
Compatible Reward Inverse Reinforcement Learning

Alberto Maria Metelli
DEIB
Politecnico di Milano, Italy
albertomaria.metelli@polimi.it

Matteo Pirotta
Sequel Team
Inria Lille, France
matteo.pirotta@inria.fr

Marcello Restelli
DEIB
Politecnico di Milano, Italy
marcello.restelli@polimi.it

Abstract

Inverse Reinforcement Learning (IRL) is an effective approach to recover a reward function that explains the behavior of an expert by observing a set of demonstrations. This paper is about a novel model-free IRL approach that, differently from most of the existing IRL algorithms, does not require to specify a function space where to search for the expert’s reward function. Leveraging on the fact that the policy gradient needs to be zero for any optimal policy, the algorithm generates a set of basis functions that span the subspace of reward functions that make the policy gradient vanish. Within this subspace, using a second-order criterion, we search for the reward function that penalizes the most a deviation from the expert’s policy. After introducing our approach for finite domains, we extend it to continuous ones. The proposed approach is empirically compared to other IRL methods both in the (finite) Taxi domain and in the (continuous) Linear Quadratic Gaussian (LQG) and Car on the Hill environments.

1 Introduction

Imitation learning aims to learn to perform a task by observing only expert’s demonstrations. We consider the settings where only expert’s demonstrations are given, no information about the dynamics and the objective of the problem is provided (e.g., reward) or ability to query for additional samples. The main approaches solving this problem are behavioral cloning [1] and inverse reinforcement learning [2]. The former recovers the demonstrated policy by learning the state-action mapping in a supervised learning way, while inverse reinforcement learning aims to learn the reward function that makes the expert optimal. Behavioral Cloning (BC) is simple, but its main limitation is the intrinsic goal, i.e., to replicate the observed policy. This task has several limitations: it requires a huge amount of data when the environment (or the expert) is stochastic [3]; it does not provide good generalization or a description of the expert’s goal. On the contrary, Inverse Reinforcement Learning (IRL) accounts for generalization and transferability by directly learning the reward function. This information can be transferred to any new environment in which the features are well defined. As a consequence, IRL allows recovering the optimal policy a posteriori, even under variations of the environment. IRL has received a lot of attention in literature and has succeeded in several applications [e.g., 4, 5, 6, 7, 8]. However, BC and IRL are tightly related by the intrinsic relationship between reward and optimal policy. The reward function defines the space of optimal policies and to recover the reward it is required to observe/recover the optimal policy. The idea of this paper, and of some recent paper [e.g., 9, 8, 3], is to exploit the synergy between BC and IRL.

Unfortunately, also IRL approaches present issues. First, several IRL methods require solving the forward problem as part of an inner loop [e.g., 4, 5]. Literature has extensively focused on removing this limitation [10, 11, 9] in order to scale IRL to real-world applications [12, 3, 13]. Second, IRL methods generally require designing the function space by providing *features* that capture the structure of the reward function [e.g., 4, 14, 5, 10, 15, 9]. This information, provided in addition to expert’s demonstrations, is critical for the success of the IRL approach. The issue of designing the function

space is a well-known problem in supervised learning, but it is even more critical in IRL since a wrong choice might prevent the algorithm from finding good solutions to the IRL problem [2, 16], especially when linear reward models are considered. The importance of incorporating feature construction in IRL has been known in literature since a while [4] but, as far as we know, it has been explicitly addressed only in [17]. Recently, IRL literature, by mimicking supervised learning one, has focused on exploiting neural network capability of automatically constructing relevant features out of the provided data [12, 8, 13]. By exploiting a “black-box” approach, these methods do not take advantage of the structure of the underlying Markov decision process (in the phase of feature construction).

We present an IRL algorithm that constructs reward features directly from expert’s demonstrations. The proposed algorithm is model-free and does not require solving the forward problem (i.e., finding an optimal policy given a candidate reward function) as an inner step. The Compatible Reward Inverse Reinforcement Learning (CR-IRL) algorithm builds a reward function that is *compatible* with the expert’s policy. It mixes BC and IRL in order to recover the “optimal” and most “informative” reward function in the space spanned by the recovered features. Inspired by the gradient-minimization IRL approach proposed in [9], we focus on the space of reward functions that makes the policy gradient of the expert vanish. Since a zero gradient is only a necessary condition for optimality, we consider a second order optimality criterion based on the policy Hessian to rank the reward functions and finally select the best one (i.e., the one that penalizes the most a deviation from the expert’s policy).

2 Algorithm Overview

A Markov Decision Process (MDP) [18] is defined as $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, R, \gamma, \mu)$ where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{P}(s'|s, a)$ is a Markovian transition model that defines the conditional distribution of the next state s' given the current state s and the current action a , $\gamma \in [0, 1]$ is the discount factor, $R(s, a)$ is the expected reward for performing action a in state s and μ is the distribution of the initial state. The optimal policy π^* is the policy that maximizes the discounted sum of rewards $\mathbb{E}[\sum_{t=0}^{+\infty} \gamma^t R(s_t, a_t) | \pi, \mathcal{M}]$.

CR-IRL takes as input a parametric policy space $\Pi_{\Theta} = \{\pi_{\theta} : \theta \in \Theta \subseteq \mathbb{R}^k\}$ and a set of rewardless trajectories from the expert policy π^E , denoted by $\mathcal{D} = \{(s_{\tau_i,0}, a_{\tau_i,0}, \dots, s_{\tau_i,T(\tau_i)}, a_{\tau_i,T(\tau_i)})\}$, where $s_{\tau_i,t}$ is the t -th state in trajectory τ_i and $i = 1, \dots, N$. CR-IRL is a non-iterative algorithm that recovers a reward function for which the expert is optimal without requiring to specify a reward function space. It starts building the features $\{\phi_i\}$ of the value function that are compatible with policy π^E , i.e., that make the policy gradient vanish (Phase 1, see Sec. 3). This step requires a parametric representation $\pi_{\theta^E} \in \Pi_{\Theta}$ of the expert’s policy which can be obtained through behavioral cloning.¹ The choice of the policy space Π_{Θ} influences the size of the functional space used by CR-IRL for representing the value function (and the reward function) associated with the expert’s policy. In order to formalize this notion, we introduce the *policy rank*, a quantity that represents the ability of a parametric policy to reduce the dimensions of the approximation space for the value function of the expert’s policy. Once these value features have been built, they can be transformed into reward features $\{\psi_i\}$ (Phase 2 see Sec. 4) by means of the Bellman equation [18] (model-based) or reward shaping [19] (model-free). All the rewards spanned by the features $\{\psi_i\}$ satisfy the first-order necessary optimality condition [20], but we are not sure about their nature (minima, maxima or saddle points). The final step is thus to recover a reward function that is maximized by the expert’s policy (Phase 3 see Sec. 5). This is achieved by considering a second-order optimality condition, with the idea that we want the reward function that penalizes the most a deviation from the parameters of the expert’s policy π_{θ^E} . This criterion is similar in spirit to what done in [2, 4, 14], where the goal is to identify the reward function that makes the expert’s policy better than any other policy by a margin. The algorithmic structure is reported in Alg. 1.

IRL literature usually considers two different settings: optimal or sub-optimal expert. This distinction is necessary when a fixed reward space is provided. In fact, the demonstrated behavior may not be optimal under the considered reward space. In this case, the problem becomes somehow not well defined and additional “optimality” criteria are required [16]. This is not the case for CR-IRL that is able to automatically generate the space of reward functions that make the policy gradient vanish,

¹We want to stress that our primal objective is to recover the reward function since we aim to explain the motivations that guide the expert and to transfer it, not just to replicate the behavior. As explained in the introduction, we aim to exploit the synergy between BC and IRL.

thus containing also reward functions under which the recovered expert’s policy π_{θ^E} is optimal. In the rest of the paper, we will assume to have a parametric representation of the expert’s policy that we will denote for simplicity by π_{θ} .

3 Expert’s Compatible Value Features

In this section, we present the procedure to obtain the set $\{\phi_i\}_{i=1}^p$ of Expert’s CCompatible Q-features (ECO-Q) that make the policy gradient vanish² (Phase 1). We start introducing the policy gradient and the associated first-order optimality condition. We will indicate with \mathbb{T} the set of all possible trajectories, $p_{\theta}(\tau)$ the probability density of trajectory τ and $R(\tau)$ the γ -discounted trajectory reward defined as $R(\tau) = \sum_{t=0}^{T(\tau)} \gamma^t R(s_{\tau,t}, a_{\tau,t})$ that, in our settings, is obtained as a linear combination of reward features. Given a policy π_{θ} , the expected γ -discounted return for an infinite horizon MDP is:

$$J(\theta) = \int_{\mathcal{S}} d_{\mu}^{\pi_{\theta}}(s) \int_{\mathcal{A}} \pi_{\theta}(a|s) R(s, a) da ds = \int_{\mathbb{T}} p_{\theta}(\tau) R(\tau) d\tau,$$

where $d_{\mu}^{\pi_{\theta}}$ is the γ -discounted future state occupancy [21]. If π_{θ} is differentiable w.r.t. the parameter θ , the gradient of the expected reward (*policy gradient*) [21, 22, 23] is:

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{S}} \int_{\mathcal{A}} d_{\mu}^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a) da ds = \int_{\mathbb{T}} p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) R(\tau) d\tau, \quad (1)$$

where $d_{\mu}^{\pi_{\theta}}(s, a) = d_{\mu}^{\pi_{\theta}}(s) \pi_{\theta}(a|s)$ is the γ -discounted future state-action occupancy, which represents the expected discounted number of times action a is executed in state s given μ as initial state distribution and following policy π_{θ} . When π_{θ} is an optimal policy in the class of policies $\Pi_{\Theta} = \{\pi_{\theta} : \theta \in \Theta \subseteq \mathbb{R}^k\}$ then θ is a *stationary point* of the expected return and thus $\nabla_{\theta} J(\theta) = \mathbf{0}$ (*first-order necessary conditions for optimality* [20]).

We assume the space $\mathcal{S} \times \mathcal{A}$ to be a Hilbert space [24] equipped with the weighted inner product:³

$$\langle f, g \rangle_{\mu, \pi_{\theta}} = \int_{\mathcal{S}} \int_{\mathcal{A}} f(s, a) d_{\mu}^{\pi_{\theta}}(s, a) g(s, a) ds da. \quad (2)$$

When π_{θ} is optimal for the MDP, $\nabla_{\theta} \log \pi_{\theta}$ and $Q^{\pi_{\theta}}$ are *orthogonal* w.r.t. the inner product (2). We can exploit the orthogonality property to build an approximation space for the Q-function. Let $G_{\pi_{\theta}} = \{\nabla_{\theta} \log \pi_{\theta} \alpha : \alpha \in \mathbb{R}^k\}$ the subspace spanned by the gradient of the log-policy π_{θ} . From equation (1) finding an approximation space for the Q-function is equivalent to find the orthogonal complement of the subspace $G_{\pi_{\theta}}$, which in turn corresponds to find the null space of the functional:

$$\mathcal{G}_{\pi_{\theta}}[\phi] = \langle \nabla_{\theta} \log \pi_{\theta}, \phi \rangle_{\mu, \pi_{\theta}}. \quad (3)$$

We define an *Expert’s CCompatible Q-feature* as any function ϕ making the functional (3) null. This space $G_{\pi_{\theta}}^{\perp} := \text{null}(\mathcal{G}_{\pi_{\theta}})$ represents the Hilbert subspace of the features for the Q-function that are compatible with the policy π_{θ} in the sense that any Q-function optimized by policy π_{θ} can be expressed as a linear combination of those features. Section 3.2 and 3.3 describe how to compute the ECO-Q from samples in finite and continuous MDPs, respectively. The dimension of $G_{\pi_{\theta}}^{\perp}$ is typically very large since the number k of policy parameters is significantly smaller than the number of state-action pairs. A formal discussion of this issue for finite MDPs is presented in the next section.

3.1 Policy rank

The parametrization of the expert’s policy influences the size of $G_{\pi_{\theta}}^{\perp}$. Intuition suggests that the larger the number k of the parameters the more the policy is *informative* to infer the Q-function and so the reward function. This is motivated by the following rationale. Consider representing the expert’s policy using two different policy models such that one model is a superclass of the other one (for instance, assume to use linear models where the features used in the simpler model are a subset of the features used by policies in the other model). All the reward functions that make the policy gradient

²Notice that any linear combination of the ECO-Q also satisfies the first-order optimality condition.

³The inner product as defined is clearly symmetric, positive definite and linear, but there could be state-action pairs never visited, i.e., $d_{\mu}^{\pi_{\theta}}(s, a) = 0$, making $\langle f, f \rangle_{\mu, \pi_{\theta}} = 0$ for non-zero f . To ensure the properties of the inner product, we assume to compute it only on visited state-action pairs.

vanish with the rich policy model, do the same with the simpler model, while the vice versa does not hold. This suggests that complex policy models are able to reduce more the space of optimal reward function w.r.t. simpler models. This notion plays an important role for finite MDPs, i.e., MDPs where the state-action space is finite. We formalize the ability of a policy to infer the characteristics of the MDP with the concept of *policy rank*.

Definition 1. Let π_θ a policy with k parameters belonging to the class Π_Θ and differentiable in θ . The policy rank is the dimension of the space of the linear combinations of the partial derivatives of π_θ w.r.t. θ :

$$\text{rank}(\pi_\theta) = \dim(\Gamma_{\pi_\theta}), \quad \Gamma_{\pi_\theta} = \{\nabla_\theta \pi_\theta \alpha : \alpha \in \mathbb{R}^k\}.$$

A first important note is that the policy rank depends not only on the policy model Π_Θ but also on the *value* of the parameters of the policy π_θ . So the policy rank is a property of the *policy* not of the *policy model*. The following bound on the policy rank holds (the proof can be found in App. A.1).

Proposition 1. Given a finite MDP \mathcal{M} , let π_θ a policy with k parameters belonging to the class Π_Θ and differentiable in θ , then: $\text{rank}(\pi_\theta) \leq \min \{k, |\mathcal{S}||\mathcal{A}| - |\mathcal{S}|\}$.

From an intuitive point of view this is justified by the fact that $\pi_\theta(\cdot|s)$ is a probability distribution. As a consequence, for all $s \in \mathcal{S}$ the probabilities $\pi_\theta(a|s)$ must sum up to 1, removing $|\mathcal{S}|$ degrees of freedom. This has a relevant impact on the algorithm since it induces a lower bound on the dimension of the orthogonal complement $\dim(G_{\pi_\theta}^\perp) \geq \max \{|\mathcal{S}||\mathcal{A}| - k, |\mathcal{S}|\}$, thus even the most flexible policy (i.e., a policy model with a parameter for each state-action pair) cannot determine a unique reward function that makes the expert's policy optimal, leaving $|\mathcal{S}|$ degrees of freedom. It follows that it makes no sense to consider a policy with more than $|\mathcal{S}||\mathcal{A}| - |\mathcal{S}|$ parameters. The generalization capabilities enjoyed by the recovered reward function are deeply related to the choice of the policy model. Complex policies (many parameters) would require finding a reward function that explains the value of all the parameters, resulting in a possible overfitting, whereas a simple policy model (few parameters) would enforce generalization as the imposed constraints are fewer.

3.2 Construction of ECO-Q in Finite MDPs

We now develop in details the algorithm to generate ECO-Q in the case of finite MDPs. From now on we will indicate with $|\mathcal{D}|$ the number of distinct state-action pairs *visited* by the expert along the available trajectories. When the state-action space is finite the inner product (2) can be written in matrix notation as:

$$\langle \mathbf{f}, \mathbf{g} \rangle_{\mu, \pi_\theta} = \mathbf{f}^T \mathbf{D}_\mu^{\pi_\theta} \mathbf{g},$$

where \mathbf{f}, \mathbf{g} and $\mathbf{d}_\mu^{\pi_\theta}$ are real vectors with $|\mathcal{D}|$ components and $\mathbf{D}_\mu^{\pi_\theta} = \text{diag}(\mathbf{d}_\mu^{\pi_\theta})$. The term $\nabla_\theta \log \pi_\theta$ is a $|\mathcal{D}| \times k$ real matrix, thus finding the null space of the functional (3) is equivalent to finding the null space of the matrix $\nabla_\theta \log \pi_\theta^T \mathbf{D}_\mu^{\pi_\theta}$. This can be done for instance through SVD which allows to obtain a set of orthogonal basis functions Φ . Given that the weight vector $d_\mu^{\pi_\theta}(s, a)$ is usually unknown, it needs to be estimated. Since the policy π_θ is known, we need to estimate just $d_\mu^{\pi_\theta}(s)$, as $d_\mu^{\pi_\theta}(s, a) = d_\mu^{\pi_\theta}(s) \pi_\theta(a|s)$. A Monte Carlo estimate exploiting the expert's demonstrations in \mathcal{D} is:

$$\hat{d}_\mu^{\pi_\theta}(s) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T(\tau_i)} \gamma^t \mathbb{1}(s_{\tau_i, t} = s). \quad (4)$$

3.3 Construction of ECO-Q in Continuous MDPs

To extend the previous approach to the continuous domain we assume that the state-action space is equipped with the Euclidean distance. Now we can adopt an approach similar to the one exploited to extend Proto-Value Functions (PVF) [25, 26] to infinite observation spaces [27]. The problem is treated as a discrete one considering only the state-action pairs visited along the collected trajectories. A Nyström interpolation method is used to approximate the value of a feature in a non-visited state-action pair as a weighted mean of the values of the closest k features. The weight of each feature is computed by means of a Gaussian kernel placed over the Euclidean space $\mathcal{S} \times \mathcal{A}$:

$$\mathcal{K}((s, a), (s', a')) = \exp\left(-\frac{1}{2\sigma_S^2} \|s - s'\|_2^2 - \frac{1}{2\sigma_A^2} \|a - a'\|_2^2\right), \quad (5)$$

where σ_S and σ_A are respectively the state and action bandwidth. In our setting this approach is fully equivalent to a kernel k-Nearest Neighbors regression.

4 Expert’s Compatible Reward Features

The set of ECO-Q basis functions allows representing the optimal value function under the policy π_θ . In this section, we will show how it is possible to exploit ECO-Q functions to generate basis functions for the reward representation (Phase 2). In principle, we can use the Bellman equation to obtain the reward from the Q-function but this approach requires the knowledge of the transition model (see App. B). The reward can be recovered in a model-free way by exploiting optimality-invariant reward transformations.

Reversing the Bellman equation [e.g., 10] allows finding the reward space that generates the estimated Q-function. However, IRL is interested in finding just a reward space under which the expert’s policy is optimal. This problem can be seen as an instance of reward shaping [19] where the authors show that the space of all the reward functions sharing the same optimal policy is given by:

$$R'(s, a) = R(s, a) + \gamma \int_S \mathcal{P}(s'|s, a) \chi(s') ds' - \chi(s),$$

where $\chi(s)$ is a state-dependent potential function. A smart choice [19] is to set $\chi = V^{\pi_\theta}$ under which the new reward space is given by the advantage function: $R'(s, a) = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) = A^{\pi_\theta}(s, a)$. Thus the expert’s advantage function is an admissible reward optimized by the expert’s policy itself. This choice is, of course, related to using Q^{π_θ} as reward. However, the advantage function encodes a more local and more transferable information w.r.t. the Q-function.

The space of reward features can be recovered through matrix equality $\Psi = (\mathbf{I} - \tilde{\pi}_\theta)\Phi$, where $\tilde{\pi}_\theta$ is a $|\mathcal{D}| \times |\mathcal{D}|$ matrix obtained from π_θ repeating the row of each visited state a number of times equal to the number of distinct actions performed by the expert in that state. Notice that this is a simple linear transformation through the expert’s policy. The specific choice of the state-potential function has the advantage to improve the learning capabilities of any RL algorithm [19]. This is not the only choice of the potential function possible, but it has the advantage of allowing model-free estimation.

Once the ECO-R basis functions have been generated, they can be used to feed any IRL algorithm that represents the expert’s reward through a linear combination of basis functions. In the next section, we propose a new method based on the optimization of a second-order criterion that favors reward functions that significantly penalize deviations from the expert’s policy.

5 Reward Selection via Second-Order Criterion

Any linear combination of the ECO-R $\{\psi_i\}_{i=1}^p$ makes the gradient vanish, however in general this is not sufficient to ensure that the policy parameter θ is a maximum of $J(\theta)$. Combinations that lead to minima or saddle points should be discarded. Furthermore, provided that a subset of ECO-R leading to maxima has been selected, we should identify a single reward function in the space spanned by this subset of features (Phase 3). Both these requirements can be enforced by imposing a second-order optimality criterion based on the policy Hessian that is given by [28, 29]:

$$\mathcal{H}_\theta J(\theta, \omega) = \int_{\mathbb{T}} p_\theta(\tau) \left(\nabla_\theta \log p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)^T + \mathcal{H}_\theta \log p_\theta(\tau) \right) R(\tau, \omega) d\tau,$$

where ω is the reward weight and $R(\tau, \omega) = \sum_{i=1}^p \omega_i \sum_{t=0}^{T(\tau)} \gamma^t \psi_i(s_{\tau,t}, a_{\tau,t})$.

In order to retain only maxima we need to impose that the Hessian is negative definite. Furthermore, we aim to find the reward function that best represents the optimal policy parametrization in the sense that even a slight change of the parameters of the expert’s policy induces a significant degradation of the performance. Geometrically this corresponds to find the reward function for which the expected return locally represents the sharpest hyper-paraboloid. These requirements can be enforced using a Semi-Definite Programming (SDP) approach where the objective is to *minimize the maximum eigenvalue* of the Hessian whose eigenvector corresponds to the direction of minimum curvature (*maximum eigenvalue optimality* criterion). This problem is not appealing in practice due to its high computational burden. Furthermore, it might be the case that the strict negative definiteness constraint is never satisfied due to blocked-to-zero eigenvalues (for instance in presence of policy parameters that do not affect the policy performance). In these cases, we can consider maximizing an index of the overall concavity. The trace of the Hessian, being the sum of the eigenvalues, can be used for this purpose. This problem can be still defined as a SDP problem (*trace optimality* criterion). See App. C for details.

Trace optimality criterion, although less demanding w.r.t. the eigenvalue-based one, still displays performance degradation as the number of basis functions increases due to the negative definiteness constraint. Solving the semidefinite programming problem of one of the previous optimality criteria is unfeasible for almost all the real world problems. We are interested in formulating a non-SDP problem, which is a surrogate of the trace optimality criterion, that can be solved more efficiently (*trace heuristic* criterion). In our framework, the reward function can be expressed as a linear combination of the ECO-R so we can rewrite the Hessian as $\mathcal{H}_\theta J(\theta, \omega) = \sum_{i=1}^p \omega_i \mathcal{H}_\theta J_i(\theta)$ where $J_i(\theta)$ is the expected return considering as reward function ψ_i . We assume that the ECO-R are orthonormal in order to compare them.⁴ The main challenge is how to select the weight ω in order to get a (sub-)optimal trace minimizer that preserves the negative semidefinite constraint. From Weyl's inequality, we get a feasible solution by retaining only the ECO-Rs yielding a semidefinite Hessian and switching sign to those with positive semidefinite Hessian. Our heuristic consists in looking for the weights ω that minimize the trace in this reduced space (in which all ECO-R have a negative semidefinite Hessian). Notice that in this way we can loose the optimal solution since the trace minimizer might assign a non-zero weight to a ECO-R with indefinite Hessian. For brevity, we will indicate with $tr_i = \text{tr}(\mathcal{H}_\theta J_i(\theta))$ and tr the vector whose components are tr_i . SDP is no longer needed:

$$\min_{\omega} \quad \omega^T \text{tr} \quad \text{s.t.} \quad \|\omega\|_2^2 = 1. \quad (6)$$

The constraint $\|\omega\|_2^2 = 1$ ensures that, when the ECO-R are orthonormal, the resulting ECO-R has Euclidean norm one. This is a convex programming problem with linear objective function and quadratic constraint, the closed form solution can be found with Lagrange multipliers: $\omega_i = -\frac{tr_i}{\|\text{tr}\|_2}$ (see App. A.2 for the derivation). Refer to Algorithm 1 for a complete overview of CR-IRL (the computational analysis of CR-IRL is reported in App. E).

CR-IRL does not assume to know the state space \mathcal{S} and the action space \mathcal{A} , thus the recovered reward is defined only in the state-action pairs visited by the expert along the trajectories in \mathcal{D} . When the state and action spaces are known, we can complete the reward function also for unexplored state-action pairs assigning a penalized reward (e.g., a large negative value), otherwise the penalization can be performed online when the recovered reward is used to solve the forward RL problem.

6 Related Work

There has been a surge of recent interest in improving IRL in order to make it more appealing for real-world applications. We highlight the lines of works that are more related to this paper.

We start investigating how IRL literature has faced the problem of designing a suitable reward space. Almost all the IRL approaches share the necessity to define *a priori* a set of handcrafted features,

⁴A normalization condition is necessary since the magnitude of the trace of a matrix can be arbitrarily changed by multiplying the matrix by a constant.

Input: $\mathcal{D} = \{(s_{\tau_i,0}, a_{\tau_i,0}, \dots, s_{\tau_i,T(\tau_i)}, a_{\tau_i,T(\tau_i)})\}_{i=1}^N$ a set of expert's trajectories and parametric expert's policy π_θ .	
Output: trace heuristic ECO-R, $R^{\text{tr-heu}}$.	
Phase 1	
1. Estimate $d_\mu^{\pi_\theta}(s)$ for the visited state-action pairs using Eq. (4) and compute $d_\mu^{\pi_\theta}(s, a) = d_\mu^{\pi_\theta}(s) \pi_\theta(a s)$.	
2. Collect $d_\mu^{\pi_\theta}(s, a)$ in the $ \mathcal{D} \times \mathcal{D} $ diagonal matrix $\mathbf{D}_\mu^{\pi_\theta}$ and $\nabla_\theta \log \pi_\theta(s, a)$ in the $ \mathcal{D} \times k$ matrix $\nabla_\theta \log \pi_\theta$.	
3. Get the set of ECO-Q by computing the null space of matrix $\nabla_\theta \log \pi_\theta^T \mathbf{D}_\mu^{\pi_\theta}$ through SVD: $\Phi = \text{null}(\nabla_\theta \log \pi_\theta^T \mathbf{D}_\mu^{\pi_\theta})$.	
4. Get the set of ECO-R by applying reward shaping to the set of ECO-Q: $\Psi = (\mathbf{I} - \tilde{\pi}_\theta) \Phi$.	
5. Apply SVD to orthogonalize Ψ .	Phase 2
6. Estimate the policy Hessian for each ECO-R ψ_i , $i = 1, \dots, p$ using equation: ^a	
$\hat{\mathcal{H}}_\theta J_i(\theta) = \frac{1}{N} \sum_{j=1}^N (\nabla_\theta \log p_\theta(\tau_j) \nabla_\theta \log p_\theta(\tau_j)^T + \mathcal{H}_\theta \log p_\theta(\tau_j)) (\psi_i(\tau_j) - b).$	
7. Discard the ECO-R having indefinite Hessian, switch sign for those having positive semidefinite Hessian, compute the traces of each Hessian and collect them in the vector tr .	
8. Compute the trace heuristic ECO-R as:	
$R^{\text{tr-heu}} = \Psi \omega \quad , \quad \omega = -\text{tr} / \ \text{tr}\ _2.$	
9. (Optional) Apply penalization to unexplored state-action pairs.	Phase 3
^a The optimal baseline b is provided in [30, 31].	

Alg 1: CR-IRL algorithm.

spanning the approximation space of the reward functions. While a good set of basis functions can greatly simplify the IRL problem, a bad choice may significantly harm the performance of any IRL algorithm. The Feature construction for Inverse Reinforcement Learning (FIRL) algorithm [17], as far as we know, is the only approach that explicitly incorporates the feature construction as an inner step. FIRL alternates between optimization and fitting phases. The optimization phase aims to recover a reward function—from the current feature set as a linear projection—such that the associated optimal policy is consistent with the demonstrations. In the fitting phase new features are created (using a regression tree) in order to better explain regions where the old features were too coarse. The method proved to be effective achieving also (features) transfer capabilities. However, FIRL requires the MDP model to solve the forward problem and the complete optimal policy for the fitting step in order to evaluate the consistency with demonstrations.

Recent works have indirectly coped with the feature construction problem by exploiting neural networks [12, 3, 13]. Although effective, the black-box approach does not take into account the MDP structure of the problem. RL has extensively investigated the feature construction for the forward problem both for value function [25, 26, 32, 33] and policy [21] features. In this paper, we have followed this line of work mixing concepts deriving from policy and value fields. We have leveraged on the policy gradient theorem and on the associated concept of compatible functions to derive ECO-Q features. First-order necessary conditions have already been used in literature to derive IRL algorithm [9, 34]. However, in both the cases the authors assume a fixed reward space under which it may not be possible to find a reward for which the expert is optimal. Although there are similarities, this paper exploits first-order optimality to recover the reward basis while the “best” reward function is selected according to a second-order criterion. This allows recovering a more robust solution overcoming uncertainty issues raised by the use of the first-order information only.

7 Experimental results

We evaluate CR-IRL against some popular IRL algorithms both in discrete and in continuous domains: the Taxi problem (discrete), the Linear Quadratic Gaussian and the Car on the Hill environments (continuous). We provide here the most significant results, the full data are reported in App. D.

7.1 Taxi

The Taxi domain is defined in [35]. We assume the expert plays an ϵ -Boltzmann policy with fixed ϵ :

$$\pi_{\theta, \epsilon}(a|s) = (1 - \epsilon) \frac{e^{\theta_a^T \zeta_s}}{\sum_{a' \in \mathcal{A}} e^{\theta_{a'}^T \zeta_s}} + \frac{\epsilon}{|\mathcal{A}|},$$

where the policy features ζ_s are the following state features: current location, passenger location, destination location, whether the passenger has already been pick up.

This test is meant to compare the learning speed of the reward functions recovered by the considered IRL methods when a Boltzmann policy ($\epsilon = 0$) is trained with REINFORCE [22]. To evaluate the robustness to imperfect experts, we introduce a noise (ϵ) in the optimal policy. Figure 2 shows that CR-IRL, with 100 expert’s trajectories, outperforms the true reward function in terms of convergence speed regardless the exploration level. Behavioral Cloning (BC), obtained by recovering the maximum likelihood ϵ -Boltzmann policy ($\epsilon = 0, 0.1$) from expert’s trajectories, is very susceptible to noise.

We compare also the second-order criterion of CR-IRL to single out the reward function with Maximum Entropy IRL (ME-IRL) [6] and Linear Programming Apprenticeship Learning (LPAL) [5] using as reward features the set of ECO-R (comparisons with different sets of features is reported in App. D.2). We can see in Figure 2 that ME-IRL does not perform well when $\epsilon = 0$, since the transition model is badly estimated. The convergence speed remains very slow also for $\epsilon = 0.1$, since ME-IRL does not guarantee that the recovered reward is a maximum of J . LPAL provides as output an apprenticeship policy (not a reward function) and, like BC, is very sensitive to noise and to the quality of the estimated transition model.

7.2 Linear Quadratic Gaussian Regulator

We consider the one-dimensional Linear Quadratic Gaussian regulator [36] with an expert playing a Gaussian policy $\pi_K(\cdot|s) \sim \mathcal{N}(Ks, \sigma^2)$, where K is the parameter and σ^2 is fixed.

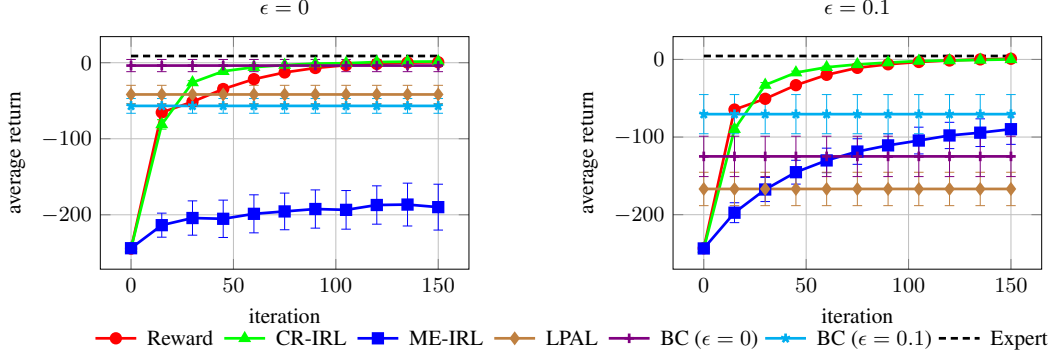


Figure 2: Average return of the Taxi problem as a function of the number of iterations of REINFORCE.

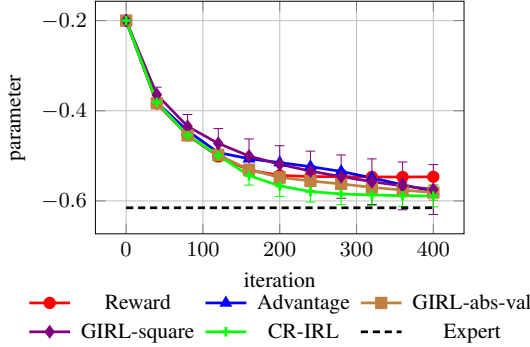


Figure 3: Parameter value of LQG as a function of the number of iterations of REINFORCE.

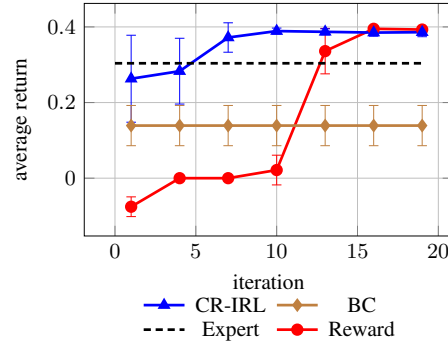


Figure 4: Average return of Car on the Hill as a function of the number of FQI iterations.

We compare CR-IRL with GIRL [9] using two linear parametrizations of the reward function: $R(s, a, \omega) = \omega_1 s^2 + \omega_2 a^2$ (GIRL-square) and $R(s, a, \omega) = \omega_1 |s| + \omega_2 |a|$ (GIRL-abs-val). Figure 3 shows the parameter (K) value learned with REINFORCE using a Gaussian policy with variance $\sigma^2 = 0.01$. We notice that CR-IRL, fed with 20 expert's trajectories, converges closer and faster to the expert's parameter w.r.t. to the true reward, advantage function and GIRL with both parametrizations.

7.3 Car on the Hill

We further experiment CR-IRL in the continuous Car on the Hill domain [37]. We build the optimal policy via FQI [37] and we consider a noisy expert's policy in which a random action is selected with probability $\epsilon = 0.1$. We exploit 20 expert's trajectories to estimate the parameters \mathbf{w} of a Gaussian policy $\pi_{\mathbf{w}}(a|s) \sim \mathcal{N}(y_{\mathbf{w}}(s), \sigma^2)$ where the mean $y_{\mathbf{w}}(s)$ is a radial basis function network (details and comparison with $\epsilon = 0.2$ in appendix D.4). The reward function recovered by CR-IRL does not necessary need to be used only with policy gradient approaches. Here we compare the average return as a function of the number of iterations of FQI, fed with the different recovered rewards. Figure 4 shows that FQI converges faster to optimal policies when coped with the reward recovered by CR-IRL rather than with the original reward. Moreover, it overcomes the performance of the policy recovered via BC.

8 Conclusions

We presented an algorithm, CR-IRL, that leverages on the policy gradient to recover, from a set of expert's demonstrations, a reward function that explains the expert's behavior and penalizes deviations. Differently from large part of IRL literature, CR-IRL does not require to specify a priori an approximation space for the reward function. The empirical results show (quite unexpectedly) that the reward function recovered by our algorithm allows learning policies that outperform both behavioral cloning and those obtained with the true reward function (learning speed). Furthermore, the Hessian trace heuristic criterion, when applied to ECO-R, outperforms classic IRL methods.

Acknowledgments

This research was supported in part by French Ministry of Higher Education and Research, Nord-Pas-de-Calais Regional Council and French National Research Agency (ANR) under project ExTra-Learn (n.ANR-14-CE24-0010-01).

References

- [1] Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5):469–483, 2009.
- [2] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *ICML*, pages 663–670, 2000.
- [3] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *NIPS*, pages 4565–4573, 2016.
- [4] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *ICML*, page 1. ACM, 2004.
- [5] Umar Syed, Michael H. Bowling, and Robert E. Schapire. Apprenticeship learning using linear programming. In *ICML*, volume 307 of *ACM International Conference Proceeding Series*, pages 1032–1039. ACM, 2008.
- [6] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. In *AAAI*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- [7] Nathan D. Ratliff, David Silver, and J. Andrew Bagnell. Learning to search: Functional gradient techniques for imitation learning. *Autonomous Robots*, 27(1):25–53, 2009.
- [8] Jonathan Ho, Jayesh K. Gupta, and Stefano Ermon. Model-free imitation learning with policy optimization. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 2760–2769. JMLR.org, 2016.
- [9] Matteo Pirota and Marcello Restelli. Inverse reinforcement learning through policy gradient minimization. In *AAAI*, pages 1993–1999, 2016.
- [10] Edouard Klein, Bilal Piot, Matthieu Geist, and Olivier Pietquin. A cascaded supervised learning approach to inverse reinforcement learning. In *ECML/PKDD (I)*, volume 8188 of *Lecture Notes in Computer Science*, pages 1–16. Springer, 2013.
- [11] Bilal Piot, Matthieu Geist, and Olivier Pietquin. Boosted and reward-regularized classification for apprenticeship learning. In *AAMAS*, pages 1249–1256. IFAAMAS/ACM, 2014.
- [12] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 49–58. JMLR.org, 2016.
- [13] Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Andrew Sendonaris, Gabriel Dulac-Arnold, Ian Osband, John Agapiou, Joel Z. Leibo, and Audrunas Gruslys. Learning from demonstrations for real world reinforcement learning. *CoRR*, abs/1704.03732, 2017.
- [14] Nathan D. Ratliff, J. Andrew Bagnell, and Martin Zinkevich. Maximum margin planning. In *ICML*, volume 148 of *ACM International Conference Proceeding Series*, pages 729–736. ACM, 2006.
- [15] Julien Audiffren, Michal Valko, Alessandro Lazaric, and Mohammad Ghavamzadeh. Maximum entropy semi-supervised inverse reinforcement learning. In *IJCAI*, pages 3315–3321. AAAI Press, 2015.
- [16] Gergely Neu and Csaba Szepesvári. Training parsers by inverse reinforcement learning. *Machine Learning*, 77(2-3):303–337, 2009.
- [17] Sergey Levine, Zoran Popovic, and Vladlen Koltun. Feature construction for inverse reinforcement learning. In *NIPS*, pages 1342–1350. Curran Associates, Inc., 2010.
- [18] Martin L Puterman. Markov decision processes: Discrete stochastic dynamic programming. 1994.
- [19] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. 99:278–287, 1999.

- [20] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York, 2006.
- [21] Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, pages 1057–1063. The MIT Press, 1999.
- [22] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [23] Jan Peters and Stefan Schaal. Reinforcement learning of motor skills with policy gradients. *Neural Networks*, 21(4):682–697, 2008.
- [24] Wendelin Böhmer, Steffen Grünewälder, Yun Shen, Marek Musial, and Klaus Obermayer. Construction of approximation spaces for reinforcement learning. *Journal of Machine Learning Research*, 14(1):2067–2118, 2013.
- [25] Sridhar Mahadevan. Proto-value functions: Developmental reinforcement learning. In *ICML*, pages 553–560. ACM, 2005.
- [26] Sridhar Mahadevan and Mauro Maggioni. Proto-value functions: A laplacian framework for learning representation and control in markov decision processes. *Journal of Machine Learning Research*, 8(Oct):2169–2231, 2007.
- [27] Sridhar Mahadevan, Mauro Maggioni, Kimberly Ferguson, and Sarah Osentoski. Learning representation and control in continuous markov decision processes. In *AAAI*, volume 6, pages 1194–1199, 2006.
- [28] Sham Kakade. A natural policy gradient. In *NIPS*, pages 1531–1538. MIT Press, 2001.
- [29] Thomas Furnstun and David Barber. A unifying perspective of parametric policy search methods for markov decision processes. In *Advances in neural information processing systems*, pages 2717–2725, 2012.
- [30] Giorgio Manganini, Matteo Pirotta, Marcello Restelli, and Luca Bascetta. Following newton direction in policy gradient with parameter exploration. In *Neural Networks (IJCNN), 2015 International Joint Conference on*, pages 1–8. IEEE, 2015.
- [31] Simone Parisi, Matteo Pirotta, and Marcello Restelli. Multi-objective reinforcement learning through continuous pareto manifold approximation. *Journal Artificial Intelligence Research*, 57:187–227, 2016.
- [32] Ronald Parr, Christopher Painter-Wakefield, Lihong Li, and Michael L. Littman. Analyzing feature generation for value-function approximation. In *ICML*, volume 227 of *ACM International Conference Proceeding Series*, pages 737–744. ACM, 2007.
- [33] Amir Massoud Farahmand and Doina Precup. Value pursuit iteration. In *NIPS*, pages 1349–1357, 2012.
- [34] Peter Englert and Marc Toussaint. Inverse kkt–learning cost functions of manipulation tasks from demonstrations. In *Proceedings of the International Symposium of Robotics Research*, 2015.
- [35] Thomas G Dietterich. Hierarchical reinforcement learning with the maxq value function decomposition. *J. Artif. Intell. Res.(JAIR)*, 13:227–303, 2000.
- [36] Peter Dorato, Vito Cerone, and Chaouki Abdallah. *Linear Quadratic Control: An Introduction*. Krieger Publishing Co., Inc., Melbourne, FL, USA, 2000.
- [37] Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556, 2005.
- [38] C-L Hwang and Abu Syed Md Masud. *Multiple objective decision making-methods and applications: a state-of-the-art survey*, volume 164. Springer Science & Business Media, 2012.
- [39] Jose M. Vidal and José M Vidal. Fundamentals of multiagent systems. 2006.
- [40] Emre Mengi, E Alper Yildirim, and Mustafa Kilic. Numerical optimization of eigenvalues of hermitian matrix functions. *SIAM Journal on Matrix Analysis and Applications*, 35(2):699–724, 2014.
- [41] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.

A Proofs and derivations

A.1 Proof of Proposition 1

Proof. $\Gamma_{\pi_\theta} = \{\nabla_\theta \pi_\theta \alpha : \alpha \in \mathbb{R}^k\}$ is a subspace of $\mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ generated by k vectors, thus its dimension cannot be larger than k . For each state $s \in \mathcal{S}$ it holds that $\pi_\theta(\cdot|s)$ is a probability density function thus:

$$\sum_{a \in \mathcal{A}} \pi_\theta(a|s) = 1, \quad \forall s \in \mathcal{S}.$$

For all parameters $\theta_j, j = 1, \dots, k$ the partial derivatives of the previous equation result in:

$$\sum_{a \in \mathcal{A}} \frac{\partial \pi_\theta}{\partial \theta_j}(a|s) = 0, \quad \forall s \in \mathcal{S}, \quad j = 1, 2, \dots, k.$$

Let us consider $h \in \Gamma_{\pi_\theta}$, by definition of Γ_{π_θ} , it can be written as a linear combination of the partial derivatives of π_θ :

$$h = \sum_{j=1}^k \alpha_j \frac{\partial \pi_\theta}{\partial \theta_j},$$

so for all the states s :

$$\sum_{a \in \mathcal{A}} h(a|s) = 0, \quad \forall s \in \mathcal{S}.$$

These are $|\mathcal{S}|$ linearly independent equations, thus the dimension of Γ_{π_θ} cannot be larger than $|\mathcal{S}||\mathcal{A}| - |\mathcal{S}|$. \square

A.2 Derivation of closed form solution of problem (6)

For two symmetric matrices A and B , it holds from Weyl's inequality that $\lambda_{\max}(A) + \lambda_{\max}(B) \leq \lambda_{\max}(A + B)$, thus any linear combination of semidefinite matrices with positive coefficients for the negative semidefinite ones and negative coefficients for the positive semidefinite ones is guaranteed to yield a negative semidefinite matrix. In principle, nothing can be ensured for the indefinite matrices, without looking to the magnitude of the eigenvalues. In the following derivation, we discard ECO-Rs yielding indefinite Hessian and we assume that the remaining yield negative semidefinite Hessian (for the ECO-Rs that yield a positive semidefinite Hessian we just need to switch the sign).

Proof. The problem can be solved in closed form using Lagrange multipliers. The Lagrangian function is:

$$\mathcal{L}(\omega, \alpha) = \sum_{i=1}^p \omega_i \text{tr}_i + \alpha \left(\sum_{i=1}^p \omega_i^2 - 1 \right),$$

where α is the Lagrange multiplier. We impose that partial derivatives of \mathcal{L} w.r.t. ω_j and α vanish, thus:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \omega_j} &= \text{tr}_j + 2\alpha \omega_j = 0, \quad j = 1, 2, \dots, p, \\ \frac{\partial \mathcal{L}}{\partial \alpha} &= \sum_{i=1}^p \omega_i^2 - 1 = 0. \end{aligned}$$

From the first equation, for $\alpha \neq 0$, we obtain an expression of ω_j as a function of α :

$$\omega_j = -\frac{\text{tr}_j}{2\alpha}, \quad j = 1, 2, \dots, p.$$

The case $\alpha = 0$ is uninteresting since it is feasible only if all traces are null. By substitution we obtain an expression for α :

$$\alpha = \pm \frac{1}{2} \sqrt{\sum_{i=1}^p \text{tr}_i^2} = \pm \frac{1}{2} \|\mathbf{tr}\|_2.$$

Since we are looking for a minimum, provided that all traces are non positive, the objective function is minimized for non negative weights:

$$\omega_j = -\frac{tr_j}{\sqrt{\sum_{i=1}^p tr_i^2}} = -\frac{tr_j}{\|\mathbf{tr}\|_2}, \quad j = 1, 2, \dots, p.$$

□

B Model-Based Construction of Reward Features

The relation between the reward and the Q-function is given by the Bellman equation:

$$Q^{\pi_\theta}(s, a) = R(s, a) + \gamma \int_S \int_{\mathcal{A}} \mathcal{P}(s'|s, a) \pi_\theta(a'|s') Q^{\pi_\theta}(s', a') ds' da', \quad (7)$$

which can be expressed for the finite case in matrix form as $\mathbf{q}^{\pi_\theta} = \mathbf{r} + \gamma \mathbf{P}^{\pi_\theta} \mathbf{q}^{\pi_\theta}$. Similarly to [10], we can obtain the reward function by reversing the equation: $\mathbf{r} = (\mathbf{I} - \gamma \mathbf{P}^{\pi_\theta}) \mathbf{q}^{\pi_\theta}$. Thus the set of ECO-R $\{\psi_i\}_{i=1}^p$ are obtained by means of the linear mapping $\Psi = (\mathbf{I} - \gamma \mathbf{P}^{\pi_\theta}) \Phi$ of the ECO-Q $\{\phi_i\}_{i=1}^p$. This transformation however requires the knowledge of the transition model \mathcal{P} which in most of the real cases cannot be assumed. As shown in [10], in order to exploit only expert's demonstrations it is necessary to resort to heuristics that may be hard to compute in practice or additional samples for exploration of the system dynamics.

C Second-Order Optimality Criteria

In this Appendix, we report the SDP formulation of the maximum eigenvalue and trace optimality criteria introduced in Section 5.

C.1 Maximum eigenvalue optimality criterion

The minimization of the curvature in the flattest direction of the expected return corresponds analytically to the minimization of the maximum eigenvalue of the Hessian. The problem can be formulated using semidefinite programming:

$$\min_{\omega} \quad \lambda_{\max}(\mathcal{H}_\theta J(\theta, \omega)) \quad \text{s.t.} \quad \mathcal{H}_\theta J(\theta, \omega) + \epsilon \mathbf{I} \preceq 0,$$

where $\epsilon \geq 0$ is a threshold to ensure the Hessian is (strictly) negative definite; a further normalization constraint over the weights (like $\|\omega\|_2 = 1$) might be added to ensure boundedness.

This optimization problem is for sure feasible (for sufficiently small ϵ) since, as already seen, the true reward function and the advantage function make the policy parametrization optimal. However, in most of the cases, it is impractical to solve, at least for two reasons. First, the computational effort is enormous even for few policy parameters. Second, it might be the case that the strict negative definiteness constraint is never satisfied due to blocked-to-zero eigenvalues. This problem arises quite often and is related to the presence of “useless” policy parameters that even if modified do not affect the policy performance. In those cases $\lambda_{\max}(\mathcal{H}_\theta J(\theta, \omega))$ would be zero no matter which reward weights ω are selected.

C.2 Trace optimality criterion

In order to account for the possible presence of blocked-to-zero eigenvalues, we need to consider also the other eigenvalues of the Hessian matrix instead of the maximum one only. Minimizing the trace of the Hessian can be formulated in the following semidefinite programming problem:

$$\min_{\omega} \quad \text{tr}(\mathcal{H}_\theta J(\theta, \omega)) \quad \text{s.t.} \quad \mathcal{H}_\theta J(\theta, \omega) + \epsilon \mathbf{I} \preceq 0,$$

Here the negative definiteness constraint is substantial to ensure the optimality, since there might be ECO-R with negative trace but positive eigenvalues. Moreover, in presence of blocked-to-zero eigenvalues, ϵ should be set to zero to ensure feasibility.

C.3 Multi-Objective interpretation of the second-order criteria

The intuitive idea behind the second-order criteria we presented in the previous sections consists in preferring rewards that penalize the most deviations from the expert’s policy. This notion can be formalized, in the general case, as the multi-objective SDP problem of minimizing the vector of eigenvalues of the policy Hessian:

$$\min_{\omega} \lambda(\omega) = (\lambda_1(\omega), \lambda_2(\omega), \dots, \lambda_k(\omega)) \quad \text{s.t.} \quad \mathcal{H}_{\theta} J(\theta, \omega) + \epsilon \mathbf{I} \preceq 0,$$

where $\lambda_i(\omega) = \lambda_i(\mathcal{H}_{\theta} J(\theta, \omega))$ for $i = 1, 2, \dots, k$ is the i -th largest eigenvalue of the policy Hessian. Clearly, among all possible feasible solutions we seek for the (strict) Pareto optimal ones ω^P , i.e., those for which there exists no feasible weight vector ω such that $\lambda_i(\omega) \leq \lambda_i(\omega^P)$, $i = 1, 2, \dots, k$ with at least one strict inequality. It is well-known that the Pareto frontier (the set of all the Pareto optimal solutions) cannot be computed efficiently in most of the cases. A standard approach to tackle this problem is *scalarization* [38]. Scalarization consists in formulating a single-objective optimization problem such that the optimal solutions to the single-objective problem are Pareto optimal solutions to the multi-objective problem and vice versa. A common choice is linear scalarization, that consists in combining the multiple objectives via a linear function:

$$L(\lambda(\omega), \beta) = \sum_{i=1}^k \beta_i \lambda_i(\omega) = \beta^T \lambda(\omega). \quad (8)$$

It can be proved that the minimizer of the scalar objective $L(\lambda(\omega), \beta)$ is a Pareto optimal solution for every value of $\beta \geq 0$. However, this scalarization is guaranteed to yield all Pareto optimal solutions only when the Pareto frontier is convex, which is not our case as, for instance, $\lambda_k(\omega)$ is concave.

It is worth noting that the maximum eigenvalue and the trace optimality criteria fall in this formulation. The former is obtained by setting $\beta_1 = 1$ and $\beta_i = 0$, $i = 2, 3, \dots, k$, while the latter is obtained by setting $\beta_i = 1$, $i = 1, 2, \dots, k$. Therefore, we are sure that those criteria provide Pareto optimal solutions. Furthermore, the two criteria correspond to popular solution concepts in multi-agent decision theory [39]. The maximum eigenvalue optimality criterion recovers the egalitarian social welfare solution, i.e., the solution that minimizes the maximum unsatisfaction, whereas the trace optimality criterion seeks for the utilitarian solution, i.e., the solution that maximizes the sum of profits. The trace heuristic, however, does not guarantee that the recovered solution is Pareto optimal.

Furthermore, the single-objective function (8) is not convex for all values of β . It can be proved that a sufficient condition for convexity is that $\beta_1 \geq \beta_2 \geq \dots \geq \beta_k \geq 0$ [40], to which both maximum eigenvalue and trace optimality criteria comply.

D Experimental results

D.1 Estimators

For the model-based IRL algorithms, the transition model is estimated using the available expert’s trajectories via the Monte Carlo estimator:

$$\hat{\mathcal{P}}(s'|s, a) = \frac{\sum_{i=1}^N \sum_{t=0}^{T(\tau_i)-1} \mathbb{1}(s_{\tau_i, t+1} = s', a_{\tau_i, t} = a, s_{\tau_i, t} = s)}{\sum_{i=1}^N \sum_{t=0}^{T(\tau_i)-1} \mathbb{1}(a_{\tau_i, t} = a, s_{\tau_i, t} = s)}. \quad (9)$$

The Kullback-Leibler divergence between the expert’s policy π_{θ}^E and the recovered policy $\pi_{\hat{\theta}}$ is estimated from samples as:

$$\hat{d}_{KL}(\pi_{\theta}^E, \pi_{\hat{\theta}}) = \sum_{i=1}^M \frac{1}{T(\tau_i)} \sum_{t=0}^{T(\tau_i)} \log \left(\frac{\pi_{\theta}^E(a_{\tau_i, t} | s_{\tau_i, t})}{\pi_{\hat{\theta}}(a_{\tau_i, t} | s_{\tau_i, t})} \right), \quad (10)$$

where the M trajectories are collected with the expert’s policy π_{θ}^E .

D.2 Taxi

The environment corresponds to a 5x5 grid in which there are 4 locations, labeled by different letters. The job of the taxi driver is to pick up the passenger at one location and drop him off in another. You

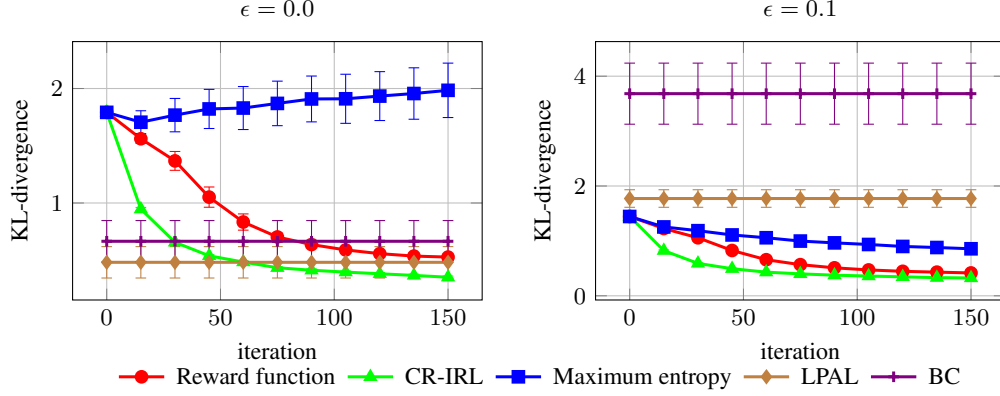


Figure 5: Kullback-Leibler divergence between expert’s policy and learned policy in the Taxi problem as a function of the number of iterations of REINFORCE.

receive +20 points for a successful dropoff and lose 1 point for every timestep it takes. There is also a 10 point penalty for illegal pick-up and drop-off actions.

The training of the Boltzmann policy is performed exploiting REINFORCE [22] with Adam [41] (learning rate: 0.008, iterations: 1000).

Besides the average return, shown in section 7.1, we compare the recovered reward functions in terms of distance between the expert’s policy and the learned policy (Kullback-Leibler divergence) estimated with (10). Figure 5 confirms the faster convergence of CR-IRL w.r.t. the true reward and ME-IRL. We observe that LPAL is able to recover a policy that is more similar to the expert’s policy w.r.t. the one produced by BC, however the average return is worse (this is a consequence of the fact that LPAL does not learn well to perform pick-up and drop-off actions that are associated with the largest rewards).

We also test the model-based version of CR-IRL, where the transition model is estimated with (9). In Figure 6 we notice that the usage of model-based ECO-R instead of model-free ECO-R has no relevant impact on CR-IRL and LPAL, while ME-IRL benefits from the model-based features only when the expert is not noisy.

The trace heuristic can be also used when a given set of features (different from ECO-R) is provided. Since the usage of the Hessian makes sense only when all the considered features are a stationary point of the policy gradient we first need to remove the orthogonal projections over the space spanned by the gradient log-policy. We compare (Figure 7) the learning performance of the set of the first 100 Proto-Value Functions (PVF) [25] when linearly combined with ME-IRL, LPAL and our trace heuristic. We notice that LPAL outperforms both the maximum entropy and the Hessian approaches regardless of the exploration level. Trace heuristic is slightly more effective w.r.t. maximum entropy when the expert is not noisy. Overall the ECO-R yields a better performance w.r.t. the PVFs.

The last set of experiments is aimed to evaluate the effect of the number of expert’s trajectories on the average return of the recovered reward functions (Figure 8). The experiment is performed with model-free ECO-R. We notice that CR-IRL is susceptible to the number of expert’s trajectories only when $\epsilon = 0.1$: the expert demonstrated a suboptimal behavior and this is more likely when the trajectories are many. This reflects on the estimation of the reward function that does not penalize suboptimal actions performed by the expert. LPAL shows the expected behavior, improving the average return as the number of trajectories increases. Also BC improves overall with the number of trajectories, more effectively when the expert is not noisy.

D.3 Linear Quadratic Gaussian

The Linear Quadratic Gaussian (LQG) corresponds to the discrete time control problem in a continuous state-action space. The state and reward equations are given by:

$$\mathbf{s}_{t+1} = \mathbf{A}\mathbf{s}_t + \mathbf{B}\mathbf{a}_t + \boldsymbol{\eta}_t \quad r_t = -\mathbf{s}_t^T \mathbf{Q} \mathbf{s}_t - \mathbf{a}_t^T \mathbf{R} \mathbf{a}_t,$$

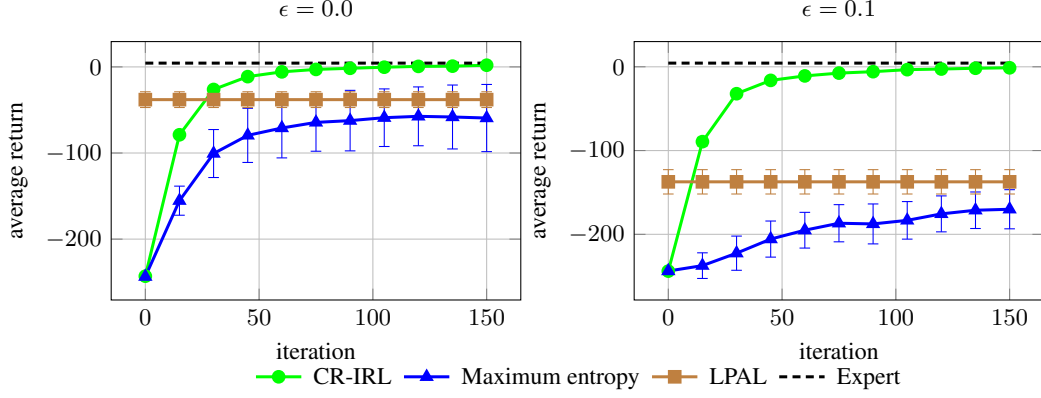


Figure 6: Average return of the Taxi problem as a function of the number of iterations of REINFORCE for reward functions recovered from model-based ECO-Rs.

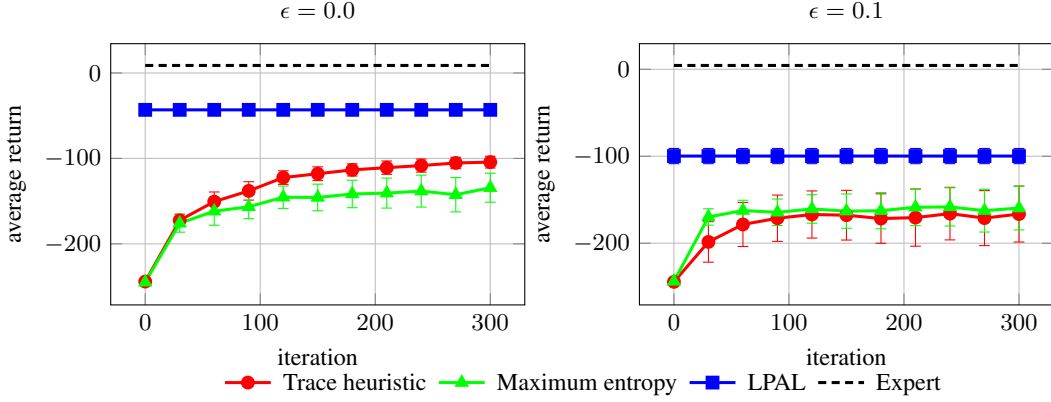


Figure 7: Average return of the Taxi problem as a function of the number of iterations of REINFORCE with the reward function obtained from PVFs.

where \mathbf{A} , \mathbf{B} , \mathbf{Q} and \mathbf{R} are coefficient matrices and $\boldsymbol{\eta}_t$ is a noise process assumed to be a Gaussian white noise $\boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda})$ with uncorrelated components $\boldsymbol{\Lambda} = \lambda^2 \mathbf{I}$. The goal of the agent is to reach as soon as possible the origin since it receives, at each time step, a penalization proportional to the magnitude of the state and the action. The optimal control policy in steady state conditions is the linear controller $\mathbf{a}_t = \mathbf{K}\mathbf{s}_t$ where matrix \mathbf{K} can be found by solving the Riccati equation [36]. The values of the coefficients are reported in Table 1.

Table 1: Coefficients of the LQG experiments. \mathbf{K}^* is the optimal parameter.

\mathbf{A}	\mathbf{B}	\mathbf{Q}	\mathbf{R}	λ^2	\mathbf{K}^*
1.0	1.0	0.9	0.9	0.1	-0.61525

The training is performed exploiting REINFORCE [22] with Adam [41] (learning rate: 0.002, iterations: 600).

We investigate the effect of the number of expert’s demonstrations on the performance of the considered IRL algorithms. Figure 9 reports the convergence parameter value and the average return for a variable number of expert’s trajectories. We notice that CR-IRL is less sensitive w.r.t. to GIRL since, even if few expert’s trajectories are provided, CR-IRL penalizes the regions that the expert did not visit, avoiding the agent to fall into low-reward regions. Differently from the discrete case, in the continuous case the penalization must be “soft”, i.e., the more a non-visited state-action pair is far from the visited ones the more it is penalized.

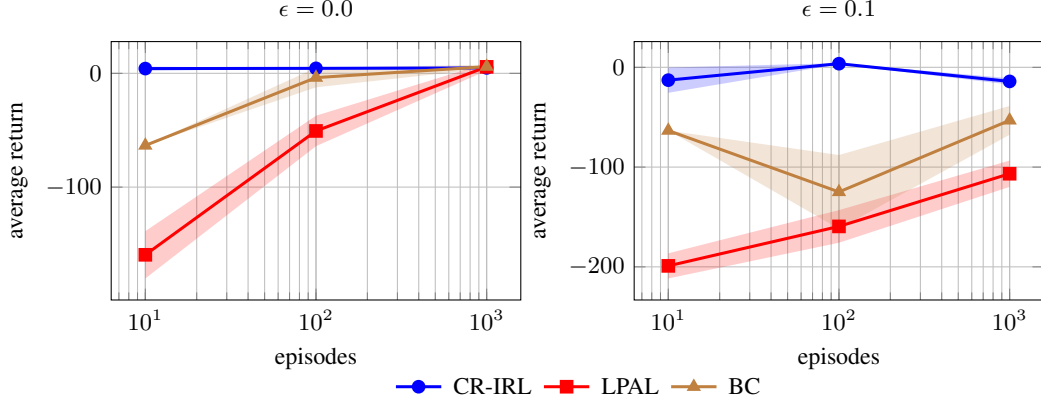


Figure 8: Average return of the Taxi problem as a function of the number of expert's trajectories.

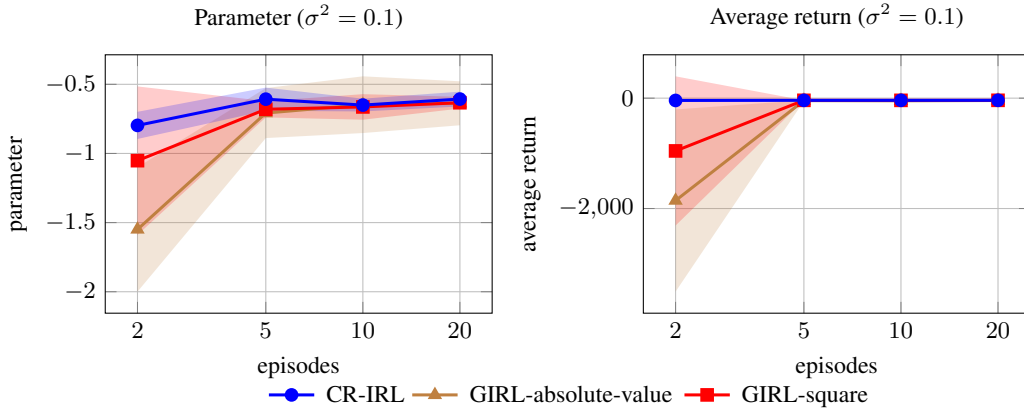


Figure 9: Parameter and average return of LQG as a function of the number of expert's trajectories.

Table 2 and Table 3 show the values of the gradient and the Hessian computed for the different recovered reward functions. GIRL, by construction, yields the reward functions with the smallest gradient in absolute value. On the contrary CR-IRL provides a reward function with larger gradient variance; this is justified by the fact that the Hessian is a large negative number making the policy parameter value a very unstable point for the expected return. Figure 10 shows the average return as a function of the number of iterations for two different values of the policy variance (σ^2).

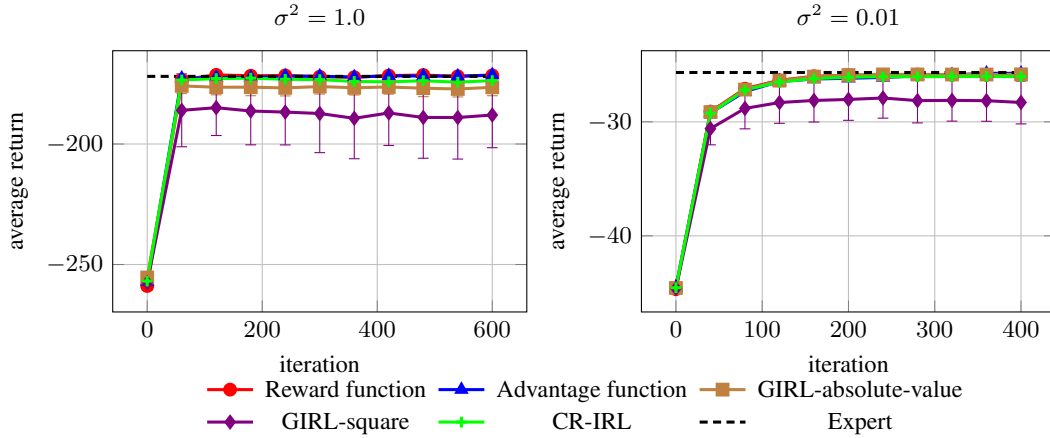


Figure 10: Average return of LQG as a function of the number of iterations of REINFORCE.

Table 2: Gradient values of the recovered reward functions of LQG for different values of σ^2 .

Reward function	$\sigma^2 = 1$	$\sigma^2 = 0.01$
True reward	-0.09203 ± 0.5076	0.03185 ± 0.08893
Advantage function	-0.09356 ± 0.4983	0.9083 ± 1.467
CR-IRL	2.822 ± 10.77	-0.7514 ± 55.45
GIRL-abs-val	$(-2.587 \pm 1.422)\text{e-}15$	$(1.858 \pm 6.068)\text{e-}14$
GIRL-square	$(-1.256 \pm 1.534)\text{e-}14$	$(-2.144 \pm 5.232)\text{e-}14$

Table 3: Hessian values of the recovered reward functions of LQG for different values σ^2 .

Reward function	$\sigma^2 = 1$	$\sigma^2 = 0.01$
True reward	-3.064 ± 7.262	-6.371 ± 5.847
Advantage function	-3.601 ± 6.521	-223.1 ± 114.4
CR-IRL	-1854 ± 147.5	-28702 ± 2586
GIRL-abs-val	3.572 ± 12.33	-17.25 ± 32.84
GIRL-square	-0.3941 ± 8.218	-6.213 ± 4.956

We also investigate the effect of the number of neighbors in the extension of the recovered reward function on unvisited state-action pairs. Figure 11 shows the parameter value and the average return as a function of the number of iterations (the bandwidth of the gaussian kernel are fixed to $\sigma_S = \sigma_A = 4$) and different values of the expert’s policy variance σ^2 . The KNN interpolation allows reconstructing the reward even in regions of the state-action space which are very far from the region visited by the expert, however, in those regions, the recovered reward is not reliable. We compared the performance of the recovered ECO-R as it is and the performance achieved by penalizing the state-action pairs far from the expert exploration region. In the former case (dashed lines) the convergence is slow and displays a high variance, while in the latter case (solid lines) the learning curve converges faster to the optimal parameter. The number of neighbors has no relevant impact, at least for sufficiently small values, when penalization is applied; while when no penalization is applied the effect is more visible. In particular, when the expert is explorative the recovered reward is smoother (many state-action pairs are visited) so 1NN yields the best performance. As the expert becomes less explorative, a large number of neighbors (50NN) performs better since the recovered reward is more discontinuous. We notice that the penalization becomes more beneficial when the expert’s policy is exploitative ($\sigma^2 = 0.01$).

D.4 Car on the Hill

The policy used to collect samples is a mix of the deterministic optimal policy provided by FQI (ExtraTree estimators: 50, metric: mean squared error) and a random policy:

$$\pi_\epsilon^E(a|\mathbf{s}) = (1 - \epsilon)\pi_{FQI}^E(a|\mathbf{s}) + \epsilon\pi_{\text{random}}(a|\mathbf{s}),$$

where ϵ is the exploration level. We collected 20 episodes from policy π_ϵ^E and we used them to fit the parameters \mathbf{w} of a Gaussian policy $\pi_{\mathbf{w}}(a|\mathbf{s}) \sim \mathcal{N}(y_{\mathbf{w}}(\mathbf{s}), \sigma^2)$ where $y_{\mathbf{w}}(\mathbf{s})$ is a radial basis function network:

$$y_{\mathbf{w}}(\mathbf{s}) = \sum_{i=1}^k w_i e^{-\delta \|\mathbf{s} - \mathbf{s}_i\|^2}. \quad (11)$$

The centers are uniformly distributed in the (position-velocity) state-space (Table 4). The parameters \mathbf{w} are learned by minimizing the KL-divergence between $\pi_{\mathbf{w}}$ the expert’s policy π_ϵ^E .

We investigate also the effect of the exploration level ϵ on the performance of the recovered algorithms. Figure 12 shows that CR-IRL outperforms BC for both exploration levels.

In Figure 13 we compare the trajectories of the expert’s policy, the maximum likelihood policy (BC) and the policy computed via FQI from the reward recovered by CR-IRL. We can see that when the expert is deterministic the trajectories are almost overlapping. On the contrary, when the exploration rate ϵ increases we can see that some expert’s trajectories fail to reach the profitable absorbing state.

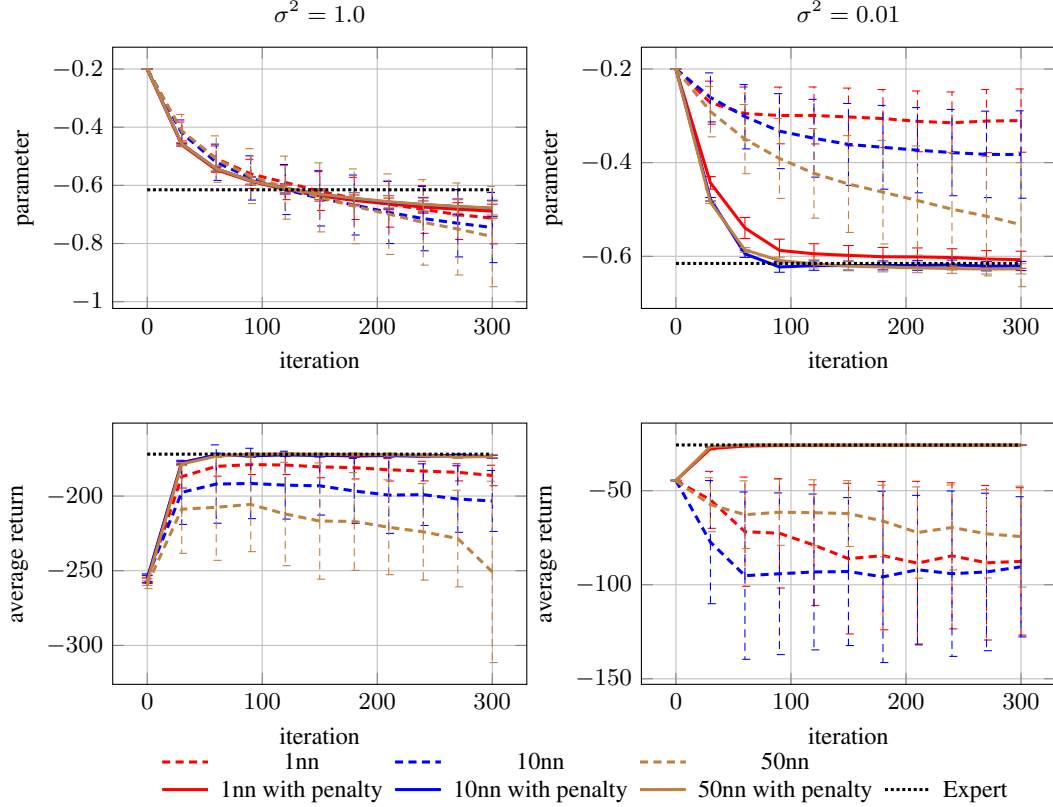


Figure 11: Parameter value and average return of LQG as a function of the number of iterations of REINFORCE with different number of neighbors.

This is a consequence of the fact that a random action is taken with probability ϵ . BC is almost always able to get the final positive reward but not optimally in terms of number of decision epochs required. Finally, CR-IRL, even if trained with a noisy expert, recovers a reward function that induces a policy which is near-optimal, as all trajectories get to the positive final reward in almost the minimum number of steps.

Table 4: Parameters of the Car on the Hill policy.

k	σ^2	δ
40×40	0.01	0.01

E Computational analysis

The main sources of computational complexity in CR-IRL are the SVD, the computation of the Hessian and its eigenvalue computation for each ECO-R. For a $m \times n$ matrix it is well known that SVD has complexity $O(\min\{mn^2, nm^2\})$. We apply SVD twice: the first time it is applied to matrix $\nabla_{\theta} \log \pi_{\theta}^T \mathbf{D}_{\mu}^{\pi_{\theta}}$ of dimension $k \times |\mathcal{D}|$ and the second time to matrix Ψ of dimension at most $|\mathcal{D}| \times |\mathcal{D}|$. Assuming $k \leq |\mathcal{D}|$, the second application dominates yielding a complexity of $O(|\mathcal{D}|^3)$. The computation of the Hessian is linear in the number of ECO-R (at maximum $|\mathcal{D}|$) and in the number of samples ($|\mathcal{D}|$) and quadratic in the number of parameters k , thus the complexity is $O(k^2 |\mathcal{D}|^2)$. Finally the computation of the eigenvalues is cubic in the order of the matrix (k) and has to be done for every ECO-R, yielding a cost of $O(k^3 |\mathcal{D}|)$, which is dominated by $O(k^2 |\mathcal{D}|^2)$ as $k \leq |\mathcal{D}|$. Overall the complexity of CR-IRL is $O(\max\{k^2 |\mathcal{D}|^2, |\mathcal{D}|^3\})$.

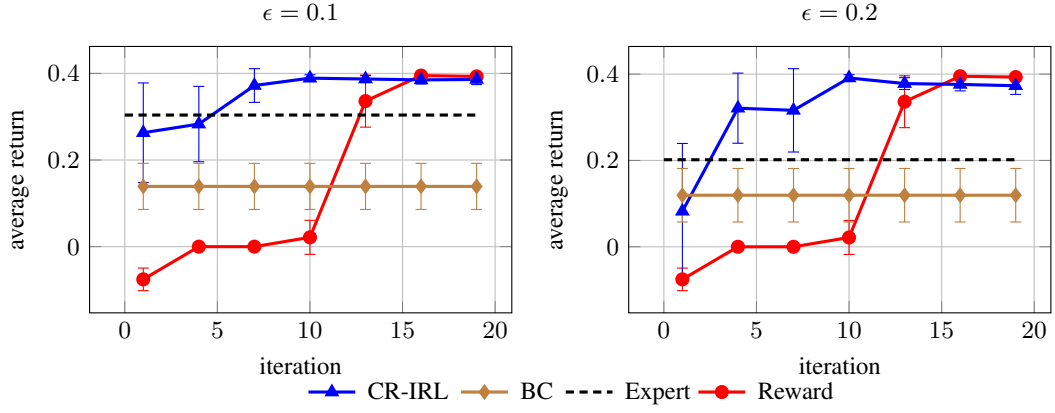


Figure 12: Average return of Car on the Hill problem varying the value of exploration ϵ .

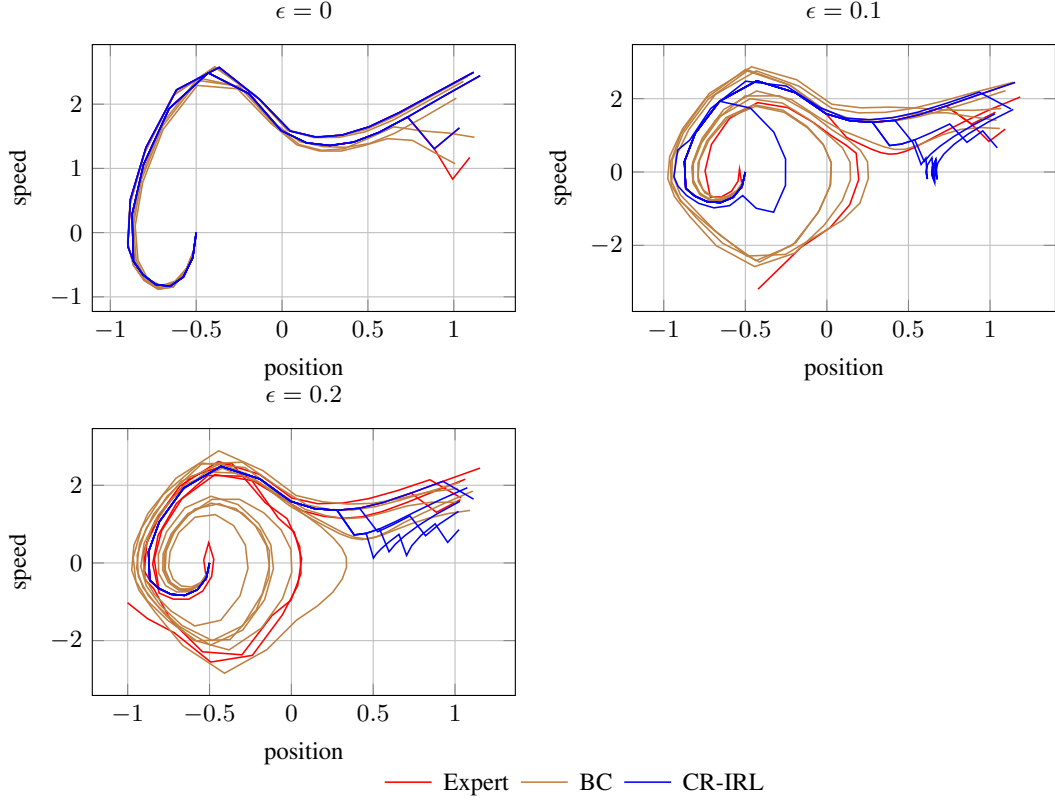


Figure 13: Trajectories of the expert's policy, the ML policy and the policy computed via FQI from the reward recovered by CR-IRL for different values of ϵ .