
Kernel Inverse Optimization

Youyuan Long

Delft Center for Systems and Control
Delft University of Technology
The Netherlands
Y.Long-2@student.tudelft.nl

Pedro Zattoni Scroccaro

Delft Center for Systems and Control
Delft University of Technology
The Netherlands
P.ZattoniScroccaro@tudelft.nl

Tolga Ok

Delft Center for Systems and Control
Delft University of Technology
The Netherlands
T.Ok@tudelft.nl

Peyman Mohajerin Esfahani

Delft Center for Systems and Control
Delft University of Technology
The Netherlands
P.MohajerinEsfahani@tudelft.nl

Abstract

In Inverse Optimization (IO), we hypothesize that experts solve an optimization problem to compute their decisions. If we can reconstruct this optimization problem using the expert’s input and decision data, then the expert’s behavior can be emulated. In this paper, we introduce a novel Kernel Inverse Optimization (KIO) model, which combines IO with kernel methods. To address the scalability issues of the resulting model, we propose a Sequential Selection Optimization (SSO) algorithm to solve the KIO problem. Finally, we demonstrate the strong generalization ability of the KIO and the efficiency of the SSO algorithm for imitation learning applications.

1 Introduction

Inverse Optimization (IO) is distinct from traditional optimization problems where we typically seek the optimal decision variables by providing an objective function and a set of constraints. In contrast, inverse optimization works “in reverse” by determining what the optimization problem is when an optimal solution is given. This approach is based on the assumption that, when individuals address a problem, there exists an underlying optimization problem in their minds, and their decision is an optimal solution to this problem. This is called the Forward Optimization Problem (FOP), which is parametric in the exogenous signal \hat{s} , and has \hat{u} as an optimal solution. Therefore IO involves deducing such an FOP from a dataset of exogenous signal and decision pairs $\{(\hat{s}_i, \hat{u}_i)\}_{i=1,\dots,N}$. After estimating an FOP from data, we can use it for tasks in the absence of an expert, by solving the constructed FOP and executing the corresponding optimal action. The field of IO has garnered widespread attention, giving rise to numerous studies encompassing both theoretical and applied research. Application domains include vehicle routing [1, 2, 3], transportation system modeling [4, 5], portfolio optimization [6, 7, 8], power systems [9, 10, 11], electric vehicle charging problems [12], network design [13], healthcare problems [14] as well as controller design [15]. For a more detailed discussion on different applications of IO, please refer to the recent survey paper [16].

IO can be categorized into classic IO and data-driven IO. In classic IO, only a single signal-decision pair is considered, where the decision is assumed to be the optimal solution of the FOP (i.e., there is no noise), and different classes of FOPs have been studied, such as linear conic problems [17, 18, 19]. However, in real-world applications, there are usually many observations of signal-decision pairs, and due to the presence of noise, it is usually unreasonable to assume noiseless data (i.e., that all observed decisions are optimal w.r.t. a single true data-generating FOP). Additionally, for complex

tasks, the chosen FOP may only approximate the task, not allowing for perfect replication of the observed behavior from the expert. Such cases are called data-driven IO problems. In such scenarios, a loss function is usually used to compute the discrepancy between observed data and the decision generated by the learned FOP. Examples of loss functions include the *2-norm distance loss* [20], *suboptimality loss* [6], *variational inequality loss* [5], *KKT loss* [21] and , *augmented suboptimality loss* [22].

In data-driven IO, the objective function of FOP is typically non-linear with respect to exogenous signal \hat{s} . Hence, classical methods that deduce a FOP based on linear function classes may oversimplify the problem and lead to suboptimal solutions. One effective approach to addressing the expressibility issue in data-driven IO problems is the introduction of kernel methods. These methods have been extensively studied within the context of IO [23, 5] and have shown promising results for scaling IO to address practical problems. The application of kernel methods in IO allows for the exploration of a broader class of optimization problems, thereby enhancing the model’s ability to generalize from observed decisions to unseen situations. Specifically, using a kernelized approach facilitates the embedding of decision data into a richer feature space, enabling the deduction of a FOP that not only fits the training data but also exhibits strong generalization capabilities.

Contributions. We list the contributions of this work as follows:

- (1) **Kernelized IO Formulation** We propose a novel Kernel Inverse Optimization (KIO) model based on suboptimality loss [6]. The proposed approach leverages kernel methods to enable IO models to operate on infinite-dimensional feature spaces, which allows KIO to outperform existing imitation learning (IL) algorithms on complex continuous control tasks in low-data regimes.
- (2) **Sequential Selection Optimization Algorithm** To address the quadratic computational complexity of the proposed KIO model, we introduce the Sequential Selection Optimization (SSO) algorithm that is inspired by coordinate descent style updates. This algorithm selectively optimizes components of the decision variable, greatly enhancing the efficiency and scalability, while provably converging to the same solution of our proposed KIO model.
- (3) **Open Source Code** To foster reproducibility and further research, we provide open source implementation of the proposed KIO model and the SSO algorithm, along with the source code of the experiments conducted in this study.

Notation \mathbb{R}_+^n denotes the space of n -dimensional non-negative vectors. The identity square matrix with dimension n is denoted by I_n . For a symmetric matrix Q , the inequality $Q \succeq 0$ (respectively, $Q \succ 0$) means that Q is positive semi-definite (respectively, positive definite). The trace of a matrix Q is denoted as $\text{Tr}(Q)$. Given a vector $x \in \mathbb{R}^n$, we use the shorthand notation $\|x\|_Q^2 := x^\top Q x$. Symmetric block matrices are described by the upper diagonal elements while the lower diagonal elements are replaced by "*". The Frobenius norm of matrix Q is denoted as $\|Q\|_F$ and the Kronecker product is denoted as \otimes . The notation Q_{ij} represents the element in the i -th row and j -th column of the matrix Q . The Euclidean inner product of x and y is denoted as $x^\top y$.

2 Preliminaries

2.1 Inverse Optimization

In general, to solve a data-driven IO problem we need to design two crucial components: the Forward Optimization Problem (FOP) and the loss function. Specifically, the FOP corresponds to the optimization problem we will “fit” to the observed dataset $\hat{\mathcal{D}} = \{(\hat{s}_i, \hat{u}_i)\}_{i=1}^N$ (in this paper, we use the “hat” notation (e.g., \hat{s}) for the objects dependent on the dataset). That is, our goal is to find a parameter vector $\theta \in \mathbb{R}^p$, such that the FOP

$$\text{FOP}(\theta, \hat{s}) := \min_{u \in \mathcal{U}(\hat{s})} F_\theta(\hat{s}, u), \quad (1)$$

replicates (as close as possible) the data, we do it by minimizing a loss function, akin to classical empirical risk minimization problems. In this work, we focus on quadratic FOPs with linear constraints and continuous decision variables:

$$F_\theta(\hat{s}, u) := u^\top \theta_{uu} u + 2\phi(\hat{w})^\top \theta_{su} u \quad \text{and} \quad \mathcal{U}(\hat{s}) := \left\{ u \in \mathbb{R}^n : \hat{M}u \leq \hat{W} \right\}, \quad (2)$$

where $\theta := (\theta_{uu}, \theta_{su})$, $\theta_{uu} \in \mathbb{R}^{n \times n}$, $\theta_{su} \in \mathbb{R}^{m \times n}$, $\hat{s} := (\hat{w}, \hat{M}, \hat{W})$, $\hat{w} \in \mathbb{R}^q$, $\hat{M} \in \mathbb{R}^{d \times n}$, $\hat{W} \in \mathbb{R}^d$, and $\phi : \mathbb{R}^q \rightarrow \mathbb{R}^m$ is the feature function that maps the feature \hat{w} to a higher-dimensional feature space to enhance the model's capacity.

To learn θ , we solve a regularized loss minimization problem using the Suboptimality Loss [6]

$$\min_{\theta \in \Theta} k\mathcal{R}(\theta) + \frac{1}{N} \sum_{i=1}^N \max_{u_i \in \mathcal{U}(\hat{s}_i)} \{F_\theta(\hat{s}_i, \hat{u}_i) - F_\theta(\hat{s}_i, u_i)\} \quad (3)$$

where $\Theta := \{\theta = (\theta_{uu}, \theta_{su}) : \theta_{uu} \succeq I_n\}$, $\mathcal{R}(\theta) := \|\theta_{uu}\|_F^2 + \|\theta_{su}\|_F^2$, and k is a positive regularization parameter. Notice that the constraint $\theta_{uu} \succeq I_n$ prevents the trivial solution $\theta_{uu} = \theta_{su} = 0$ and guarantees the resulting FOP is a convex optimization problem. Moreover, since F_θ is linear in θ , the optimization program (3) is convex w.r.t. θ , and it can be reformulated from the “minimax” form to a single minimization problem, where the reformulation is based on dualizing the inner maximization problems of (3) and combining the resulting minimization problems.

Proposition 1 (LMI reformulation [15]). *For the hypothesis function and feasible set in (2), the optimization program (3) is equivalent to*

$$\begin{aligned} \min_{\theta, \lambda_i, \gamma_i} \quad & k\mathcal{R}(\theta) + \frac{1}{N} \sum_{i=1}^N \left(F_\theta(\hat{s}_i, \hat{u}_i) + \frac{1}{4} \gamma_i + \hat{W}_i^\top \lambda_i \right) \\ \text{s.t.} \quad & \theta = (\theta_{uu}, \theta_{su}), \theta_{uu} \succeq I_n, \lambda_i \in \mathbb{R}_+^d, \gamma_i \in \mathbb{R} \quad \forall i \leq N \\ & \begin{bmatrix} \theta_{uu} & \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \gamma_i \end{bmatrix} \succeq 0 \quad \forall i \leq N. \end{aligned} \quad (4)$$

2.2 The Kernel Method

The kernel method is a powerful technique used in machine learning and statistics for handling non-linear relationships between variables. The fundamental idea behind the kernel method is to implicitly map input data into a higher-dimensional space without explicitly computing the transformation, thus enabling the algorithms to capture complex patterns and non-linear relationships without heavy computational burden [23]. The kernel method generalizes the hypothesis class of the optimization problem to nonlinear and infinite dimensional class based on *reproducing kernel Hilbert space* (RKHS). Moore-Aronszajn's reproducing kernels theory implies that for any proper kernel function k , there exists a RKHS and a corresponding feature map ϕ with the reproducing properties.

A kernel function $\kappa : \mathcal{H} \times \mathcal{H}$ is defined as

$$\kappa(\hat{w}_i, \hat{w}_j) := \phi(\hat{w}_i)^\top \phi(\hat{w}_j), \quad (5)$$

where $\phi : \mathbb{R}^q \rightarrow \mathcal{H}$, and \mathcal{H} is a dot product space, a.k.a., *feature space* [24]. Define the Gram matrix $K \in \mathbb{R}^{N \times N}$ so that the element in the i 'th row and j 'th column equals $\kappa(\hat{w}_i, \hat{w}_j)$. κ is called a positive-definite kernel if the Gram matrix K is positive-definite [25]. By Mercer's Theorem [26], κ is a proper kernel (having a corresponding feature map function ϕ) if and only if its Gram matrix K is positive-semidefinite. The kernel method has found numerous applications, including Support Vector Machines [27], Kernel Principal Component Analysis [28], and Kernel Linear Discriminant Analysis [29].

3 Kernel Inverse Optimization

In this section, we extend the Inverse Optimization (IO) model proposed by [15] to kernel-methods. We opt to obtain a reformulation that does not depend on the feature mapping ϕ , which is possibly computationally intractable. Hence, we kernelize (4) by substituting the dot product term $\phi(x')^\top \phi(x)$ by the kernel function $\kappa(x', x)$ as a result of the reproducing property [30]. In this work, we study the proposed KIO model with the radial basis function (RBF) kernel, but our results hold for any proper kernel function. The resulting reformulation is shown in Theorem 1.

116 **Theorem 1** (Kernel reformulation). *The Lagrangian dual of the optimization program (4) is*

$$\begin{aligned}
\min_{P, \Lambda_i, \Gamma_i} \quad & \frac{1}{4k} \left\| \left(\sum_{i=1}^N \frac{\hat{u}_i \hat{u}_i^\top}{N} - \Lambda_i \right) - P \right\|_F^2 - \text{Tr}(P) \\
& + \frac{1}{k} \sum_{i=1}^N \sum_{j=1}^N \kappa(\hat{w}_i, \hat{w}_j) \left(\frac{\hat{u}_i}{N} - 2\Gamma_i \right)^\top \left(\frac{\hat{u}_j}{N} - 2\Gamma_j \right) \\
\text{s.t.} \quad & P \succeq 0, \Lambda_i \succeq 0, \Gamma_i \in \mathbb{R}^n \quad \forall i \leq N \\
& \frac{\hat{W}_i}{N} - 2\hat{M}_i \Gamma_i \geq 0 \quad \forall i \leq N \\
& \begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \frac{1}{4N} \end{bmatrix} \succeq 0 \quad \forall i \leq N,
\end{aligned} \tag{6}$$

117 where $\kappa(\hat{w}_i, \hat{w}_j) = \phi(\hat{w}_i)^\top \phi(\hat{w}_j)$ is the kernel function. The primal variables θ_{uu} and θ_{su} can be
118 recovered using

$$\theta_{uu} = \frac{P - \sum_{i=1}^N \left(\Lambda_i - \frac{\hat{u}_i \hat{u}_i^\top}{N} \right)}{2k} \quad \text{and} \quad \theta_{su} = \frac{\sum_{i=1}^N \phi(\hat{w}_i) \left(2\Gamma_i^\top - \frac{\hat{u}_i^\top}{N} \right)}{k}. \tag{7}$$

119 *Proof.* See Appendix A. □

120 Notice that the complexity of the optimization program (6) does not depend on the dimensionality of
121 the feature vector $\phi(\hat{w}_i)$. Consequently, this allows us to use kernels generated even from infinite-
122 dimensional feature spaces, e.g., the Gaussian (a.k.a. radial basis function) kernel $\kappa(\hat{w}_i, \hat{w}_j) =$
123 $\exp(-\gamma \|\hat{w}_i - \hat{w}_j\|_2^2)$. Program (6) is a convex optimization problem, and can be solved using
124 off-the-shelf solvers, such as MOSEK [31]. Once solved, we can recover the optimal primal variables
125 θ_{uu}^* and θ_{su}^* from the optimal dual variables Λ_i^*, Γ_i^* and P^* using (7). Notice that the dimensionality
126 of θ_{su}^* depends on ϕ , thus, it can be an infinite dimensional matrix. However, recall that our ultimate
127 goal is to learn an FOP, and to use it to replicate the behavior observed in the data, and combining (7)
128 into (1), we have that, for the signal $\hat{s}_{\text{new}} = (\hat{w}_{\text{new}}, \hat{M}_{\text{new}}, \hat{W}_{\text{new}})$, the resulting FOP is

$$\min_{\hat{M}_{\text{new}} u \leq \hat{W}_{\text{new}}} \frac{1}{2k} u^\top \left(P^* - \sum_{i=1}^N \left(\Lambda_i^* - \frac{\hat{u}_i \hat{u}_i^\top}{N} \right) \right) u + \frac{2}{k} \left(\sum_{i=1}^N \kappa(\hat{w}_{\text{new}}, \hat{w}_i) \left(2\Gamma_i^* - \frac{\hat{u}_i}{N} \right) \right)^\top u \tag{8}$$

129 which again does not depend on the dimensionality of ϕ , but only depends on the kernel function
130 κ . However, a major difference between solving the IO problem using the kernel reformulation
131 of Theorem 1 and other classical IO approaches (e.g., [6, 15, 22]) is that the resulting FOP (8)
132 explicitly depends on the entire training dataset $\{(\hat{s}_i, \hat{u}_i)\}_{i=1}^N$. These models are sometimes called
133 *nonparametric models* [32], which means that the number of parameters of the model (in our case,
134 P^*, Λ_i^* , and Γ_i^* for all $i \leq N$) depends on the size of the training dataset.

135 In the class of cost function studied in this paper (2), θ_{uu} can be interpreted as a matrix that penalizes
136 the components of the decision vector u . However, for many problems, it turns out that the expert
137 that generates the data equally penalizes each dimension of u , or equivalently, uses $\theta_{uu} = I_n$. This is
138 the case, for instance, in the Gymnasium MuJoCo environments [33], where the reward settings for
139 all tasks have the same penalty on each dimension of the decision vector. Intuitively, this means that
140 the expert trained under such reward settings tries to reduce the magnitude of the decision vector for
141 each dimension uniformly, rather than deliberately decreasing it for a specific dimension. Therefore,
142 assuming $\theta_{uu} = I_n$ as prior knowledge can reduce the model complexity and lead to faster training,
143 without deteriorating the performance of the learned model.

144 **Corollary 1** (Kernel reformulation for $\theta_{uu} = I_n$). *The Lagrangian dual of the optimization program*
 145 *(4) with $\theta_{uu} = I_n$ is*

$$\begin{aligned}
 \min_{\Lambda_i, \Gamma_i: \forall i \in S} \quad & \frac{1}{k} \sum_{i=1}^N \sum_{j=1}^N \kappa(\hat{w}_i, \hat{w}_j) \left(\frac{\hat{u}_i}{N} - 2\Gamma_i \right)^\top \left(\frac{\hat{u}_j}{N} - 2\Gamma_j \right) + \sum_{i=1}^N \text{Tr}(\Lambda_i) \\
 \text{s.t.} \quad & \Lambda_i \succeq 0, \Gamma_i \in \mathbb{R}^n \quad \forall i \in S \\
 & \frac{\hat{W}_i}{N} - 2\hat{M}_i\Gamma_i \geq 0 \quad \forall i \in S \\
 & \begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \frac{1}{4N} \end{bmatrix} \succeq 0 \quad \forall i \in S,
 \end{aligned} \tag{9}$$

146 where $\kappa(\hat{w}_i, \hat{w}_j) = \phi(\hat{w}_i)^\top \phi(\hat{w}_j)$ is the kernel function and $S = \{1, 2, 3, \dots, N\}$. The primal
 147 variable θ_{su} can be recovered using (7).

148 S represents a index set, where all decision variables whose indices belong to S will be optimized,
 149 while decision variables whose indices do not belong to S will retain their original values and be
 150 treated as constants. In Corollary 1, all variables will be optimized, hence S is the set comprising
 151 natural numbers from 1 to N . The concept of index set S introduced in Problem (9) is to align with
 152 the sub-optimization problems based on coordinate descent method outlined in Section 4.

153 In the following section, we will focus on algorithms to solve problem (9). However, all ideas also
 154 apply to the general problem (6).

155 4 Sequential Selection Optimization

156 Solving the kernel IO problem (9) involves optimizing a Semidefinite program (SDP), which may
 157 become prohibitively costly if the number of semidefinite constraints and optimization variables
 158 becomes too large. In our case, the size of the SPD grows linearly with the size of the training dataset
 159 N . For instance, in our experiments, solving (9) with $N = 20000$ using CVXPY [34], requires up to
 160 256GB of memory. Therefore, in this section, we propose a coordinate descent-type algorithm to
 161 find an approximate solution to (9) by iteratively optimizing only a subset of the coordinates at each
 162 iteration of the algorithm, keeping all other coordinates fixed. We define a pair of variables Λ_i and Γ_i
 163 as the i -th coordinate, denoted as $\{\Lambda_i, \Gamma_i\}$. Note that in problem (9), each coordinate is decoupled
 164 in the constraints which is the reason coordinate descent framework is applicable here. We call this
 165 method *Sequential Selection Optimization* (SSO), and we show it in Algorithm 1.

Algorithm 1 Sequential Selection Optimization (SSO)

- 1: Initialize $\{\Lambda_i, \Gamma_i\}_{i=1}^N$
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: Select a batch of p coordinates $S = \{a_i\}_{i=1}^p$, where $a_i \in [1, \dots, N]$
 - 4: Update $\{\Lambda_{a_i}, \Gamma_{a_i}\}_{i=1}^p$ based on (9) with S
 - 5: **end for**
-

166 Here, we explain each step of Algorithm 1: (i) Initialization of the optimization variables. In general,
 167 the variables are usually initialized randomly or set to 0 or 1. These methods are simple, but can not
 168 provide a good initial guess. (ii) Selection of a batch of p coordinates. The most natural approach
 169 to choose p coordinates is to select them cyclically. Alternatively, we can select the coordinates at
 170 random at each iteration (not necessarily with equal probability). Lastly, we can choose coordinates
 171 greedily, choosing the components corresponding to the greatest descent or the ones with the largest
 172 gradient or subgradient at the current iteration [35]. (iii) Solving the KIO subproblem to update the
 173 selected coordinates. The mathematical expression of the subproblem is problem (9) with S , where
 174 S is a set containing the indices of the coordinates that need to be updated. Note that, the coordinates
 175 whose indices $\notin S$ are fixed. Therefore the number of quadratic terms in (9) scales with $|S|^2$ but not
 176 N^2 .

177 Next, we propose two heuristics to accelerate the convergence rate of the SSO algorithm. Namely,
 178 we present a heuristic method for choosing which coordinates (line 3 of Algorithm 1) to optimize and
 179 a warm-up trick to improve the initialization of the optimization variables (line 1 of Algorithm 1).

4.1 Heuristic for choosing coordinates

At each iteration of Algorithm 1, intuitively, the largest improvement will be made by updating the “least optimal” set of p variables. One way to evaluate their degree of suboptimality is to choose the variables with the most significant violation of the Karush-Kuhn-Tucker (KKT) conditions of the primal version of the program (9) (i.e., (4) with $\theta_{uu} = I_n$), which is inspired by the Sequential Minimal Optimization (SMO) method from [36].

Proposition 2. For optimal decision variables of problem (9), the coordinate, $\{\Lambda_i, \Gamma_i\}$, that satisfies

$$\frac{\hat{W}_i}{N} - 2\hat{M}_i\Gamma_i > 0 \quad (10)$$

should also satisfies

$$\text{Tr} \left(\begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \frac{1}{4N} \end{bmatrix} \begin{bmatrix} I_n & 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \|2\theta_{su}^\top \phi(\hat{w}_i)\|_2^2 \end{bmatrix} \right) = 0. \quad (11)$$

Proof. See Appendix B. \square

The Proposition 2 is based on KKT conditions. Based on condition (11), we can define the KKT violation condition

$$\text{kkt_violator}(i) := \left| \text{Tr} \left(\begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \frac{1}{4N} \end{bmatrix} \begin{bmatrix} I_n & 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \|2\theta_{su}^\top \phi(\hat{w}_i)\|_2^2 \end{bmatrix} \right) \right|. \quad (12)$$

Finally, from the violation condition (12), we can establish the following heuristic method to construct the set S in Algorithm 1: given the current values of $\{\Lambda_i, \Gamma_i\}_{i=1}^N$ at iteration t , we choose the p coordinates that satisfy condition (10) with the maximum KKT violation (12). However, in practice, at each iteration, we additionally select some random coordinates to update. This ensures that all coordinates have the chance to be updated, including those that initially do not meet the criteria specified in condition (10), which would otherwise never be updated.

4.2 Warm-up trick for improved initialization

Another component of Algorithm 1 that may have a great practical impact is how the optimization variables $\{\Lambda_i, \Gamma_i\}_{i=1}^N$ are initialized. A poor initial guess (e.g., $\Lambda_i = \Gamma_i = 0$) can lead to slow solver convergence or even result in numerical instability. In this part, we propose a simple warm-up trick that leads to a better initialization of the optimization variables. First, we divide the original dataset $\hat{\mathcal{D}}$ into n sub datasets $\hat{\mathcal{D}}_1, \dots, \hat{\mathcal{D}}_n$ and solve the n small problems (9) with respect to these sub datasets ($N = |\hat{\mathcal{D}}_i|$ and S is the set of indices of all the data in $\hat{\mathcal{D}}_i$). Then we concatenate the optimal solutions of these n solved small problems to form an initial guess. Here we emphasize that even when the original problem (9) is intractable due to a large training dataset (i.e., large N), each subproblem is tractable for a small enough batch size $|\hat{\mathcal{D}}_i|$, and its solutions are still feasible w.r.t. (9).

5 Numerical Experiments

5.1 Performance evaluation

In this evaluation, KIO is utilized in the simplified version (9), incorporating a Gaussian kernel, tested on continuous control datasets from the D4RL benchmark [37] and we use SSO Algorithm 1 to train the model. In each task, the model’s performance is assessed through 100 episodes of testing, and the score¹ for KIO is the average obtained over these 100 episodes. The parentheses following KIO scores represent the amount of data used.

For comparison, four additional agents are selected in this experiment. **IO** is the inverse optimization model without kernel method, introduced in Proposition 1. To demonstrate the effect of the kernel method, both the **KIO** and **IO** models utilize identical datasets. The scores of two behaviour cloning agents **BC(TD3+BC)**[38] and **BC(CQL)**[39] are obtained from two offline reinforcement learning papers. These papers have implemented their respective behavior cloning agents using

¹Regarding the definition of the score for one episode, we recommend readers to consult the official documentation of Gymnasium [33] and the D4RL paper [37].

the corresponding datasets in D4RL, serving as baselines to compare against their proposed offline reinforcement learning algorithms. From the original paper, **BC(TD3+BC)** and **BC(CQL)** are evaluated over 10 seeds and 3 seeds respectively. The **Teacher** is the agent responsible for generating the dataset and serves as the target for our imitation learning.

Task	KIO	IO	BC(TD3+BC)[38]	BC(CQL)[39]	Teacher
Hopper-expert	109.9(5k)	31.8	111.5	109.0	108.5
Hopper-medium	50.2(5k)	20.6	30.0	29.0	44.3
Walker2d-expert	108.5(10k)	0.9	56.0	125.7	107.1
Walker2d-medium	74.6(5k)	0.0	11.4	6.6	62.1
Halfcheetah-expert	84.4(10k)	-1.7	105.2	107.0	88.1
Halfcheetah-medium	39.0(5k)	-3.1	36.6	36.1	40.7

Table 1: Performance of KIO, IO, two Behaviour Cloning (BC) agents and the Teacher agent on MuJoCo tasks from D4RL benchmark, on the normalized return metric. The numbers inside the parentheses represent the amount of data used, and the score for KIO and IO in every task is the average score over 100 episodes.

Table 1 displays the final experimental results, where KIO achieves competitive results in four out of the six tasks. In these six tasks, except for a slightly lower score in the Halfcheetah-expert task compared to the teacher agent, KIO’s scores are either close to or higher than those of the teacher agent. This indicates that the KIO model exhibits strong learning capabilities in complex control tasks. However, without the kernel method, the IO model demonstrates weak learning capabilities: achieving only low scores in the hopper task, and failing to learn anything in the other two more challenging tasks. All hyperparameters used in this experiment for KIO are listed in Appendix C.

5.2 Ablation studies

We perform ablation studies to understand the contribution of each individual component: Heuristic Coordinates Selection, abbreviated as HeuristicSelection (Section 4.1), and Warm-Up Trick (Section 4.2). We present our results in Figure 1 in which we compare the performance of removing each component from SSO (All model hyperparameters remain unchanged as shown in Appendix C).

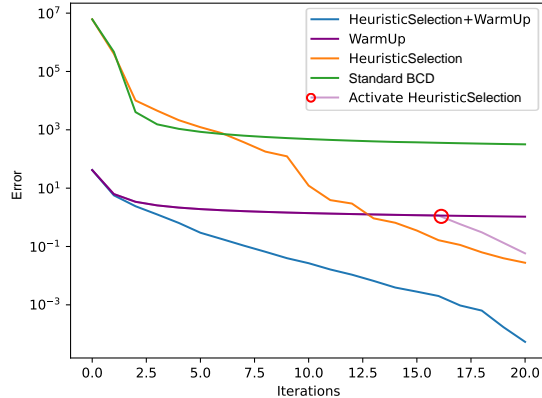


Figure 1: Convergence curves on the MuJoCo Hopper task with the first 5k data points from the D4RL Hopper-expert dataset. The vertical axis represents the difference between the current objective function value and the optimal value. Sequential Selection Optimization (blue) exhibits the fastest convergence rate.

234

We use the Hopper task as the testing task and employ the first 5000 data points from the D4RL Hopper-expert dataset as the training data. The block coordinate for each iteration consists of 2500 coordinates ($|S| = 2500$). When employing the Warm-Up Trick, we partition the data into two equal parts, solving two subproblems (9) each with $|S| = 2500$. Therefore, the computational time required for the Warm-Up Trick is approximately equivalent to the time needed for two iterations of the SSO

algorithm. In Figure 1, we present the results of 20 iterations, with the vertical axis representing the difference between the current objective function value and the optimal value in problem (9).

It is evident that both Heuristic Coordinates Selection and the Warm-Up Trick have significantly accelerated the algorithm. When using the warm-up trick, the initial value of the objective function is markedly reduced. Note that without employing the Warm-Up Trick, the Heuristic (orange) curve decreases to the initial values of the SSO (blue) curve after approximately 10 iterations, whereas using this trick consumes only about the time for 2 iterations. With Heuristic Coordinates Selection, the error curves descend rapidly. It's worth noting the warm-up curve (purple), which becomes nearly flat after a few iterations. However, at the 17th iteration, when we activate the Heuristic Coordinates Selection method, the curve starts to decrease rapidly.

6 Conclusion

We introduce Kernel Inverse Optimization (KIO), an inverse optimization model utilizing kernel methods, and its simple variant along with the theoretical derivations for them. Subsequently, we introduce the Sequential Selection Optimization (SSO) algorithm for training the KIO model, which can address the memory issues by decomposing the original problem into a series of subproblems. Finally, we empirically show that KIO has strong learning capabilities in complex control tasks and SSO algorithm can achieve rapid convergence to the optimal solution within a small number of iterations.

258 A Proof of Theorem 1

259 First, let $P, \tilde{\lambda}_i$ and $\begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \alpha_i \end{bmatrix}$ be the Lagrange multiplier associated with the constraints $\theta_{uu} \succeq I_n$,
 260 $\lambda_i \in \mathbb{R}_+^d$ and $\begin{bmatrix} \theta_{uu} & \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \gamma_i \end{bmatrix} \succeq 0$ respectively, where $P, \Lambda_i \in \mathbb{R}^{n \times n}$, $\Gamma_i \in \mathbb{R}^n$,
 261 $\tilde{\lambda}_i \in \mathbb{R}^d$ and $\gamma_i \in \mathbb{R}$. Then define the Lagrangian function

$$\begin{aligned} L(\theta_{uu}, \theta_{su}, \lambda_i, \gamma_i, P, \tilde{\lambda}_i, \Lambda_i, \Gamma_i, \alpha_i) = & k\|\theta_{uu}\|_F^2 + k\|\theta_{su}\|_F^2 - \text{Tr}(P(\theta_{uu} - I_n)) - \sum_{i=1}^N \tilde{\lambda}_i^\top \lambda_i \\ & + \frac{1}{N} \sum_{i=1}^N \left(\hat{w}_i^\top \theta_{uu} \hat{w}_i + 2\phi(\hat{w}_i)^\top \theta_{su} \hat{w}_i + \frac{1}{4} \gamma_i + \hat{W}_i^\top \lambda_i \right) \\ & - \sum_{i=1}^N \text{Tr} \left(\begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \alpha_i \end{bmatrix} \begin{bmatrix} \theta_{uu} & \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \gamma_i \end{bmatrix} \right), \end{aligned} \quad (13)$$

262 and the Lagrange dual problem

$$\begin{aligned} \max_{P, \tilde{\lambda}_i, \Lambda_i, \Gamma_i, \alpha_i} \quad & \inf_{\theta_{uu}, \theta_{su}, \lambda_i, \gamma_i} L(\theta_{uu}, \theta_{su}, \lambda_i, \gamma_i, P, \tilde{\lambda}_i, \Lambda_i, \Gamma_i, \alpha_i) \\ \text{s.t.} \quad & P \succeq 0, \tilde{\lambda}_i \in \mathbb{R}_+^d, \quad \forall i \leq N \\ & \begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \alpha_i \end{bmatrix} \succeq 0, \quad \forall i \leq N. \end{aligned} \quad (14)$$

263 When the Lagrangian function (13) is at the point of the infimum with respect to $\theta_{uu}, \theta_{su}, \lambda_i, \gamma_i$, we
 264 have

$$\frac{\partial L}{\partial \theta_{uu}} = 2k\theta_{uu} + \frac{1}{N} \sum_{i=1}^N \hat{w}_i \hat{w}_i^\top - P + \sum_{i=1}^N -\Lambda_i = 0 \Rightarrow \theta_{uu} = \frac{P - \left(\sum_{i=1}^N \frac{\hat{w}_i \hat{w}_i^\top}{N} - \Lambda_i \right)}{2k}$$

$$\frac{\partial L}{\partial \theta_{su}} = 2k\theta_{su} + 2 \sum_{i=1}^N \phi(\hat{w}_i) \left(\frac{\hat{w}_i^\top}{N} - 2\Gamma_i^\top \right) = 0 \Rightarrow \theta_{su} = \frac{\sum_{i=1}^N \phi(\hat{w}_i) \left(2\Gamma_i^\top - \frac{\hat{w}_i^\top}{N} \right)}{k}$$

$$\frac{\partial L}{\partial \lambda_i} = \frac{\hat{W}_i}{N} - \tilde{\lambda}_i - 2\hat{M}_i \Gamma_i = 0 \Rightarrow \tilde{\lambda}_i = \frac{\hat{W}_i}{N} - 2\hat{M}_i \Gamma_i$$

$$\frac{\partial L}{\partial \gamma_i} = \frac{1}{4N} - \alpha_i = 0 \Rightarrow \alpha_i = \frac{1}{4N}$$

265 Finally, by substituting the expressions for $\theta_{uu}, \theta_{su}, \tilde{\lambda}_i$ and α_i into the Lagrange dual problem (14)
 266 and simplifying it ends the proof.

270 B Proof of Proposition 2

271 Note that the problem (9) is the dual problem of the problem (4) with $\theta_{uu} = I_n$. Here, we enumerate
 272 the five KKT conditions that will be employed

$$\theta_{su} = \frac{\sum_{i=1}^N \phi(\hat{w}_i) \left(2\Gamma_i^\top - \frac{\hat{w}_i^\top}{N} \right)}{k} \quad (\text{stationarity}), \quad (15a)$$

$$\tilde{\lambda}_i = \frac{\hat{W}_i}{N} - 2\hat{M}_i \Gamma_i, \quad \forall i \leq N \quad (\text{stationarity}), \quad (15b)$$

$$\tilde{\lambda}_i^\top \lambda_i = 0, \quad \forall i \leq N \quad (\text{complementary slackness}), \quad (15c)$$

$$\text{Tr} \left(\begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \frac{1}{4N} \end{bmatrix} \begin{bmatrix} I_n & \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \gamma_i \end{bmatrix} \right) = 0, \quad \forall i \leq N \quad (\text{complementary slackness}), \quad (15d)$$

$$\lambda_i \in \mathbb{R}_+^d, \quad \forall i \leq N \quad (\text{primal feasibility}). \quad (15e)$$

First, we choose coordinate i that satisfy condition (10). Then based on KKT condition (15b), we have

$$\tilde{\lambda}_i > 0. \quad (16)$$

Next, based on conditions (15c) and (15e), one can obtain

$$\lambda_i = 0. \quad (17)$$

Substituting the result (17) into condition (15d) yields

$$\text{Tr} \left(\begin{bmatrix} \Lambda_i & \Gamma_i \\ * & \frac{1}{4N} \end{bmatrix} \begin{bmatrix} I_n & 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \gamma_i \end{bmatrix} \right) = 0. \quad (18)$$

γ_i is the decision variable of problem (4) with $\theta_{uu} = I_n$. Next, let's solve for its expression. By utilizing the Schur complement, we can prove the following two constraints are equivalent

$$\begin{bmatrix} I_n & \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \\ * & \gamma_i \end{bmatrix} \succeq 0 \Leftrightarrow \gamma_i \geq \left\| \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \right\|_2^2.$$

Therefore, problem (4) with $\theta_{uu} = I_n$ can be equivalently expressed as

$$\begin{aligned} \min_{\theta_{su}, \gamma_i, \lambda_i} \quad & k \|\theta_{su}\|_F^2 + \frac{1}{N} \sum_{i=1}^N \left(2\hat{u}_i^\top \theta_{su}^\top \phi(\hat{w}_i) + \frac{1}{4} \gamma_i + \hat{W}_i^\top \lambda_i \right) \\ \text{s.t.} \quad & \lambda_i \in \mathbb{R}_+^d, \gamma_i \in \mathbb{R}, \quad \forall i \leq N \\ & \gamma_i \geq \left\| \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \right\|_2^2, \quad \forall i \leq N, \end{aligned} \quad (19)$$

Here, the variable γ_i is highlighted. It can be easily proven that when γ_i attain its optimal values, the equality in the last constraint should hold: $\gamma_i = \left\| \hat{M}_i^\top \lambda_i + 2\theta_{su}^\top \phi(\hat{w}_i) \right\|_2^2$. Note that $\lambda_i = 0$ (17), then the expression of γ_i is

$$\gamma_i = \left\| 2\theta_{su}^\top \phi(\hat{w}_i) \right\|_2^2. \quad (20)$$

Substituting (20) into (18), one can obtain condition (11).

C Hyperparameters of KIO

Environment	Dataset	k	$scalar$	p
Hopper-expert	0:5k	1e-6	200N	0.5N
Hopper-medium	0:5k	1e-6	200N	0.5N
walker2d-expert	0:5k+10k:15k	1e-5	50N	0.5N
walker2d-medium	15k:20k	1e-5	50N	0.5N
Halfcheetah-expert	325k:330k+345k:350k	5e-6	50N	0.5N
Halfcheetah-medium	10k:15k	5e-6	50N	0.5N

Table 2: KIO Environment Specific Parameters. (N is the size of the dataset and p is the number of coordinates that are updated in one iteration of SSO)

D Extra experiments about SSO

Figure 2 illustrates the convergence performance of the SSO algorithm across six distinct tasks, with the vertical axis representing the error between the current objective function value and the optimal objective function value in problem (9). The SSO algorithm demonstrates fast convergence speed. By the 10th iteration, the errors for all tasks are below 0.1, and by the 20th iteration, the errors have further diminished to approximately $1e-4$ for all tasks. Table 3 presents the optimal objective function values and task scores obtained by the centralized algorithm, Splitting Conic Solver (SCS) [40], and the distributed algorithm, Sequential Selection Optimization (SSO), for solving the problem (9). In this experiment, SCS is employed to directly address this large-scale problem (9) with $|S| = N$. The corresponding solution is evaluated 100 times, and the average is taken as the final score. Meanwhile, SSO addresses the large-scale problem by solving a series of subproblems. After each iteration, we also assess the current solution 100 times, taking the average as the current score. After 20 iterations, there are 20 corresponding scores, and we select the highest score along with the objective function value from the last iteration as the output. It is evident that SSO and SCS yield nearly identical optimal objective function values. However, except for the Halfcheetah-medium task, SSO achieved higher scores across all other tasks.

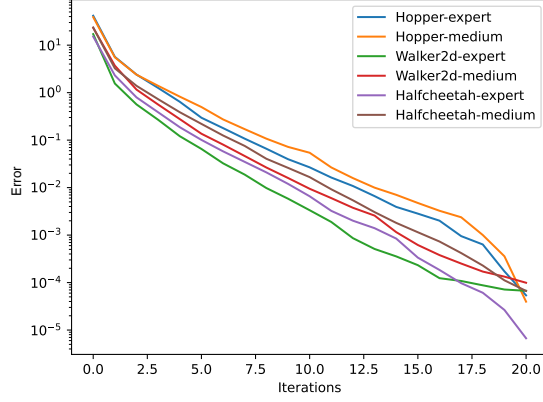


Figure 2: Convergence curves of SSO on the OpenAI gymnasium continuous control tasks. The horizontal axis represents the number of iterations, while the vertical axis denotes the difference between the current objective function value and the optimal objective function value. By the 20th iteration, these problems have converged closely to their optimal values.

Task	SCS		SSO	
	Obj Value	Score	Obj Value	Score
Hopper-expert	185.219	109.9	185.220	110.2
Hopper-medium	218.761	50.2	218.761	51.8
Walker2d-expert	140.121	108.5	140.121	109.2
Walker2d-medium	151.117	74.6	151.117	74.9
Halfcheetah-expert	165.041	84.4	165.041	83.8
Halfcheetah-medium	188.184	39.0	188.184	39.7

Table 3: Final Objective Function Value and Score (Average Return over 100 evaluations) for Splitting Conic Solver (SCS) [40] and Sequential Selection Optimization (SSO). The ultimate Objective Function Values of the two algorithms are nearly identical, yet across the majority of tasks, the SSO exhibits a slightly higher Score compared to the SCS.

References

- [1] Lu Chen, Yuyi Chen, and André Langevin. An inverse optimization approach for a capacitated vehicle routing problem. *European Journal of Operational Research*, 295(3):1087–1098, 2021.
- [2] Mikael Rönnqvist, Gunnar Svenson, Patrik Flisberg, and Lars-Erik Jönsson. Calibrated route finder: Improving the safety, environmental consciousness, and cost effectiveness of truck routing in sweden. *Interfaces*, 47(5):372–395, 2017.
- [3] Pedro Zattoni Scroccaro, Piet van Beek, Peyman Mohajerin Esfahani, and Bilge Atasoy. Inverse optimization for routing problems. *arXiv preprint arXiv:2307.07357*, 2023.
- [4] Michael Patriksson. *The traffic assignment problem: models and methods*. Courier Dover Publications, 2015.
- [5] Dimitris Bertsimas, Vishal Gupta, and Ioannis Ch Paschalidis. Data-driven estimation in equilibrium using inverse optimization. *Mathematical Programming*, 153:595–633, 2015.
- [6] Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the wasserstein metric: performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1-2):115–166, 2018.
- [7] Shi Yu, Haoran Wang, and Chaosheng Dong. Learning risk preferences from investment portfolios using inverse optimization. *Research in International Business and Finance*, 64:101879, 2023.

- [8] Jonathan Yu-Meng Li. Inverse optimization of convex risk functions. *Management Science*, 67(11):7113–7141, 2021.
- [9] John R Birge, Ali Hortaçsu, and J Michael Pavlin. Inverse optimization for the recovery of market structure from market outcomes: An application to the miso electricity market. *Operations Research*, 65(4):837–855, 2017.
- [10] Ricardo Fernández-Blanco, Juan Miguel Morales, and Salvador Pineda. Forecasting the price-response of a pool of buildings via homothetic inverse optimization. *Applied Energy*, 290:116791, 2021.
- [11] Javier Saez-Gallego, Juan M Morales, Marco Zugno, and Henrik Madsen. A data-driven bidding model for a cluster of price-responsive consumers of electricity. *IEEE Transactions on Power Systems*, 31(6):5001–5011, 2016.
- [12] Ricardo Fernández-Blanco, Juan Miguel Morales, Salvador Pineda, and Álvaro Porras. Inverse optimization with kernel regression: Application to the power forecasting and bidding of a fleet of electric vehicles. *Computers & Operations Research*, 134:105405, 2021.
- [13] András Faragó, Áron Szentesi, and Balázs Szviatovszki. Inverse optimization in high-speed networks. *Discrete Applied Mathematics*, 129(1):83–98, 2003.
- [14] Turgay Ayer. Inverse optimization for assessing emerging technologies in breast cancer screening. *Annals of operations research*, 230:57–85, 2015.
- [15] Syed Adnan Akhtar, Arman Sharifi Kolarijani, and Peyman Mohajerin Esfahani. Learning for control: An inverse optimization approach. *IEEE Control Systems Letters*, 6:187–192, 2021.
- [16] Timothy CY Chan, Rafid Mahmood, and Ian Yihang Zhu. Inverse optimization: Theory and applications. *Operations Research*, 2023.
- [17] Ravindra K Ahuja and James B Orlin. Inverse optimization. *Operations research*, 49(5):771–783, 2001.
- [18] Zahed Shahmoradi and Taewoo Lee. Quantile inverse optimization: Improving stability in inverse linear programming. *Operations research*, 70(4):2538–2562, 2022.
- [19] Garud Iyengar and Wanmo Kang. Inverse conic programming with applications. *Operations Research Letters*, 33(3):319–330, 2005.
- [20] Anil Aswani, Zuo-Jun Shen, and Auyon Siddiq. Inverse optimization with noisy data. *Operations Research*, 66(3):870–892, 2018.
- [21] Arezou Keshavarz, Yang Wang, and Stephen Boyd. Imputing a convex objective function. In *2011 IEEE international symposium on intelligent control*, pages 613–619. IEEE, 2011.
- [22] Pedro Zattoni Scroccaro, Bilge Atasoy, and Peyman Mohajerin Esfahani. Learning in inverse optimization: Incenter cost, augmented suboptimality loss, and algorithms. *arXiv preprint arXiv:2305.07730*, 2023.
- [23] Soroosh Shafieezadeh-Abadeh, Daniel Kuhn, and Peyman Mohajerin Esfahani. Regularization via mass transportation. *Journal of Machine Learning Research*, 20(103):1–68, 2019.
- [24] Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. Kernel methods in machine learning. *The Annals of Statistics*, 36(3):1171–1220, 2008.
- [25] Nikhil Das and Michael C Yip. Forward kinematics kernel for improved proxy collision checking. *IEEE Robotics and Automation Letters*, 5(2):2349–2356, 2020.
- [26] James Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909.
- [27] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20:273–297, 1995.

- 369 [28] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component
370 analysis. In *International conference on artificial neural networks*, pages 583–588. Springer,
371 1997.
- 372 [29] Gaston Baudat and Fatiha Anouar. Generalized discriminant analysis using a kernel approach.
373 *Neural computation*, 12(10):2385–2404, 2000.
- 374 [30] Bernhard Scholkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines,*
375 *Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001.
- 376 [31] Mosek ApS. Mosek optimization toolbox for matlab. *User’s Guide and Reference Manual,*
377 *Version*, 4:1, 2019.
- 378 [32] Kevin P. Murphy. *Probabilistic machine learning: an introduction*. MIT press, 2022.
- 379 [33] Mark Towers, Jordan K. Terry, Ariel Kwiatkowski, John U. Balis, Gianluca de Cola, Tristan
380 Deleu, Manuel Goulão, Andreas Kallinteris, Arjun KG, Markus Krimmel, Rodrigo Perez-
381 Vicente, Andrea Pierré, Sander Schulhoff, Jun Jet Tai, Andrew Tan Jin Shen, and Omar G.
382 Younis. Gymnasium, March 2023.
- 383 [34] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for
384 convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- 385 [35] Hao-Jun Michael Shi, Shenyinying Tu, Yangyang Xu, and Wotao Yin. A primer on coordinate
386 descent algorithms. *arXiv preprint arXiv:1610.00040*, 2016.
- 387 [36] John Platt. Sequential minimal optimization: A fast algorithm for training support vector
388 machines. Technical report, Microsoft, 1998.
- 389 [37] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for
390 deep data-driven reinforcement learning, 2020.
- 391 [38] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.
392 *Advances in neural information processing systems*, 34:20132–20145, 2021.
- 393 [39] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for
394 offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–
395 1191, 2020.
- 396 [40] Brendan O’Donoghue, Eric Chu, Neal Parikh, and Stephen Boyd. Conic optimization via
397 operator splitting and homogeneous self-dual embedding. *Journal of Optimization Theory and*
398 *Applications*, 169(3):1042–1068, June 2016.