in with a mixture of learnt and oracle policy (Lines 1–10). The oracle provides the value to imitate (Lines 12–13) - if the one-step-reward is used the algorithm is referred to REWARDAGG else if an oracle is used QVALAGG. Data is appended to the original dataset and used to train an updated learner $\hat{\pi}_{i+1}$ (Lines 15–17).

C. Imitation of Clairvoyant One-Step-Reward

The strategy to imitate the clairvoyant one-step-reward is employed in Problem HIDDEN-UNC. This is motivated by the observation that the utility function satisfies a property of adaptive submodularity (as mentioned in Section II-C). For such problems, greedily selecting actions with highest expected reward has near-optimality guarantees. This implies the following Lemma

Lemma 2. *The performance of the hallucinating oracle $\tilde{\pi}_{\text{OR}}$ is near-optimal w.r.t the optimal policy π^* .*

$$J(\tilde{\pi}_{\text{OR}}) \geq \left(1 - \frac{1}{e}\right) J(\pi^*)$$

This property can then be used to obtain a near-optimality bound for the learnt policy

Theorem 1. *N iterations of REWARDAGG, collecting m regression examples per iteration guarantees that with probability at least $1 - \delta$*

$$\begin{aligned} J(\hat{\pi}) \geq & \left(1 - \frac{1}{e}\right) J(\pi^*) \\ & - 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} \\ & - O\left(\frac{T^2 \log T}{\alpha N}\right) \end{aligned}$$

where ε_{reg} is the empirical average online learning regret on the training regression examples collected over the iterations and $\varepsilon_{\text{class}}$ is the empirical regression regret of the best regressor in the policy class.

D. Imitation of Clairvoyant Reward-To-Go

Problem HIDDEN-CON does not possess the adaptive submodularity property. However there is empirical evidence to suggest that the hallucinating oracle performance $J(\tilde{\pi}_{\text{OR}})$ is sufficiently high [19]. Hence the learnt policy has a performance guarantee with respect to the hallucinating oracle

Theorem 2. *N iterations of QVALAGG, collecting m regression examples per iteration guarantees that with probability*

at least $1 - \delta$

$$\begin{aligned} J(\hat{\pi}) \geq & J(\tilde{\pi}_{\text{OR}}) \\ & - 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} \\ & - O\left(\frac{T^2 \log T}{\alpha N}\right) \end{aligned}$$

where ε_{reg} is the empirical average online learning regret on the training regression examples collected over the iterations and $\varepsilon_{\text{class}}$ is the empirical regression regret of the best regressor in the policy class.

V. EXPERIMENTAL RESULTS

A. Implementation Details

1) *Problem Details:* The utility function \mathcal{F} is selected to be a fractional coverage function (similar to [14]) which is defined as follows. The world map ϕ is represented as a voxel grid representing the surface of a 3D model. The sensor measurement $\mathcal{H}(v, \phi)$ at node v is obtained by ray-casting on this 3D model. A voxel of the model is said to be ‘covered’ by a measurement received at a node if a point lies in that voxel. The coverage of a path ξ is the fraction of covered voxels by the union of measurements received when visiting each node of the path. The travel cost function \mathcal{T} is chosen to be the euclidean distance. The values of total time step T and travel budget B vary with problem instances and are specified along with the results.

2) *Learning Details:* The tuple (s, a, ψ) is mapped to a vector of features $f = [f_{\text{IG}}^T \ f_{\text{mot}}^T]^T$. The feature vector f_{IG} is a vector of information gain metrics as described in [14]. f_{mot} encodes the relative rotation and translation required to visit a node. Random forest regression is used as a function approximator. The oracle used is the generalized cost benefit algorithm (GCB) [35].

3) *Baseline:* The baseline policies are a class of information gain heuristics discussed in [14] augmented with a motion penalization term when applied to Problem HIDDEN-CON. The heuristics are remarkably effective, however, their performance depends on the distribution of objections in a world map.

B. Adaptation to Different Distributions

We created a set of 2D exploration problems to gain a better understanding of the learnt policies and baseline heuristics. The problem was HIDDEN-UNC, the dataset comprises of 2D binary world maps, uniformly distributed nodes and a simulated laser. The problem details are $T = 30$ and $|\mathcal{A}| = 300$. The train size is 100, test size is 100. REWARDAGG is executed for 10 iterations.

The overall conclusion is that on changing the datasets the performance of the heuristics vary widely while the learnt policies outperform both heuristics. Interestingly, REWARDAGG outperforms REWARDFT- this is probably due to the generalization across time-steps. We now analyze each dataset.

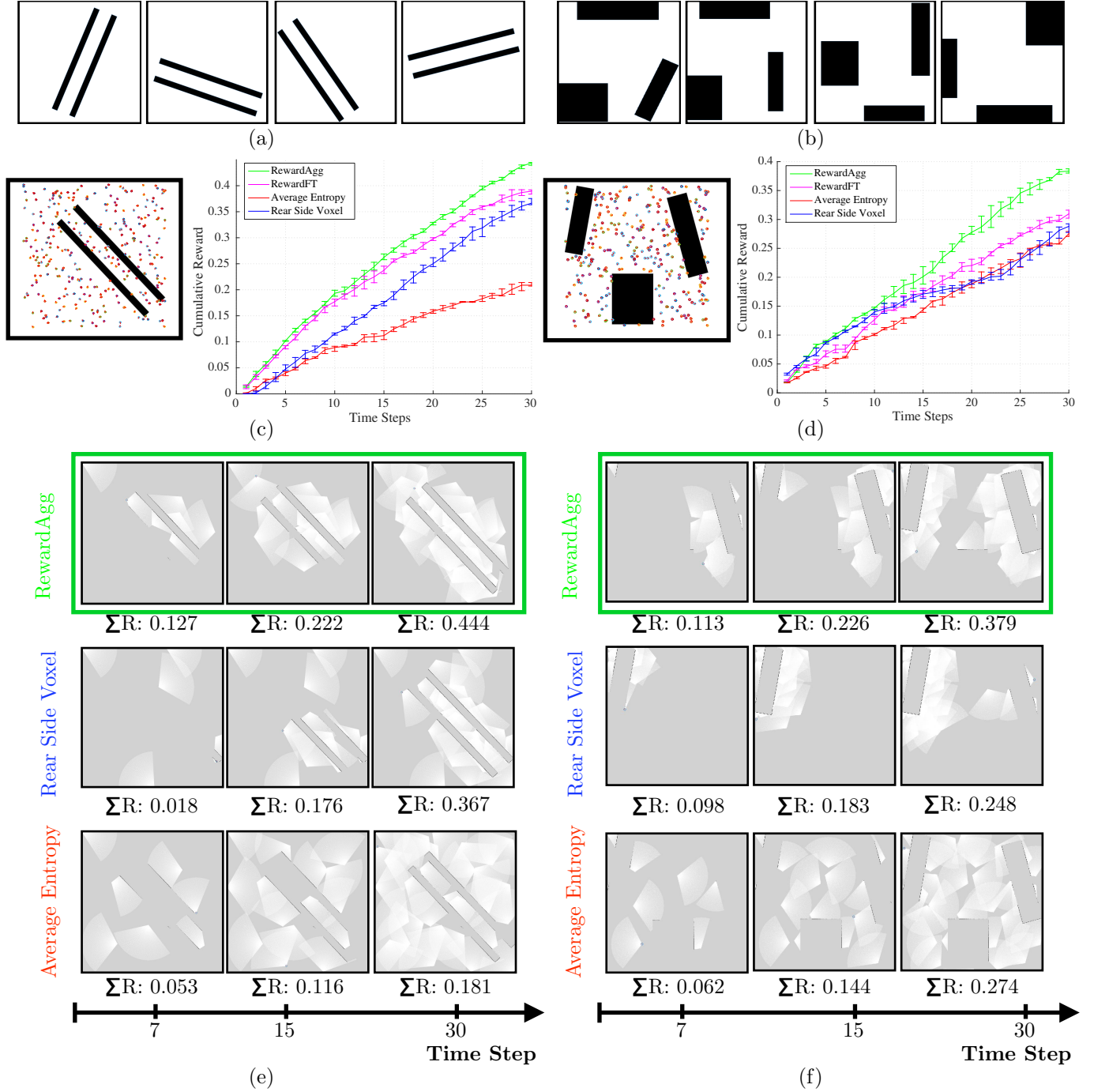


Fig. 3: Case study of Problem HIDDEN-UNC using REWARDAGG, REWARDFT and baseline heuristics. Two different datasets of 2D exploration are considered - (a) dataset 1 (parallel lines) and (b) dataset 2 (distributed blocks). Problem details are: $T = 30$, $|\mathcal{A}| = 300$, 100 train and 100 test maps. A sample test instance along with a plot of cumulative reward with time steps for different policies is shown in (c) and (d). The error bars show 95% confidence intervals. (e) and (f) show snapshots of the execution at time steps 7, 15 and 30.

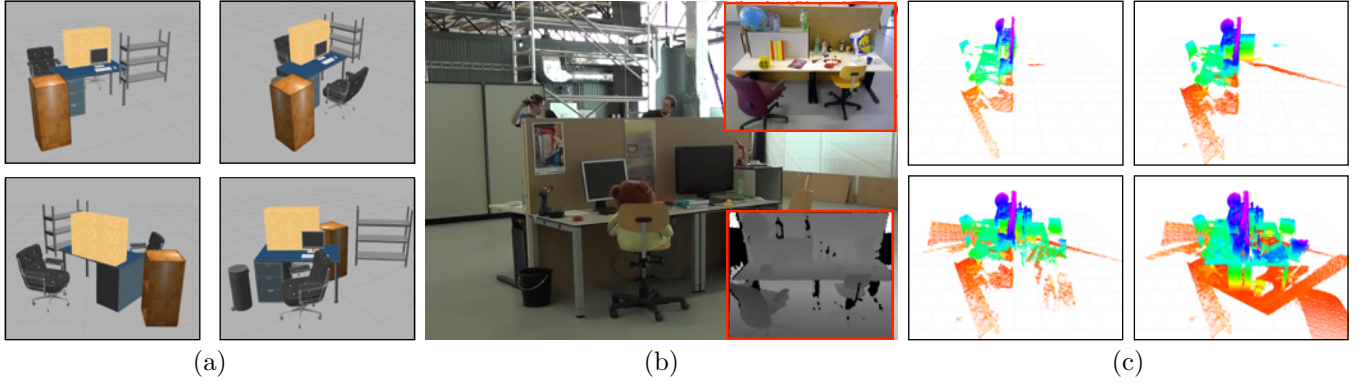


Fig. 4: Comparison of QVALAGG with baseline heuristics on a 3D exploration problem where training is done on simulated world maps and testing is done on a real dataset of an office workspace. The problem details are: $T = 10$, $B = 12$, $|\mathcal{A}| = 50$. (a) Samples from 100 simulated worlds resembling an office workspace created in Gazebo. (b) Real dataset collected by [33] using a RGBD camera. (c) Snapshots of execution of QVALAGG heuristic at time steps 1, 3, 5, 9.

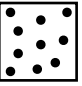
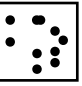
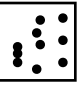


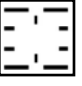









Dataset	Sample World Maps	Problem	RewardAgg / QvalAgg	Average Entropy	Occlusion Aware	Unobserved Voxels	Rear Side Voxels	Rear Side Entropy
Poisson Forest of Circular Discs (2D)	  	HIDDEN-UNC	(0.58, 0.61)	(0.59, 0.62)	(0.49, 0.53)	(0.39, 0.45)	(0.53, 0.55)	(0.42, 0.47)
		HIDDEN-CON	(0.54, 0.59)	(0.54, 0.59)	(0.42, 0.46)	(0.34, 0.41)	(0/37, 0.43)	(0.39, 0.44)
Tabular World of Rectilinear Blocks (2D)	  	HIDDEN-UNC	(0.43, 0.53)	(0.31, 0.35)	(0.20, 0.26)	(0.28, 0.35)	(0.35, 0.44)	(0.25, 0.31)
		HIDDEN-CON	(0.27, 0.33)	(0.26, 0.29)	(0.18, 0.23)	(0.21, 0.28)	(0.18, 0.24)	(0.21, 0.27)
Bookshelves and Tables (3D)	  	HIDDEN-UNC	(0.14, 0.31)	(0.01, 0.04)	(0.01, 0.04)	(0.01, 0.04)	(0.01, 0.22)	(0.01, 0.19)
		HIDDEN-CON	(0.05, 0.24)	(0.01, 0.04)	(0.01, 0.04)	(0.01, 0.04)	(0.01, 0.22)	(0.01, 0.19)
Cluttered Construction Site (3D)	  	HIDDEN-UNC	(0.14, 0.20)	(0.01, 0.12)	(0.01, 0.09)	(0.01, 0.09)	(0.01, 0.11)	(0.01, 0.10)
		HIDDEN-CON	(0.08, 0.12)	(0.01, 0.12)	(0.01, 0.09)	(0.01, 0.09)	(0.01, 0.11)	(0.01, 0.10)
Office Desk and Chairs (3D)	  	HIDDEN-UNC	(0.69, 0.80)	(0.46, 0.59)	(0.51, 0.63)	(0.51, 0.63)	(0.59, 0.67)	(0.61, 0.72)
		HIDDEN-CON	(0.55, 0.72)	(0.46, 0.59)	(0.48, 0.63)	(0.48, 0.63)	(0.43, 0.52)	(0.41, 0.53)

Fig. 5: Results for Problems HIDDEN-UNC and HIDDEN-CON on a spectrum of 2D and 3D exploration problems. The train size is 100 and test size is 10. Numbers are the confidence bounds (for 95% CI) of cumulative reward at the final time step. Algorithm with the highest median performance is emphasized in bold.

1) *Dataset 1: Parallel Lines*: Fig. 3a shows a dataset created by applying random affine transformations to a pair of parallel lines. This dataset is representative of information being concentrated in a particular fashion. Fig. 3c shows a comparison of REWARDAGG, REWARDFT with baseline heuristics. While Rear Side Voxel outperforms Average Entropy, REWARDAGG outperforms both. Fig. 3e shows progress of each. Average Entropy explores the whole world without focusing, Rear Side Voxel exploits early while QVALAGG trades off exploration and exploitation.

2) *Dataset 2: Distributed Blocks*: Fig. 3b shows a dataset created by randomly distributing rectangular blocks around the periphery of the map. This dataset is representative of information being distributed around. Fig. 3c shows that Rear Side Voxel saturates early, Average Entropy eventually overtaking it while REWARDAGG outperforms all. Fig. 3e shows that Rear Side Voxel gets stuck exploiting an island of information. Average Entropy takes broader sweeps of the area thus gaining

more information about the world. QVALAGG shows a non-trivial behavior exploiting one island before moving to another.

C. Train on Synthetic, Test on Real

To show the practical usage of our pipeline, we show a scenario where a policy is trained on synthetic data and tested on a real dataset. Fig. 4a shows some sample worlds created in Gazebo to represent an office desk environment on which QVALAGG is trained. Fig. 4b shows a dataset of an office desk collected by TUM Computer Vision Group [33]. The dataset is parsed to create a pair of pose and registered point cloud which can then be used to evaluate different algorithms. Fig. 4c shows QVALAGG learns a desk exploring policy by circumnavigating around the desk. This shows the powerful generalization capabilities of the approach.

D. Spectrum of 2D / 3D exploration problems

We evaluate the framework on a spectrum of 2D / 3D exploration problems on synthetic worlds as shown in

Fig. 5. For Problem HIDDEN-UNC, REWARDAGG is employed along with baseline heuristics. For Problem HIDDEN-CON, QVALAGG is employed with baseline heuristic augmented with motion penalization. The train size is 100 and test size is 10. We see that the learnt policies outperform all heuristics on most datasets by exploiting the distribution of objects in the world particularly in Problem HIDDEN-UNC. This is indicative of the power of the hallucinating oracle. However, the Poisson forest datasets stand out - where given the ergodic distribution, the Average Entropy heuristic performs best.

VI. CONCLUSION

We present a novel framework for learning information gathering policies via imitation learning of clairvoyant oracles. We presented analysis that establishes an equivalence to online imitation of hallucinating oracles thereby explaining the success of such policies. The framework is validated on a spectrum of 2D and 3D exploration problems.

There are several key directions for future research. Firstly, analysis from Chen et al. [3], where the authors show that under mild assumptions POMDPs can be reduced to sequence of MDPs, is promising for obtaining better regret guarantees. [23] also offer alternatives to adaptive submodularity that could improve guarantees for HIDDEN-CON. Secondly, instead of learning a policy, learning a surrogate utility function that is solved online with a routing TSP might lead to improved performance. Finally, we are looking to apply our framework to POMDP problems where hindsight optimization have shown success (e.g learning to grasp Koval et al. [19]).

APPENDIX

We begin with defining a clairvoyant oracle

Definition 3 (Clairvoyant Oracle). *Given a distribution of world map $P(\phi)$, a clairvoyant oracle $\pi_{\text{OR}}(s, \phi)$ is a policy that maps state s and world map ϕ to a feasible action $a \in \mathcal{A}_{\text{feas}}(s, \phi)$ such that it approximates the optimal MDP policy, $\pi_{\text{OR}} \approx \pi_{\text{MDP}} = \arg \max_{\pi \in \Pi} J(\pi)$.*

Our approach is to imitate the oracle during training. This implies that we train a policy $\hat{\pi}$ by solving the following optimization problem

$$\hat{\pi} = \arg \max_{\tilde{\pi} \in \tilde{\Pi}} \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s|\tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi, \tilde{\pi}(s, \psi))] \quad (12)$$

The learnt policy imitates a clairvoyant oracle that has access to more information (world map ϕ compared to belief ψ). Hence, the realizability error of the policy is due to two terms - firstly the information mismatch and secondly the expressiveness of feature space. This realizability error can be hard to bound making it difficult to bound the performance of the learnt policy. This motivates us to introduce a hypothetical construct, a *hallucinating oracle*, to alleviate the information mismatch.

Definition 4 (Hallucinating Oracle). *Given a prior distribution of world map $P(\phi)$ and rollouts by a policy $\tilde{\pi}$, a hallucinating oracle $\tilde{\pi}_{\text{OR}}$ computes the instantaneous posterior distribution over world maps and takes the action with the highest expected value.*

$$\tilde{\pi}_{\text{OR}} = \arg \max_{a \in \mathcal{A}} \mathbb{E}_{\phi' \sim P(\phi'|\psi, \tilde{\pi}, t)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a)] \quad (13)$$

We will now show that by imitating a clairvoyant oracle we effectively imitate the corresponding hallucinating oracle

Lemma 3. *The **offline** imitation of **clairvoyant** oracle (12) is equivalent to sampling **online** a world from the posterior distribution and executing a **hallucinating** oracle as shown*

$$\hat{\pi} = \arg \max_{\tilde{\pi} \in \tilde{\Pi}} \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s|\tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\tilde{\pi}_{\text{OR}}}(s, \phi, \tilde{\pi}(s, \psi))]$$

by using the fact that $\mathbb{E}_{\phi \sim P(\phi)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi, a)] = \mathbb{E}_{\substack{\phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\tilde{\pi}_{\text{OR}}}(s, \phi, a)]$

Proof: Given a belief ψ and rollout policy $\tilde{\pi}$, the action value function for the hallucinating oracle $Q_{T-t+1}^{\tilde{\pi}_{\text{OR}}}(s, \phi, a)$ is defined as

$$Q_{T-t+1}^{\tilde{\pi}_{\text{OR}}}(s, \phi, a) = \mathbb{E}_{\phi' \sim P(\phi'|\psi, \tilde{\pi}, t)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a)] \quad (14)$$

We will now show $\mathbb{E}_{\phi \sim P(\phi)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi, a)] = \mathbb{E}_{\substack{\phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\tilde{\pi}_{\text{OR}}}(s, \phi, a)]$.

Taking expectation $\mathbb{E}_{\substack{\phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [\cdot]$ on both sides of (14) we have

$$\mathbb{E}_{\substack{\phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi, a)] = \mathbb{E}_{\substack{\phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [\mathbb{E}_{\phi' \sim P(\phi'|\psi, \tilde{\pi}, t)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a)]] \quad (15)$$

Evaluating the right hand side of (15)

$$\begin{aligned} & \mathbb{E}_{\substack{\phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [\mathbb{E}_{\phi' \sim P(\phi'|\psi, \tilde{\pi}, t)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a)]] \\ &= \sum_{\phi} \sum_{\psi} \sum_{\phi'} P(\phi) P(\psi | \phi, \tilde{\pi}, t) P(\phi' | \psi, \tilde{\pi}, t) Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a) \\ &= \sum_{\phi} \sum_{\psi} \sum_{\phi'} P(\phi) \frac{P(\phi | \psi, \tilde{\pi}, t) P(\psi | \tilde{\pi}, t)}{P(\phi | \tilde{\pi}, t)} \frac{P(\psi | \phi', \tilde{\pi}, t) P(\phi' | \tilde{\pi}, t)}{P(\psi | \tilde{\pi}, t)} Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a) && \text{(Bayes rule)} \\ &= \sum_{\phi} \sum_{\psi} \sum_{\phi'} P(\phi | \psi, \tilde{\pi}, t) P(\psi | \phi', \tilde{\pi}, t) P(\phi') Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a) && \text{(Cond. Indep.)} \\ &= \sum_{\phi} P(\phi') Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a) \sum_{\psi} \sum_{\phi'} P(\phi | \psi, \tilde{\pi}, t) P(\psi | \phi', \tilde{\pi}, t) P(\phi') \\ &= \sum_{\phi} P(\phi') Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi', a) && \text{(Marginalizing } \phi, \psi) \\ &= \mathbb{E}_{\phi \sim P(\phi)} [Q_{T-t+1}^{\pi_{\text{OR}}}(s, \phi, a)] \end{aligned}$$

This implies that

$$\mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s|\tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, \tilde{\pi}(s, \psi))] = \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s|\tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi|\phi, \tilde{\pi}, t)}} [Q_{T-t+1}^{\pi\text{OR}}(s, \phi, \tilde{\pi}(s, \psi))]$$

■

Now that we have established that we are imitating a reference policy that is not clairvoyant (hallucinating oracle), we can invoke the analysis used by Ross et al. [27] for the algorithm AGGREGATE. A key difference is that AGGREGATE defines the problem as minimizing cost to go, while we define it as maximizing reward to go. Instead of re-deriving the analysis, we will instead define a term by term mapping between our framework and AGGREGATE and then use the proven bounds.

TABLE II: Mapping between QVALAGG and AGGREGATE

Term	QVALAGG	AGGREGATE
Objective	Reward: $R(s, \phi, a)$	Cost: $C(s, \phi, a)$ $C(s, \phi, a) = 1 - R(s, \phi, a)$
State Value Function	$V_T^{\tilde{\pi}}(s, \phi) = \sum_{i=1}^T \mathbb{E}_{s_i \sim P(s' s, \tilde{\pi}, i)} [R(s_i, \phi, \tilde{\pi}(s_i, \phi))]$	$\bar{V}_T^{\tilde{\pi}}(s, \phi) = \sum_{i=1}^T \mathbb{E}_{s_i \sim P(s' s, \tilde{\pi}, i)} [C(s_i, \phi, \tilde{\pi}(s_i, \phi))]$ $\bar{V}_T^{\tilde{\pi}}(s, \phi) = T - V_T^{\tilde{\pi}}(s, \phi)$
Action Value Function	$Q_T^{\tilde{\pi}}(s, \phi, a)$	$\bar{Q}_T^{\tilde{\pi}}(s, \phi, a) = T - Q_T^{\tilde{\pi}}(s, \phi, a)$
Policy Value	$J(\tilde{\pi}) = \mathbb{E}_{s \sim P(s), \phi \sim P(\phi)} [V_T^{\tilde{\pi}}(s, \phi)]$	$\bar{J}(\tilde{\pi}) = \mathbb{E}_{s \sim P(s), \phi \sim P(\phi)} [\bar{V}_T^{\tilde{\pi}}(s, \phi)]$ $\bar{J}(\tilde{\pi}) = T - J(\tilde{\pi})$
Loss of Policy	$l_i(\tilde{\pi}) = \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s \tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi \phi, \tilde{\pi}, t)}} [\max_{a \in \mathcal{A}} Q_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, a) - Q_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, \tilde{\pi}(s, \psi))]$	$\bar{l}_i(\tilde{\pi}) = \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s \tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi \phi, \tilde{\pi}, t)}} [\bar{Q}_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, \tilde{\pi}) - \min_{a \in \mathcal{A}} \bar{Q}_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, a)]$ $\bar{l}_i = l_i$
Classification Error of Best Policy	$\varepsilon_{\text{class}} = \min_{\tilde{\pi} \in \bar{\Pi}} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s \tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi \phi, \tilde{\pi}, t)}} [\max_{a \in \mathcal{A}} Q_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, a) - Q_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, \tilde{\pi})]$	$\bar{\varepsilon}_{\text{class}} = \min_{\tilde{\pi} \in \bar{\Pi}} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\substack{t \sim U(1:T), \\ s \sim P(s \tilde{\pi}, t), \\ \phi \sim P(\phi), \\ \psi \sim P(\psi \phi, \tilde{\pi}, t)}} [\bar{Q}_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, \tilde{\pi}) - \min_{a \in \mathcal{A}} \bar{Q}_{T-t+1}^{\tilde{\pi}\text{OR}}(s, \phi, a)]$ $\bar{\varepsilon}_{\text{class}} = \varepsilon_{\text{class}}$
Regret	$\varepsilon_{\text{reg}} = \frac{1}{N} \left[\sum_{i=1}^N l_i(\hat{\pi}_i) - \min_{\pi \in \bar{\Pi}} \sum_{i=1}^N l_i(\pi) \right]$	$\bar{\varepsilon}_{\text{reg}} = \frac{1}{N} \left[\sum_{i=1}^N \bar{l}_i(\hat{\pi}_i) - \min_{\pi \in \bar{\Pi}} \sum_{i=1}^N \bar{l}_i(\pi) \right]$ $\bar{\varepsilon}_{\text{reg}} = \varepsilon_{\text{reg}}$

The AGGREGATE proofs also defines a bound \bar{Q}_{\max}^* which is the performance bound of reference policy. This is mapped as follows

$$\begin{aligned} \bar{Q}_{\max}^* &= \max_{(s, a, t)} t - Q_t^{\tilde{\pi}\text{OR}}(s, \phi, a) \\ &\leq T \end{aligned} \tag{16}$$

Lemma 4. N iterations of QVALAGG for the infinite sample case

$$J(\hat{\pi}) \geq J(\tilde{\pi}_{\text{OR}}) - T[\varepsilon_{\text{class}} + \varepsilon_{\text{reg}}] - O\left(\frac{T^2 \log T}{\alpha N}\right)$$

Proof: Theorem 2.1 in [27] states

$$\bar{J}(\hat{\pi}) \leq \bar{J}(\tilde{\pi}_{\text{OR}}) + T[\bar{\varepsilon}_{\text{class}} + \bar{\varepsilon}_{\text{reg}}] + O\left(\frac{\bar{Q}_{\max}^* T \log T}{\alpha N}\right)$$

Substituting terms from Table. II and (16) we get

$$\begin{aligned}
\bar{J}(\hat{\pi}) &\leq \bar{J}(\tilde{\pi}_{\text{OR}}) + T[\bar{\varepsilon}_{\text{class}} + \bar{\varepsilon}_{\text{reg}}] + O\left(\frac{\bar{Q}_{\max}^* T \log T}{\alpha N}\right) \\
T - J(\hat{\pi}) &\leq T - J(\tilde{\pi}_{\text{OR}}) + T[\bar{\varepsilon}_{\text{class}} + \bar{\varepsilon}_{\text{reg}}] + O\left(\frac{T^2 \log T}{\alpha N}\right) \\
J(\hat{\pi}) &\geq J(\tilde{\pi}_{\text{OR}}) - T[\varepsilon_{\text{class}} + \varepsilon_{\text{reg}}] - O\left(\frac{T^2 \log T}{\alpha N}\right)
\end{aligned}$$

Lemma 5. N iterations of QVALAGG, collecting m regression examples per iteration guarantees that with probability at least $1 - \delta$

$$J(\hat{\pi}) \geq J(\tilde{\pi}_{\text{OR}}) - 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} - O\left(\frac{T^2 \log T}{\alpha N}\right)$$

where $\varepsilon_{\text{reg}}^{\hat{}}$ is the empirical average online learning regret on the training regression examples collected over the iterations and $\varepsilon_{\text{class}}^{\hat{}}$ is the empirical regression regret of the best regressor in the policy class.

Proof: Theorem 2.2 in [27] states

$$\bar{J}(\hat{\pi}) \leq \bar{J}(\tilde{\pi}_{\text{OR}}) + 2\sqrt{|\mathcal{A}|T} \sqrt{\bar{\varepsilon}_{\text{class}} + \bar{\varepsilon}_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} + O\left(\frac{\bar{Q}_{\max}^* T \log T}{\alpha N}\right)$$

Substituting terms from Table. II and (16) we get

$$\begin{aligned}
\bar{J}(\hat{\pi}) &\leq \bar{J}(\tilde{\pi}_{\text{OR}}) + 2\sqrt{|\mathcal{A}|T} \sqrt{\bar{\varepsilon}_{\text{class}} + \bar{\varepsilon}_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} + O\left(\frac{\bar{Q}_{\max}^* T \log T}{\alpha N}\right) \\
T - J(\hat{\pi}) &\leq T - J(\tilde{\pi}_{\text{OR}}) + 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} + O\left(\frac{T^2 \log T}{\alpha N}\right) \\
J(\hat{\pi}) &\geq J(\tilde{\pi}_{\text{OR}}) - 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} - O\left(\frac{T^2 \log T}{\alpha N}\right)
\end{aligned}$$

Lemma 6. The performance of the hallucinating oracle $\tilde{\pi}_{\text{OR}}$ is near-optimal w.r.t the optimal policy π^* .

$$J(\tilde{\pi}_{\text{OR}}) \geq \left(1 - \frac{1}{e}\right) J(\pi^*)$$

Proof:

Problem HIDDEN-UNC is maximization of an adaptive submodular function subject to cardinality constraints. An apdative greedy approach to solving this problem has a near-optimality guarantee as proved by Golovin and Krause [7]. We will establish an equivalence between the hallucinating one-step-reward oracle and the greedy algorithm.

The greedy algorithm selects a node to visit that has the highest expected marginal gain under the conditional distribution of world maps given the observations received so far. If the history of vertices visited and measurements received are $\{v_j\}_{j=1}^{i-1}$ and $\{y_j\}_{j=1}^{i-1}$, the greedy algorithm selects element v_i with the highest expected marginal gain

$$v_i = \arg \max_{v \in \mathcal{V}} \mathbb{E}_{\phi \sim P(\phi | \{v_j\}_{j=1}^{i-1}, \{y_j\}_{j=1}^{i-1})} [\Delta_{\mathcal{F}}(v | \{v_j\}_{j=1}^{i-1}, \phi)] \quad (17)$$

Golovin and Krause [7] prove the following: Let \mathcal{F} be a utility function satisfying the properties of a non-negative adaptive monotonone and adaptive submodular set function with respect to distribution $P(\phi)$. Let σ^* be the optimal solution to Problem HIDDEN-UNC. Let σ_{greedy} be the solution obtained by applying the greedy algorithm (17). Let $\mathbb{E}_{\phi \sim P(\phi)} [\mathcal{F}(\sigma, \phi)]$ be the expected utility of a policy when evaluated on a distribution of world maps $P(\phi)$. Then the following guarantee holds

$$\mathbb{E}_{\phi \sim P(\phi)} [\mathcal{F}(\sigma_{\text{greedy}}, \phi)] \geq (1 - e^{-1}) \mathbb{E}_{\phi \sim P(\phi)} [\mathcal{F}(\sigma^*, \phi)]$$

We now show the equivalence of hallucinating one-step-reward oracle to greedy algorithm. The hallucinating one-step-reward oracle is

$$\tilde{\pi}_{\text{OR}} = \arg \max_{a \in \mathcal{A}} \mathbb{E}_{\phi' \sim P(\phi' | \psi, \tilde{\pi}, t)} [R(s, \phi', a)] \quad (18)$$

This can be expanded as

$$\begin{aligned}
\tilde{\pi}_{\text{OR}} &= \arg \max_{a \in \mathcal{A}} \mathbb{E}_{\phi' \sim P(\phi' | \psi, \tilde{\pi}, t)} [R(s, \phi', a)] \\
&= \arg \max_{a \in \mathcal{A}} \mathbb{E}_{\phi' \sim P(\phi' | \psi, \tilde{\pi}, t)} \left[\frac{\Delta_{\mathcal{F}}(a | s, \phi')}{\mathcal{F}(\mathcal{A}, \phi')} \right] \\
&= \arg \max_{a \in \mathcal{A}} \mathbb{E}_{\phi' \sim P(\phi' | \psi, \tilde{\pi}, t)} [\Delta_{\mathcal{F}}(a | s, \phi')]
\end{aligned}$$

Since $s = \{v_j\}_{j=1}^t$, $\mathcal{A} = \mathcal{V}$ and $\psi_t = \{(v_j, y_j)\}_{j=1}^t$

$$\tilde{\pi}_{\text{OR}} = \arg \max_{v \in \mathcal{V}} \mathbb{E}_{\phi' \sim P(\phi' | \{v_j\}_{j=1}^t, \{y_j\}_{j=1}^t)} [\Delta_{\mathcal{F}}(v | \{v_j\}_{j=1}^t, \phi')]$$

Hence the $\tilde{\pi}_{\text{OR}}$ is equivalent to the greedy algorithm and we get the following guarantee

$$J(\tilde{\pi}_{\text{OR}}) \geq \left(1 - \frac{1}{e}\right) J(\pi^*)$$

Now proving the regret guarantee is a straightforward conjunction of lemmas ■

Theorem 3. *N iterations of REWARDAGG, collecting m regression examples per iteration guarantees that with probability at least $1 - \delta$*

$$\begin{aligned}
J(\hat{\pi}) &\geq \left(1 - \frac{1}{e}\right) J(\pi^*) \\
&\quad - 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} \\
&\quad - O\left(\frac{T^2 \log T}{\alpha N}\right)
\end{aligned}$$

where ε_{reg} is the empirical average online learning regret on the training regression examples collected over the iterations and $\varepsilon_{\text{class}}$ is the empirical regression regret of the best regressor in the policy class.

Proof: In Lemma 5 use $J(\tilde{\pi}_{\text{OR}}) \geq \left(1 - \frac{1}{e}\right) J(\pi^*)$. ■

Theorem 4. *N iterations of QVALAGG, collecting m regression examples per iteration guarantees that with probability at least $1 - \delta$*

$$\begin{aligned}
J(\hat{\pi}) &\geq J(\tilde{\pi}_{\text{OR}}) \\
&\quad - 2\sqrt{|\mathcal{A}|T} \sqrt{\varepsilon_{\text{class}} + \varepsilon_{\text{reg}} + O\left(\sqrt{\log((1/\delta)/Nm)}\right)} \\
&\quad - O\left(\frac{T^2 \log T}{\alpha N}\right)
\end{aligned}$$

where ε_{reg} is the empirical average online learning regret on the training regression examples collected over the iterations and $\varepsilon_{\text{class}}$ is the empirical regression regret of the best regressor in the policy class.

Proof: This is same as Lemma 5 ■

REFERENCES

- [1] Benjamin Charrow, Gregory Kahn, Sachin Patil, Sikang Liu, Ken Goldberg, Pieter Abbeel, Nathan Michael, and Vijay Kumar. **Information-Theoretic Planning with Trajectory Optimization for Dense 3D Mapping**. In *RSS*, 2015.
- [2] Chandra Chekuri and Martin Pal. **A recursive greedy algorithm for walks in directed graphs**. In *FOCS*, 2005.
- [3] Min Chen, Emilio Frazzoli, David Hsu, and Wee Sun Lee. **POMDP-lite for Robust Robot Planning under Uncertainty**. *preprint arXiv:1602.04875*, 2016.
- [4] Yuxin Chen, Shervin Javdani, Amin Karbasi, J. Bagnell, Siddhartha Srinivasa, and Andreas Krause. **Submodular Surrogates for Value of Information**. In *AAAI*, 2015.
- [5] Yuxin Chen, S. Hamed Hassani, and Andreas Krause. **Near-optimal bayesian active learning with correlated and noisy tests**. *CoRR*, abs/1605.07334, 2016.
- [6] Sanjiban Choudhury, Ashish Kapoor, Gireeja Ranade, and Debadeepta Dey. **Learning to Gather Information via Imitation**. *preprint arXiv:1611.04180*, 2016.
- [7] Daniel Golovin and Andreas Krause. **Adaptive submodularity: Theory and applications in active learning and stochastic optimization**. *JAIR*, 2011.
- [8] Daniel Golovin, Andreas Krause, and Debajyoti Ray. **Near-optimal bayesian active learning with noisy observations**. In *NIPS*, 2010.
- [9] Anupam Gupta, Viswanath Nagarajan, and R Ravi. **Approximation algorithms for optimal decision trees and adaptive TSP problems**. In *International Colloquium on Automata, Languages, and Programming*, 2010.
- [10] Lionel Heng, Alkis Gotovos, Andreas Krause, and Marc Pollefeys. **Efficient visual exploration and coverage with a micro aerial vehicle in unknown environments**. In *ICRA*, 2015.
- [11] Geoffrey A Hollinger and Gaurav S Sukhatme. **Sampling-based Motion Planning for Robotic Information Gathering**. In *RSS*, 2013.
- [12] Geoffrey A Hollinger, Urbashi Mitra, and Gaurav S Sukhatme. **Active classification: Theory and application to underwater inspection**. *preprint arXiv:1106.5829*, 2011.
- [13] Geoffrey A Hollinger, Brendan Englot, Franz S Hover, Urbashi Mitra, and Gaurav S Sukhatme. **Active planning for underwater inspection and the benefit of adaptivity**. *IJRR*, 2012.
- [14] Stefan Isler, Reza Sabzevari, Jeffrey Delmerico, and Davide Scaramuzza. **An Information Gain Formulation for Active Volumetric 3D Reconstruction**. In *ICRA*, 2016.
- [15] Rishabh K Iyer and Jeff A Bilmes. **Submodular optimization with submodular cover and submodular knapsack constraints**. In *NIPS*, 2013.
- [16] Shervin Javdani, Matthew Klingensmith, J. Andrew Bagnell, Nancy Pollard, and Siddhartha Srinivasa. **Efficient Touch Based Localization through Submodularity**. In *ICRA*, 2013.
- [17] Shervin Javdani, Yuxin Chen, Amin Karbasi, Andreas Krause, J. Andrew Bagnell, and Siddhartha Srinivasa. **Near Optimal Bayesian Active Learning for Decision Making**. In *AISTATS*, 2014.
- [18] Gregory Kahn, Tianhao Zhang, Sergey Levine, and Pieter Abbeel. **PLATO: Policy Learning using Adaptive Trajectory Optimization**. *preprint arXiv:1603.00622*, 2016.
- [19] Michael Koval, Nancy Pollard, and Siddhartha Srinivasa. **Pre- and Post-Contact Policy Decomposition for Planar Contact Manipulation Under Uncertainty**. In *RSS*, 2014.
- [20] Andreas Krause and Daniel Golovin. **Submodular function maximization**. *Tractability: Practical Approaches to Hard Problems*, 2012.
- [21] Andreas Krause and Carlos Guestrin. **Near-optimal Observation Selection Using Submodular Functions**. In *AAAI*, 2007.
- [22] Andreas Krause, Jure Leskovec, Carlos Guestrin, Jeanne VanBriesen, and Christos Faloutsos. **Efficient sensor placement optimization for securing large water distribution networks**. *Journal of Water Resources Planning and Management*, 2008.
- [23] Zhan Wei Lim, David Hsu, and Wee Sun Lee. **Adaptive Stochastic Optimization: From Sets to Paths**. In *NIPS*, 2015.
- [24] Zhan Wei Lim, David Hsu, and Wee Sun Lee. **Adaptive informative path planning in metric spaces**. *IJRR*, 2016.
- [25] Jan Peters and Stefan Schaal. **Policy gradient methods for robotics**. In *IROS*, 2006.
- [26] Stéphane Ross and Drew Bagnell. **Efficient Reductions for Imitation Learning**. In *AISTATS*, 2010.
- [27] Stéphane Ross and J Andrew Bagnell. **Reinforcement and imitation learning via interactive no-regret learning**. *preprint arXiv:1406.5979*, 2014.
- [28] Stéphane Ross, Joelle Pineau, Sébastien Paquet, and Brahim Chaib-Draa. **Online planning algorithms for POMDPs**. *JAIR*, 2008.
- [29] David Silver and Joel Veness. **Monte-Carlo planning in large POMDPs**. In *NIPS*, 2010.
- [30] Amarjeet Singh, Andreas Krause, Carlos Guestrin, William Kaiser, and Maxim Batalin. **Efficient Planning of Informative Paths for Multiple Robots**. In *IJCAI*, 2007.
- [31] Amarjeet Singh, Andreas Krause, and William J. Kaiser. **Nonmyopic Adaptive Informative Path Planning for Multiple Robots**. In *IJCAI*, 2009.
- [32] Adhiraj Somani, Nan Ye, David Hsu, and Wee Sun Lee. **DESPT: Online POMDP planning with regularization**. In *NIPS*, 2013.
- [33] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. **A Benchmark for the Evaluation of RGB-D SLAM Systems**. In *IROS*, 2012.
- [34] Jingjin Yu, Mac Schwager, and Daniela Rus. **Correlated orienteering problem and its application to informative path planning for persistent monitoring tasks**. In *IROS*, 2014.
- [35] Haifeng Zhang and Yevgeniy Vorobeychik. **Submodular Optimization with Routing Constraints**, 2016.