

A UNIFIED THEORY OF LINEARLY SOLVABLE OPTIMAL CONTROL

KRISHNAMURTHY DVIJOTHAM
EMANUEL TODOROV*

Abstract. We present a unified theory of Linearly Solvable Optimal Control, that is, a class of optimal control problems whose solution reduces to solving a linear equation (for finite state spaces) or a linear integral equation (for continuous state spaces). The framework presented includes all previous work on linearly solvable optimal control as special cases. It includes both standard control problems and risk-sensitive control problems. The degree of risk sensitivity is a parameter of the optimal control problem and can be tuned to achieve the desired trade-off between performance and robustness (to noise/modeling errors). Linearly Solvable Optimal Control problems also possess a number of attractive properties that we explore in this paper. We show that it is possible to construct optimal control laws for new problems by combining the control laws of previously solved optimal control problems. This leads to analytical solutions for a class of non-LQG control problems. Another property is the existence of a path integral representation of the solution to the optimal control problem, which allows us to leverage approximate probabilistic inference techniques to compute optimal control laws. Further, we show that the Inverse Optimal Control problem, that is, the problem of inferring the cost function given trajectories sampled from the optimal control law, can be posed as a convex optimization problem and solved efficiently for problems in this class. We show that the risk sensitive problems can also be viewed as zero sum stochastic games, where the degree of risk averseness grows as the adversary becomes stronger. Under this interpretation, we derive a stochastic maximum principle that characterizes the most likely trajectory of the optimally controlled closed-loop system (that includes both the controller and the adversary).

Key words. Optimal Control, Linear PDEs, Efficient Algorithms, Analytical Solutions

AMS subject classifications. Optimal Control

1. Introduction. ¹ Stochastic optimal control is an elegant framework for specifying and solving control problems. One simply defines a high level cost function encoding the goals of the control task and the framework of stochastic optimal control takes care of the details. Despite this elegance, applications of stochastic optimal control to complex real-world control problems have been rare. This is principally because of the computational complexity of solving optimal control, which scales exponentially with the number of state/control variables. This is a fundamental limitation that Richard Bellman called the curse of dimensionality [3] and we cannot hope to alleviate this problem completely. However, practical control problems often have a lot of structure that can potentially be exploited to obtain an efficient solution. This suggests that a fruitful line of research is to develop a subclass of problems that are more tractable and then try to approximate real control problems with problems of this subclass. One such class of problems is Linear Quadratic Gaussian (LQG) problems where the dynamics are linear, costs are quadratic and noise additive and Gaussian. LQG control is perhaps the most widely used optimal control method in practice. However, the assumption of linear dynamics is limiting. In this paper, we aim to address this limitation and develop a general class of control problems that are computationally tractable and amenable to approximation.

We build on a series of recent interesting results on Linearly Solvable Optimal

*Computer Science and Engineering & Applied Mathematics, University of Washington, Seattle

¹A preliminary conference version of this work appeared in the proceedings of Uncertainty in Artificial Intelligence (UAI) 2011

Control, i.e., optimal control problems for which the Hamilton Jacobi Bellman (HJB) (in continuous time) [20] or the Bellman equation (BE) (in discrete time) [34] can be made **linear**. We develop a general class of linearly solvable optimal control problems (LCs) that include all previous work on linearly solvable optimal control as special cases. The dynamics in such problems can be non-linear (and even non-smooth), the costs can be non-quadratic, and the noise can be non-Gaussian. Yet the problem reduces to solving a linear equation which is a minimized and exponentially-transformed Bellman equation. To be sure, this is not nearly as tractable as an LQG problem, because the linear equation in question is a functional equation characterizing a scalar function (the exponent of the cost-to-go function) over a high-dimensional continuous state space. Nevertheless, the linearity leads to interesting properties that facilitate the development of principled approximation schemes for solving these problems efficiently. LCs include both standard control problems and risk sensitive control problems, which allow one to tune a risk-sensitivity parameter so as to tradeoff performance and robustness. We show that these problems can also be interpreted in a game-theoretic fashion, extending the well known relationships between risk sensitive and game theoretic/robust control [14] [23]. LCs use Rényi divergences [27], a family of divergences between probability measures, as control costs. To the best of our knowledge, this is the first application of Rényi divergences in control. Further, we prove several interesting properties of LCs that can be used to develop efficient approximation methods for solving LCs in high dimensional state spaces.

The paper is organized as follows. In the following section 1.1, we discuss the historical development of the theory of linearly solvable optimal control and relationships between LCs and these previous results. In section 2, we present an intuitive development of LCs, outlining the principal mathematical ideas without going into rigorous proofs. In section 3, we show how all previous work on linearly solvable optimal control (to the best of our knowledge) can be viewed as special cases of LCs, showing that the theory of LCs can be viewed as a unified theory of linearly solvable optimal control. In section 4, we provide a rigorous measure-theoretic development of LCs, providing sufficient conditions for the existence of solutions to LCs for arbitrary compact Borel state spaces.

1.1. Historical perspective. The mathematical trick that makes the HJB or Bellman equation linear is the use of an exponential transformation that maps from cost-to-go functions to so-called “desirability” functions $z = \exp(\kappa v)$. Equivalently, $v = \frac{1}{\kappa} \log(z)$, so the cost-to-go function is a logarithmically transformed desirability function. Logarithmic transformations have a long history, with old roots in physics. Hopf [19] used a logarithmic transformation to show the connection between the Burgers equation and the heat equation. The first application of logarithmic transformations to stochastic optimal control is due to Fleming [13] and Holland [18]. In [15], Fleming and Mitter showed that non-linear filtering corresponds to a stochastic optimal control problem whose HJB equation can be made linear. In [14], these ideas were generalized to risk sensitive stochastic optimal control. A summary of results obtained in this area can be found in [16] (chapter 6). The initial motivation in these works was to use techniques based on optimization and optimal control to solve inference and estimation problems. However, in [20], Kappen realized that this connection between optimal control and estimation can be exploited in the opposite direction, ie, by developing efficient approximation techniques for computing control laws based on algorithms for probabilistic inference. He developed Monte Carlo approximation schemes to approximate the optimal cost-to-go function and compute approximately

optimal control laws efficiently for continuous time optimal control problems [20].

The discrete-time results on linearly solvable optimal control [31] were motivated by the same earlier results but in a more abstract way: we asked, are there classes of linearly-solvable optimal control problems involving arbitrary dynamics? This led to the Linearly Solvable MDP (LMDP) framework. In discrete time, the trick that makes the Bellman equation linear is a well known variational property [10] of the KL divergence and has been used to derive a variational characterization of the Bayesian posterior [25].

The discrete time results for the risk-neutral case were developed first in [31]. This was followed by several papers exploring various properties of the framework and applications. See [5, 6, 24, 40] for continuous time results and [11, 12, 32–36, 41, 42] for discrete time results. Applications to robotic control and control of animated characters have been developed in [8, 29, 30, 37]. We generalized these results to the risk sensitive case in a recent conference paper [12] for finite state spaces. Our contribution here is to present a rigorous development of those results for arbitrary Borel state spaces.

The connection to stochastic games presented in section 4.2.2 is not new: Similar tricks have been used to show the existence of solutions to risk sensitive control problems by relating the solution to a discounted stochastic game [23] [9]. However, the results we develop here are slightly different since they directly relate infinite horizon stochastic games with infinite horizon risk sensitive control, without going through vanishing discount limit of discounted stochastic games. Further, we present a maximum principle characterizing the most likely trajectory of the optimally controlled system (section 4.3.4). This is a new result that does not hold in the general stochastic game setting.

2. Linearly Solvable Optimal Control : An Informal Introduction. In this section, we provide an informal introduction to the ideas contained in this paper that would be accessible to a broader audience. In section 4, we provide a rigorous measure theoretic development of these ideas that is applicable to control problems with arbitrary compact Borel state spaces.

2.1. An alternate view of control. Conventionally, we think of control signals as quantities that modify the system behavior in some pre-specified manner. In this framework it is more convenient to work with a somewhat different notion of control, which is nevertheless largely equivalent to the conventional notion and still allows us to model many problems of practical interest. To motivate this alternative view, consider a control-affine diffusion:

$$d\mathbf{x} = (\mathbf{a}(\mathbf{x}) + \mathbf{B}(\mathbf{x}) \mathbf{u}) dt + \mathbf{C}(\mathbf{x}) d\omega$$

where $\mathbf{x} \in \mathfrak{R}^n$. This equation specifies the infinitesimal change in the state \mathbf{x} , caused by a passive/uncontrolled drift term $\mathbf{a}(\mathbf{x})$, a control input \mathbf{u} scaled by a control gain $\mathbf{B}(\mathbf{x})$, and Brownian motion noise with amplitude $\mathbf{C}(\mathbf{x})$. Subject to this system dynamics, the controller seeks to minimize a cost function of the form

$$\ell(\mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{u}$$

accumulated over time. Suppose that the system is in state \mathbf{x}_t at time t . If time moves forward by h (which is sufficiently small), the probability distribution of the state \mathbf{x}_{t+h} at time $t+h$ is approximately

$$\mathbf{x}_{t+h} \sim \mathbb{P}(\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_t + h(\mathbf{a}(\mathbf{x}_t) + \mathbf{B}(\mathbf{x}_t) \mathbf{u}_t), h\mathbf{C}\mathbf{C}^T)$$

where we have discretized the stochastic process using a time step h . Thus, one way of thinking of the effect of control is that it changes the distribution of the next state from $\mathcal{N}(\mathbf{x}_t + h\mathbf{a}(\mathbf{x}_t), h\mathbf{C}\mathbf{C}^T)$ to $\mathcal{N}(\mathbf{x}_t + h(\mathbf{a}(\mathbf{x}_t) + \mathbf{B}(\mathbf{x}_t)\mathbf{u}_t), H\mathbf{C}\mathbf{C}^T)$. Thus, we can think of the controlled dynamics $\mathcal{N}(\mathbf{x}_t + h(\mathbf{a}(\mathbf{x}_t) + \mathbf{B}(\mathbf{x}_t)\mathbf{u}_t), H\mathbf{C}\mathbf{C}^T)$ as an alternate way to parameterize our control input.

The linearly solvable optimal control problems we develop in this paper will use control signals that are probability measures over the entire state space and directly specify the distribution of the next state given the current state: $\mathbf{x}' \sim u(\mathbf{x})$. This gives the controller a lot of freedom, and in principle, it could pick the controlled dynamics to be a delta distribution at the goal state (in a goal-reaching task). In order to avoid this unrealistic solution, we impose a cost on picking a controlled density: this is defined as some sort of divergence to the natural or uncontrolled or **passive** dynamics of the system. The control cost measures how different the controlled dynamics is from the uncontrolled dynamics, thereby measuring the “control effort” used. We will also impose the restriction that u has the same support as the uncontrolled dynamics, so that any transition that happened under the controlled dynamics could also have happened under the passive dynamics (but perhaps with a much lower or higher probability). This interchangeability of systematic control and random chance, is a key feature of all linearly solvable optimal control problems. Mathematically, treating control signals as probability distribution and imposing a control cost of this kind, allows us to find the minimizing control in the Bellman equation analytically. Once this is done, for the specific costs we use, the Bellman equation can be exponentiated to get a linear equation which we call the linear Bellman equation. This is the essential intuition behind all the linearly solvable optimal control problems we will describe.

2.2. Markov Decision Processes. A Markov Decision Process (MDP) is characterized by specifying state space \mathcal{X} , a control space \mathcal{U} , a transition probability kernel $\mathbb{P}(\mathbf{x}'|\mathbf{x}, \mathbf{u})$, $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, $\mathbf{u} \in \mathcal{U}$, which denotes the probability of reaching \mathbf{x}' in one time step after applying control \mathbf{u} in state \mathbf{x} , and a cost function $\ell(\mathbf{x}, \mathbf{u})$, denoting the cost of applying control \mathbf{u} in state \mathbf{x} . The objective is to design an **optimal feedback policy** $u^* : \mathcal{X} \rightarrow \mathcal{U}$ that minimizes the expected accumulated cost

$$\mathbb{E}_{\mathbf{x}_{t+1} \sim \mathbb{P}(\mathbf{x}_t, u(\mathbf{x}_t))} \left[\sum \ell(\mathbf{x}_t, u(\mathbf{x}_t)) \right].$$

The sum is over time, but can be accumulated in different ways (finite horizon, infinite horizon average cost, first exit, discounted infinite horizon). We will not get to the specifics of the formulation in this section. Section 4 contains a rigorous development with all the details.

2.3. Linearly Solvable MDPs (LMDPs). Consider a Markov Decision Process with a finite state space \mathcal{X} . In linearly solvable MDPs, the controller is allowed to directly pick the distribution of the next state $\mathbf{x}' \sim u$ given the current state \mathbf{x} . However, the choice is penalized for being different from the “passive” or “uncontrolled” dynamics of the system, $P^0(\mathbf{x})$ as $\text{KL}(u \parallel P^0(\mathbf{x}))$. This acts like a control cost, and puts a penalty on the controller for making the controlled dynamics very different from the uncontrolled dynamics. In addition to this, we allow an arbitrary state cost $\ell(\mathbf{x})$. If we plug this into the Bellman equation for the infinite horizon average cost

MDP, we get

$$\begin{aligned}
v(\mathbf{x}) + \lambda &= \min_u \ell(\mathbf{x}) + \text{KL}(u \parallel P^0(\mathbf{x})) + \mathbb{E}_u[v] \\
&= \ell(\mathbf{x}) + \min_u \mathbb{E}_u \left[\log \left(\frac{u}{P^0(\mathbf{x})} \right) + v \right] \\
&= \ell(\mathbf{x}) + \min_u \mathbb{E}_u \left[\log \left(\frac{u}{P^0(\mathbf{x}) \exp(-v)} \right) \right] \\
&= \ell(\mathbf{x}) - \log \left(\mathbb{E}_{P^0(\mathbf{x})} [\exp(-v)] \right) + \min_u \text{KL} \left(u \parallel \frac{P^0(\mathbf{x}) \exp(-v)}{\mathbb{E}_{P^0(\mathbf{x})} [\exp(-v)]} \right) \\
&= \ell(\mathbf{x}) - \log \left(\mathbb{E}_{P^0(\mathbf{x})} [\exp(-v)] \right)
\end{aligned}$$

$$\exp(-\lambda) \exp(-v(\mathbf{x})) = \exp(-\ell(\mathbf{x})) \mathbb{E}_{P^0(\mathbf{x})} [\exp(-v)]$$

where the optimal policy is $u^*(\mathbf{x}) = \frac{P^0(\mathbf{x}) \exp(-v)}{\mathbb{E}_{P^0(\mathbf{x})} [\exp(-v)]}$, since the KL divergence is minimized when the two distributions are equal. Defining $z = \exp(-v)$, we get the linear Bellman equation

$$\exp(-\lambda) z(\mathbf{x}) = \exp(-\ell(\mathbf{x})) \mathbb{E}_{P^0(\mathbf{x})} [z].$$

Since z is inversely related to v , the cost-to-go function, we call it the **desirability function**.

2.4. Generalization to Risk Sensitive Control. Risk sensitive MDPs [23] are generalization of MDPs that take into account the risk seeking or risk averse nature of the decision maker (controller). Instead of minimizing the standard expected accumulated cost criterion, we minimize

$$\frac{1}{\alpha} \log \left(\mathbb{E}_{\mathbf{x}_{t+1} \sim \mathbb{P}(\mathbf{x}_t, u(\mathbf{x}_t))} \left[\exp \left(\alpha \left(\sum \ell(\mathbf{x}_t, u(\mathbf{x}_t)) \right) \right) \right] \right)$$

where α is the degree of risk sensitivity. When $\alpha > 0$, the controller is risk averse, when $\alpha < 0$, the controller is risk seeking and in the limit $\alpha \rightarrow 0$, the above objective reduces to the standard MDP objective (known as the risk-neutral case). Risk sensitivity has been discussed in detail in decision theory literature [39]. For our purposes, we will treat it simply as a means to tune the behavior of a controller to trade-off performance and robustness. In fact, the risk sensitive formulation we will use is slightly different from the standard risk sensitive problem: In our formulation, the state cost will be scaled linearly by α , but the control cost will change nonlinearly (but monotonically) as a function of α . We will present a game theoretic interpretation (section 4.2.2) of our results that further justifies the intuition of using α to trade-off performance and robustness. The risk sensitive extension of Linearly Solvable MDPs is formulated with the following objective:

$$\min \frac{1}{\alpha} \log \left(\mathbb{E}_{\mathbf{x}_{t+1} \sim u(\mathbf{x}_t)} \left[\exp \left(\sum_t \alpha \ell(\mathbf{x}_t) + \mathbb{D}_\alpha (P^0(\mathbf{x}_t) \parallel u(\mathbf{x}_t)) \right) \right] \right)$$

where \mathbb{D}_α is the Rényi divergence, informally defined as

$$\mathbb{D}_\alpha(p \parallel q) = \frac{1}{\alpha - 1} \log \left(\int p(\mathbf{x})^\alpha q(\mathbf{x})^{1-\alpha} d\mathbf{x} \right)$$

where p, q are probability density functions. The Bellman equation for this problem (in the infinite horizon average cost setting) can be written as

$$\begin{aligned} v(\mathbf{x}) + \lambda &= \min_u \ell(\mathbf{x}) + \frac{\mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u)}{\alpha} + \frac{1}{\alpha} \log \left(\mathbb{E}_{\mathbf{x}' \sim u(\mathbf{x})} [\exp(v(\mathbf{x}'))] \right) \\ &= \ell(\mathbf{x}) + \min_u \frac{\mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u)}{\alpha} + \frac{1}{\alpha} \log \left(\mathbb{E}_{\mathbf{x}' \sim u(\mathbf{x})} [\exp(v(\mathbf{x}'))] \right) \\ &= \ell(\mathbf{x}) + \frac{1}{\alpha - 1} \log \left(\mathbb{E}_{\mathbf{x}' \sim P^0(\mathbf{x})} [\exp((\alpha - 1)v(\mathbf{x}'))] \right) \end{aligned}$$

where the minimum is achieved at $u^*(\mathbf{x}) = \frac{P^0(\mathbf{x})(z)^{\frac{1}{\alpha-1}}}{\mathbb{E}_{P^0(\mathbf{x})}[(z)^{\frac{1}{\alpha-1}}]}$ (this will be proved rigorously in section 4). Like before, we define the **desirability function** $z(\mathbf{x}) = \exp((\alpha - 1)v(\mathbf{x}))$ and get a linear bellman equation

$$\exp((\alpha - 1)\lambda) z(\mathbf{x}) = \exp((\alpha - 1)\ell(\mathbf{x})) \mathbb{E}_{P^0(\mathbf{x})} [z].$$

As α increases, the controller becomes increasingly risk averse. As $\alpha \rightarrow 0$, $\frac{\mathbb{D}_\alpha}{\alpha} \rightarrow \text{KL}$, so that we recover the risk neutral LMDP case. We call this general class of linearly solvable optimal control problems LCs.

3. LCs: A Unified View of Linearly Solvable Control.

3.1. Relationship to LMDPs. Here we show that LCs can be seen as an extension of LMDPs. Suppose μ, μ' are probability measures over \mathcal{X} . Then,

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{\text{sgn}(\alpha) \mathbb{D}_\alpha(\mu \parallel \mu')}{\alpha} &= \lim_{\alpha \rightarrow 0} \frac{1}{\alpha(\alpha - 1)} \log \left(\mathbb{E}_\mu \left[\left(\frac{d\mu'}{d\mu} \right)^{1-\alpha} \right] \right) \\ &= \lim_{\alpha \rightarrow 0} \frac{1}{2\alpha - 1} \log \left(\mathbb{E}_\mu \left[\left(\frac{d\mu'}{d\mu} \right)^{1-\alpha} \log \left(\frac{d\mu'}{d\mu} \right) \right] \right) \quad (\text{L'Hopital's rule}) \\ &= \mathbb{E}_\mu \left[\frac{d\mu'}{d\mu} \log \left(\frac{d\mu'}{d\mu} \right) \right] = \mathbb{E}_{\mu'} \left[\log \left(\frac{d\mu'}{d\mu} \right) \right] = \text{KL}(\mu' \parallel \mu). \end{aligned}$$

Thus,

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \Psi_{\mathbf{x}_{t+1} \sim u(\mathbf{x}_t)}^\alpha &\left[\left(\sum_{t=0}^{T-1} \ell_t(\mathbf{x}_t) + \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha(P^0(\mathbf{x}_t) \parallel u(\mathbf{x}_t)) \right) + \ell_f(\mathbf{x}_T) \right] \\ &= \mathbb{E}_{\mathbf{x}_{t+1} \sim u(\mathbf{x}_t)} \left[\sum_{t=0}^{T-1} \ell_t(\mathbf{x}_t) + \text{KL}(u(\mathbf{x}_t) \parallel P^0(\mathbf{x}_t)) \right]. \end{aligned}$$

This is the objective for the Linearly Solvable MDPs (section 2.3) [34]. Taking the limit of the linear Bellman Equation (4.1) as $\alpha \rightarrow 0$ gives us the Linear Bellman Equation for LMDPs. Thus, in the risk-neutral limit $\alpha \rightarrow 0$, we recover LMDPs as a special case of LCs.

3.2. Relationship to Risk Sensitive Path Integral Control. A Linearly Solvable Risk Sensitive Controlled Diffusion (LRD) is defined by specifying an Ito diffusion process of the following kind:

$$d\mathbf{x} = \mathbf{a}(\mathbf{x}) dt + \mathbf{B}(\mathbf{x})(\mathbf{u} dt + \sigma d\omega),$$

a state cost function $\ell(\mathbf{x})$, a noise level σ , a scaling factor θ and a risk factor α . The Control Problem is to design a feedback policy $u : \mathcal{X} \rightarrow \mathcal{U}$ in order to minimize

$$\log \left(\mathbb{E} \exp \left(\alpha \int \left(\ell(\mathbf{x}(t)) + \frac{\theta}{2\sigma^2} u(\mathbf{x}(t))^T u(\mathbf{x}(t)) \right) dt \right) \right)$$

where the expectation is under the stochastic process defined by the diffusion. Consider an h -step Euler discretization of the problem:

$$P^h(\mathbf{x}' | \mathbf{x}, \mathbf{u}) = \mathcal{N}(\mathbf{x} + (a(\mathbf{x}) + B(\mathbf{x}) \mathbf{u})h; \sigma h B(\mathbf{x}) B(\mathbf{x})^T)$$

Let $z^h(\mathbf{x})$ be the z function for the h -step discretization and

$$P^{0h}(\mathbf{x}' | \mathbf{x}) = \mathcal{N}(\mathbf{x} + a(\mathbf{x})h; \sigma h B(\mathbf{x}) B(\mathbf{x})^T).$$

The Bellman equation for the LC with risk factor $\alpha' = \theta\alpha$ is

$$\begin{aligned} z^h(\mathbf{x}) \exp(h(1 - \alpha')(\ell(\mathbf{x}) - \lambda)) &= \mathbb{E}_{P^{0h}(\mathbf{x}' | \mathbf{x})} [z^h(\mathbf{x}')] \\ z^h(\mathbf{x}) \frac{\exp(h(1 - \alpha')(q(\mathbf{x}) - \lambda)) - 1}{h} &= \mathbb{E} \left[\frac{z^h(\mathbf{x}') - z^h(\mathbf{x})}{h} \right] \end{aligned}$$

Taking limits as $h \rightarrow 0$, stochastic calculus indicates that [26] that the RHS becomes the generator of the passive Ito diffusion $d\mathbf{x} = \mathbf{a}(\mathbf{x})dt + \sigma \mathbf{B}(\mathbf{x})d\omega$ applied to $z = \lim_{h \rightarrow 0} z^h$. Note here that we're taking two limits simultaneously the function z^h changes with h and so does the stochastic process. Assuming that z^h converges uniformly to a limiting function at a fast enough rate, this result can be stated rigorously. Assuming the result holds, taking the limit gives us

$$z(\mathbf{x}) (\ell(\mathbf{x}) - \lambda) = \frac{\nabla_{\mathbf{x}} z(\mathbf{x})^T \mathbf{a}(\mathbf{x}) + \frac{\text{tr}(\mathbf{B}(\mathbf{x})\mathbf{B}(\mathbf{x})^T \nabla^2 z(\mathbf{x}))}{2}}{1 - \theta\alpha}$$

Using the formulas for Rényi divergence between Gaussians, we get $\mathbb{D}_\alpha \left(P^{0h}(\cdot | \mathbf{x}) \parallel P^h(\cdot | \mathbf{x}, \mathbf{u}) \right) = \frac{h\alpha' \mathbf{u}^T \mathbf{u}}{2\sigma^2} = \frac{h\theta\alpha \mathbf{u}^T \mathbf{u}}{2\sigma^2}$ so that the overall immediate cost for the LC is $h \ell(\mathbf{x}) + \frac{h\alpha' \mathbf{u}^T \mathbf{u}}{2\sigma^2}$, which is exactly the discrete-time Euler approximation of the LRD cost integrated over the time step h . All these results indicate that we're solving exactly the traditional risk sensitive control problem for the given Ito process as $h \rightarrow 0$. However, we leave a rigorous derivation for future work and conjecture that the limit exists under suitable regularity assumptions, similar to results obtained in [22]. The linear PDE obtained above coincides with that obtained in [4] where the LRD problem was first considered.

3.2.1. Linearly Solvable Controlled Diffusions (LDs). The above result also relates LCs with the initial work on continuous time linearly solvable optimal control problems by Kappen [20]. We call these problems Linearly Solvable Controlled Diffusions (LDs). LDs can be viewed either as continuous time limits of LMDPs or as risk neutral limit of LRDs. They deal with controlled diffusions just like LRDs, but the objective is just the expected accumulated cost:

$$\mathbb{E} \int_0^T \left(\ell(\mathbf{x}(t)) + \frac{\theta}{2\sigma^2} u(\mathbf{x}(t))^T u(\mathbf{x}(t)) \right) dt$$

Since we have already related LCs with LMDPs and LRDs, the relationship to LDs follows.

3.3. A Unified Picture. LCs are the most general class of linearly solvable control problems known, to the best of our knowledge. There are two dimensions along which previous results can be recovered: taking risk neutral limits and taking continuous time limits. As the degree of risk sensitivity decreases ($\alpha \rightarrow 0$), we recover Linearly Solvable MDPs (LMDPs) [34] as a special case of LCs. When we view LCs as arising from the time-discretization of Linearly Solvable Risk Sensitive Controlled Diffusions (LRDs) [4] [30], we recover LCs as a continuous time limit ($dt \rightarrow 0$). Linearly Solvable Controlled Diffusions (LDs) [20] can be recovered either as the continuous time limit of an LMDP, or as the non-adversarial limit ($\alpha \rightarrow 0$) of LRDs. The overall relationships between the various classes of linearly solvable control problems is summarized in the figure below:

$$\begin{array}{ccc} LCs & \xrightarrow{\alpha \rightarrow 0} & LMDPs \\ \downarrow dt \rightarrow 0 & & \downarrow dt \rightarrow 0 \\ LRDs & \xrightarrow{\alpha \rightarrow 0} & LDs \end{array}$$

We have given intuitive justifications of these relationships in the preceding sections. We now present a formal measure-theoretic development of LCs. Rigorous proofs of these relationships, however, are left for future work.

4. LCs: A Formal Development. We now present a rigorous measure-theoretic development of the theory of LCs.

4.1. Background and Notation. DEFINITION 1. $\text{sgn}(\alpha) = 1$ if $\alpha \geq 0$ and -1 otherwise. Let $(\Omega, \mathcal{F}, \mu)$ be a probability space. Then, we have the following definitions:

DEFINITION 2. If μ_1, μ_2 are two measures on (Ω, \mathcal{F}) , then we say that μ_1 is absolutely continuous with respect to μ_2 if

$$\forall A \in \mathcal{F}, \mu_2(A) = 0 \implies \mu_1(A) = 0.$$

and denote this by $\mu_1 \ll \mu_2$. Further, if $\mu_1 \ll \mu_2$ and $\mu_2 \ll \mu_1$, we say that these measures are mutually absolutely continuous.

DEFINITION 3. We use $E_\mu[f]$ to denote the expectation of the function f under probability measure μ . Sometimes, we also write $E_{\mathbf{x}' \sim \mu}[f(\mathbf{x}')]$, which denotes the expectation of the random variable $f(\mathbf{x}')$ where \mathbf{x}' is sampled from μ . We use $E_{\mathbf{x}_{t+1} \sim g(\mathbf{x}_t)}[f(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T)]$ to denote the expectation of f under trajectories sampled from a Markov chain $\mathbf{x}' \sim g(\mathbf{x})$.

DEFINITION 4. Define

$$\text{Pos}\{(\Omega, \mathcal{F}, \mu)\} = \{f : f : \Omega \rightarrow \mathbb{R}^+, E_\mu[f] < \infty\}.$$

For any $f \in \text{Pos}\{(\Omega, \mathcal{F}, \mu)\}$, define a new probability measure on this space as

$$(\mu \otimes f)(A) = \frac{E_\mu[f \mathbb{I}[A]]}{E_\mu[f]} \quad \forall A \in \mathcal{F}.$$

Then it can be shown that $(\mu \otimes f) \ll \mu$ and if $g \in \text{Pos}\{(\Omega, \mathcal{F}, \mu)\}$, $E_\mu[fg] < \infty$, then $E_{\mu \otimes f}[g] = \frac{E_\mu[fg]}{E_\mu[f]}$. Also, the Radon Nikodym derivative $\frac{d\mu \otimes f}{d\mu} = \frac{f}{E_\mu[f]}$ and since f is strictly positive, $\mu \ll \mu \otimes f$ so that $\mu, \mu \otimes f$ are mutually absolutely continuous.

DEFINITION 5. *Define*

$$\Psi_\mu^\theta [f] = \frac{1}{\theta} \log \left(\mathbb{E}_\mu [\exp(\theta f)] \right).$$

It can be shown that $\lim_{\theta \rightarrow 0} \Psi_\mu^\theta [f] = \mathbb{E}_\mu [f]$, so we define $\Psi_0^f [\mu] = \mathbb{E}_\mu [f]$.

DEFINITION 6. *If μ, μ' are probability measures on (Ω, \mathcal{F}) such that $\mu' \ll \mu, \mu \ll \mu'$ then define the Rényi divergence of order $\alpha \in (-\infty, \infty)$ between them to be*

$$\mathbb{D}_\alpha (\mu \parallel \mu') = \begin{cases} \frac{\text{sgn}(\alpha)}{\alpha-1} \log \left(\mathbb{E}_\mu \left[\left(\frac{d\mu'}{d\mu} \right)^{1-\alpha} \right] \right) & \text{if } \alpha \neq 1 \\ \text{KL} (\mu \parallel \mu') = \mathbb{E}_\mu \left[\log \left(\frac{d\mu}{d\mu'} \right) \right] & \text{if } \alpha = 1 \end{cases}$$

where $\frac{d\mu'}{d\mu}$ is the Radon-Nikodym derivative. This definition generalizes the standard definition Rényi divergence (which is restricted to $\alpha \geq 0$ [38]). If $\alpha < 0$, it can be shown that $\mathbb{D}_\alpha (\mu \parallel \mu') = \frac{1}{1-\alpha} \mathbb{D}_{1-\alpha} (\mu' \parallel \mu)$. This observation allows us to generalize the properties of the standard Rényi divergence proved in [38]:

- $\mathbb{D}_\alpha (\mu \parallel \mu')$ is continuous as a function of α on $\{\alpha : 0 \leq \alpha \leq 1 \text{ or } \mathbb{D}_\alpha (\mu \parallel \mu') < \infty\}$.
- $\mathbb{D}_\alpha > 0 \forall \alpha$. If $\alpha \neq 0$, $\mathbb{D}_\alpha (\mu \parallel \mu') = 0$ if and only if $\mu = \mu'$.
- \mathbb{D}_α is decreasing for $\alpha < 0$ and increasing for $\alpha > 0$.

4.2. Linearly Solvable Optimal Control Problems (LCs). DEFINITION 7.

A Linearly Solvable Optimal Control Problem (LC) is defined by specifying:

- a A compact metric space \mathcal{X} equipped with the Borel sigma algebra $\mathcal{F}[\mathcal{X}]$ of the Borel subsets of \mathcal{X} . This is the state space.
- b A passive dynamics $P^0 : \mathcal{X} \rightarrow \mathcal{P}[\mathcal{X}] = \{\mu : (\mathcal{X}, \mathcal{F}[\mathcal{X}], \mu) \text{ is a probability space}\}$.
- c A risk parameter $\alpha \in \mathfrak{R}$.
- d A problem formulation: Infinite Horizon Average Cost (IH), First Exit (FE) or Finite Horizon (FH).
- e For IH, FE, a cost function $\ell : \mathcal{X} \rightarrow \mathfrak{R}^+$ and for FH, a time dependent cost function $\ell_t : \mathcal{X} \rightarrow \mathfrak{R}^+$ for FH.
- f For FH, FE, a terminal cost function $\ell_f : \mathcal{X} \rightarrow \mathfrak{R}^+$.
- g For FH, a time horizon T and for FE, a set of terminal states $\mathcal{T} \in \mathcal{F}[\mathcal{X}]$.

The control problem for LCs can be formulated in 2 ways, all of which lead to the same linear Bellman equation. We present all formulations, since they offer complementary insights.

4.2.1. Risk Sensitive Interpretation. DEFINITION 8. *An optimal risk sensitive solution of an LC is a feedback policy $u_c^* : \mathcal{X} \rightarrow \mathcal{P}[\mathcal{X}]$ that achieves the following extremum*

$$\min_{u: \mathcal{X} \rightarrow \mathcal{P}[\mathcal{X}]} \mathcal{J}(\mathbf{x}, u)$$

Subject to $u(\mathbf{x}) \ll P^0(\mathbf{x}), P^0(\mathbf{x}) \ll u(\mathbf{x}) \quad \forall \mathbf{x}$

at every state $\mathbf{x} \in \mathcal{X}$, where $\mathcal{J}(\mathbf{x}, u)$ depends on the problem formulation as follows:

$$\begin{aligned} \text{FH: } \mathcal{J}(\mathbf{x}, u) &= \Psi_{\mathbf{x}_{t+1} \sim u(\mathbf{x}_t)}^\alpha \left[\left(\sum_{t=0}^{T-1} \ell_t(\mathbf{x}_t) + \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha (P^0(\mathbf{x}_t) \parallel u(\mathbf{x}_t)) \right) + \ell_f(\mathbf{x}_T) \right] \\ \text{FE: } \mathcal{J}(\mathbf{x}, u) &= \Psi_{\substack{\mathbf{x}_{t+1} \sim u(\mathbf{x}_t) \\ T_e = \min\{t: \mathbf{x}_t \in \mathcal{T}\}}}^\alpha \left[\left(\sum_{t=0}^{T_e-1} \ell_t(\mathbf{x}_t) + \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha (P^0(\mathbf{x}_t) \parallel u(\mathbf{x}_t)) \right) + \ell_f(\mathbf{x}_{T_e}) \right] \\ \text{IH: } \mathcal{J}(\mathbf{x}, u) &= \lim_{T \rightarrow \infty} \frac{1}{T} \Psi_{\mathbf{x}_{t+1} \sim u(\mathbf{x}_t)}^\alpha \left[\left(\sum_{t=0}^T \ell(\mathbf{x}_t) + \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha (P^0(\mathbf{x}_t) \parallel u(\mathbf{x}_t)) \right) \right] \end{aligned}$$

where the expectations are under the stochastic dynamics of the system defined by u and $\mathbf{x}_0 = \mathbf{x}$.

NOTE 1. The objective above is very similar to that of Risk-Sensitive MDPs [23]. In Risk-Sensitive MDPs, one chooses a feedback policy $\mathbf{u}_t(\mathbf{x})$ that minimizes the cost

$$\Psi_{\mathbf{x}_{t+1} \sim \mathbb{P}(\mathbf{x}_t, \mathbf{u}_t)}^\alpha \left[\sum_{t=0}^{T-1} \ell_t(\mathbf{x}_t, \mathbf{u}_t) + \ell_f(\mathbf{x}_T) \right].$$

In traditional Risk-Sensitive MDPs, the control \mathbf{u} is generally a continuous valued vector ($\in \mathfrak{R}^n$) or a discrete symbol (as in bang-bang control). In LCs, the control input is an entire probability distribution over the state space and $\mathbb{P}(\mathbf{x}, u) = u$. For LCs, we have $\ell_t(\mathbf{x}, u) = \ell_t(\mathbf{x}) + \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha (P^0(\mathbf{x}) \parallel u)$, so the control cost depends on α . In other words, traditional risk sensitive MDPs minimize linearly scaled exponentiated costs $\mathbb{E} \exp(\sum_t \alpha \ell_t(\mathbf{x}_t, \mathbf{u}_t))$ while LCs minimize

$$\mathbb{E} \exp \left(\sum_t \alpha \ell_t(\mathbf{x}_t) + \text{sgn}(\alpha) \mathbb{D}_\alpha (P^0(\mathbf{x}_t) \parallel u(\mathbf{x}_t)) \right),$$

so that the control costs are nonlinearly scaled while the state costs are linearly scaled. However, the nonlinear scaling is still monotonic in α and always has the same sign as α , which makes the resulting behavior similar to traditional risk sensitive control. We also provide a game theoretic interpretation in the next section that further justifies this non-standard formulation of risk sensitivity. Note that the traditional approach of linearly scaling both the state and control costs does not lead to a linearly solvable optimal control problem.

THEOREM 1. If the following linear BE can be solved for a measurable function $z : \mathcal{X} \rightarrow (\epsilon, 1)$ for some $\epsilon > 0$,

$$\begin{aligned} \text{FH: } z_t(\mathbf{x}) &= \exp((\alpha - 1) \ell_t(\mathbf{x})) \mathbb{E}_{P^0(\mathbf{x})} [z_{t+1}], t = 0, 1, \dots, T-1, \forall \mathbf{x} \\ z_T(\mathbf{x}) &= \exp((\alpha - 1) \ell_f(\mathbf{x})) \forall \mathbf{x} \\ \text{FE: } z(\mathbf{x}) &= \exp((\alpha - 1) \ell(\mathbf{x})) \mathbb{E}_{P^0(\mathbf{x})} [z] \quad \forall \mathbf{x} \notin \mathcal{T} \\ z(\mathbf{x}) &= \exp((\alpha - 1) \ell_f(\mathbf{x})) \forall \mathbf{x} \in \mathcal{T} \\ \text{IH: } z(\mathbf{x}) &= \lambda \exp((\alpha - 1) \ell(\mathbf{x})) \mathbb{E}_{P^0(\mathbf{x})} [z] \quad \forall \mathbf{x} \end{aligned} \quad (4.1)$$

then the LC has an optimal feedback solution given by $u^*(\mathbf{x}; t) = P^0(\mathbf{x}) \otimes (z_{t+1})^{\frac{1}{1-\alpha}}$ for FH and $u^*(\mathbf{x}) = P^0(\mathbf{x}) \otimes (z)^{\frac{1}{1-\alpha}}$ for other formulations. We call z the desirability function.

Proof. We present the proof for IH. If the BE for risk sensitive MDPs ([7] equation (2.4)) has a bounded solution, then it characterizes the optimal feedback policy (see verification theorems in [7] [9]). We will show that if the linear BE has a solution $z : \mathcal{X} \rightarrow (\epsilon, 1)$ for some $\epsilon > 0$, then the risk sensitive BE has a solution and its solution can be derived from z .

Suppose the linear BE has a solution $z : \mathcal{X} \rightarrow (\epsilon, 1)$ for some $\epsilon > 0$. Consider the function $v(\mathbf{x}) = \frac{1}{\alpha-1} \log(z(\mathbf{x}))$. Taking logarithms of the linear BE and dividing by $\alpha - 1$ gives:

$$v(\mathbf{x}) = \ell(\mathbf{x}) + \frac{1}{\alpha-1} \log \left(\mathbb{E}_{P^0(\mathbf{x})} [\exp((\alpha-1)v)] \right) \quad (4.2)$$

and v is measurable and bounded in the interval with endpoints at $\frac{\log(\epsilon)}{\alpha-1}$ and 0. By lemma 2, we know that

$$\text{sgn}(\alpha) \mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u) + \Psi_u[\alpha v] \geq -\Psi_{P^0(\mathbf{x})}^{\frac{1-\alpha}{\alpha}}[-\alpha v] \quad \forall u \in \mathcal{P}[(\mathcal{X}, \mathcal{F}[\mathcal{X}])].$$

if $\alpha > 0$ and the inequality is reversed if $\alpha < 0$. Since v is bounded, the probability measure $u^* = P^0(\mathbf{x}) \otimes \exp(-v)$ exists and by definition, satisfies $u^* \ll P^0(\mathbf{x}), P^0(\mathbf{x}) \ll u^*$. We have

$$\frac{du^*}{dP^0(\mathbf{x})} = \frac{\exp(-v)}{\mathbb{E}_{P^0(\mathbf{x})}[\exp(-v)]},$$

$$\mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u^*) = \frac{\text{sgn}(\alpha)}{\alpha-1} \log \left(\left(\mathbb{E}_{P^0(\mathbf{x})} [\exp((\alpha-1)v)] \right)^\alpha \right),$$

$$\Psi_{u^*}[\alpha v] = \log \left(\mathbb{E}_{P^0(\mathbf{x})} [\exp((\alpha-1)v)] \right).$$

$$\text{sgn}(\alpha) \mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u^*) + \Psi_{P^0(\mathbf{x})}[(\alpha-1)v] = \frac{\alpha}{\alpha-1} \Psi_{P^0(\mathbf{x})}[(\alpha-1)v] = -\Psi_{P^0(\mathbf{x})}^{\frac{(1-\alpha)}{\alpha}}[v].$$

Thus, we have

$$-\frac{1}{\alpha} \Psi_{P^0(\mathbf{x})}^{\frac{(1-\alpha)}{\alpha}}[-\alpha v] = \min_{u \ll P^0(\mathbf{x}), P^0(\mathbf{x}) \ll u} \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u) + \frac{1}{\alpha} \Psi_u[\alpha v].$$

with the minimum attained at $u = u^*$. The LHS is equal to $\Psi_{P^0(\mathbf{x})}^{\alpha-1}[v]$. Thus, (4.2) can be rewritten as

$$v(\mathbf{x}) = \ell(\mathbf{x}) + \min_{u \ll P^0(\mathbf{x}), P^0(\mathbf{x}) \ll u} \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u) + \Psi_u^\alpha[v]$$

or equivalently as

$$v(\mathbf{x}) = \min_{u \ll P^0(\mathbf{x}), P^0(\mathbf{x}) \ll u} \ell(\mathbf{x}) + \frac{\text{sgn}(\alpha)}{\alpha} \mathbb{D}_\alpha(P^0(\mathbf{x}) \parallel u) + \Psi_u^\alpha[v].$$

This is precisely the Risk Sensitive Bellman Equation ([7] equation (2.4)) with the LC cost plugged in. Thus, the Risk Sensitive Bellman Equation has a bounded

solution and hence the LC has an optimal policy given by $u^* = P^0(\mathbf{x}) \otimes \exp(-v) = P^0(\mathbf{x}) \otimes (z)^{\frac{1}{1-\alpha}}$. \square

NOTE 2. *The solution to the Linear BE for FH problems always exists and is bounded away from 0 as long as the cost $\ell_t(\mathbf{x})$ is bounded. This can be easily proven constructively, by starting at time T and using dynamic programming backward in time.*

THEOREM 2. *Consider an IH LC. Suppose that the passive dynamics is weaker Feller, that is, for every continuous and bounded function f , the function $g(\mathbf{x}) = E_{P^0(\mathbf{x})}[f]$ is also continuous and bounded. Suppose also that the state cost $\ell(\mathbf{x})$ is continuous. Then, the linear BE has a solution $z : \mathcal{X} \rightarrow (\epsilon, 1)$ for some $\epsilon > 0$.*

Proof. Consider the cone K_+ of positive continuous measurable functions on \mathcal{X} . By the assumptions, the operator defined by $\mathcal{G}[f](\mathbf{x}) = \exp((\alpha - 1)\ell(\mathbf{x}))E_{P^0(\mathbf{x})}[f]$ is weakly Feller and leaves the cone K_+ invariant. Also, it is strongly positive in the sense that for any $f \in \text{int}(K_+)$, $\mathcal{G}[f](\mathbf{x}) \in \text{int}(K_+)$. Then, by the Krein-Rutman theorem ([21], chapter 11), we can conclude that the linear operator \mathcal{G} has a unique eigenfunction $z \in \text{int}(K_+)$ corresponding to the leading eigenvalue $\lambda \in \mathbb{R}^+$. Since z is a continuous function on a compact domain, it must attain its bounds. Also, $z \in \text{int}(K_+)$. The above facts imply that z is bounded below by some number $\eta > 0$. Since any scaled version of z is also an eigenfunction, we can scale z by $1/(M + 1)$ (where M is the upper bound on z) to get a function $\tilde{z} : \mathcal{X} \rightarrow [\frac{\eta}{M+1}, 1)$. Taking $\epsilon = \frac{\eta}{2*(M+1)}$ gives us the result. \square

THEOREM 3. *Consider an FE LC. Suppose that the passive dynamics is such that $\exists T_f < \infty$ such that $P(\mathbf{x}_{T_f} \in \mathcal{T} | \mathbf{x}_0) > \delta > 0 \forall \mathbf{x}_0 \in \mathcal{X}$ under sampling from the passive dynamics $\mathbf{x}_{t+1} \sim P^0(\mathbf{x}_t)$. Further, suppose that $\alpha < 1$. Then there exists a solution $z : (\epsilon, 1)$ satisfying the linear BE for some $\epsilon > 0$.*

Proof. Consider the cone K_+ of positive continuous measurable functions on \mathcal{X} . By the assumptions, the operator defined by $\mathcal{G}[f](\mathbf{x}) = \exp((\alpha - 1)\ell(\mathbf{x}))E_{P^0(\mathbf{x})}[f - f\mathbb{I}[\mathcal{T}]]$ leaves the cone K_+ invariant. Also, the cone K_+ is solid and normal. Taking f to be the constant function 1, we have $\mathcal{G}[f] \leq \left(\sup_{\mathbf{x} \in \overline{\mathcal{X} \setminus \mathcal{T}}} \exp((\alpha - 1)\ell(\mathbf{x}))\right) f$. Since $\overline{\mathcal{X} \setminus \mathcal{T}}$ is compact and $\ell(\mathbf{x})$ is continuous, the supremum is attained and is finite, say γ . Further, since $\ell > 0, \alpha < 1$, $\exp((\alpha - 1)\ell(\mathbf{x})) < 1$ so that $\gamma < 1$. Thus, there exist $f \in \text{int}(K_+)$ such that $\mathcal{G}[f] \leq \gamma f$, $0 < \gamma < 1$. Hence, the spectral radius of the operator \mathcal{G} is smaller than 1. Hence $(I - \mathcal{G})$ is positively invertible. Define the function $g : \mathcal{X} \setminus \mathcal{T} \mapsto \mathbb{R}$ as $g(\mathbf{x}) = E_{P^0(\mathbf{x})}[\mathbb{I}[\mathcal{T}] \exp((\alpha - 1)\ell_f(\mathbf{x}))]$. The function defined by $z = (I - \mathcal{G})^{-1}g$ where $\forall \mathbf{x} \in \mathcal{X} \setminus \mathcal{T}$ and $z(\mathbf{x}) = \exp((\alpha - 1)\ell_f(\mathbf{x})) \forall \mathbf{x} \in \mathcal{T}$ is the unique solution to the linear BE. Also, from theorem 6, z has a path integral representation that shows that under the assumptions of the theorem, z is bounded above and is bounded below by some $\epsilon > 0$. \square

4.2.2. Game Theoretic Interpretation. DEFINITION 9. *In the game theoretic interpretation, an LC has 2 players, the adversary and the controller. It proceeds as follows:*

- *The system is in state \mathbf{x} at time t .*
- *The adversary and controller pick distributions $u_c, u_a \in \mathcal{P}[\mathcal{X}]$ such that $u_c \ll P^0(\mathbf{x}) \ll u_c, u_a \ll P^0(\mathbf{x}) \ll u_a$.*
- *The controller incurs cost $\ell(\mathbf{x}, u_c, u_a) = \ell(\mathbf{x}) - \mathbb{D}_{1/\alpha}(P^0(\mathbf{x}) \| u_a) + \text{KL}(u_c \| u_a)$ while the adversary incurs cost $-\ell(\mathbf{x}, u_c, u_a)$.*
- *The system transitions into a state \mathbf{x}' sampled from the distribution u_c and time advances to $t + 1$.*

NOTE 3. An LC game is a special case of a zero-sum stochastic game [2]. The passive dynamics is meant to encode the autonomous behavior of the system in the absence of controls. The cost function consists of three terms. The first is a state cost $\ell(\mathbf{x})$ encoding desirable states the controller would like to be in. The second term is a control cost for the adversary $-\mathbb{D}_{1/\alpha}(P^0(s) \parallel u_a)$ which penalizes how much the adversary changes the uncontrolled dynamics (the negative sign is because the adversary is maximizing the cost). This cost decreases monotonically with α , so the adversary becomes more powerful as α increases. The third term is a control cost for the controller, $\text{KL}(u_c \parallel u_a)$ which measures the additional change in the system dynamics due to the controller. Thus, as α increases, we expect the adversary to dominate and the controller's strategy to change so as to counter a stronger adversary. The final controlled dynamics is completely determined by the controller's choice, but the controller has to pay a cost for deviating from the choice of the adversary and the adversary has to pay a cost for deviating from the passive dynamics.

DEFINITION 10. A saddle point equilibrium of an LC is a pair of feedback policies $u_c^*, u_a^* : \mathcal{X} \rightarrow \mathcal{P}[\mathcal{U}]$ that achieve the following extremum

$$\min_{u_c} \max_{u_a} \mathcal{J}(\mathbf{x}, u_c, u_a) = \max_{u_a} \min_{u_c} \mathcal{J}(\mathbf{x}, u_c, u_a)$$

at every state $\mathbf{x} \in \mathcal{X}$, where $\mathcal{J}(\mathbf{x}, u_c, u_a)$ depends on the problem formulation:

$$\text{Finite Horizon (FH): } \mathcal{J}(\mathbf{x}, u_c, u_a) = \mathbb{E}_{u_c, u_a, \mathbf{x}_0 = \mathbf{x}} \left[\sum_{t=0}^{T-1} \ell(\mathbf{x}_t, u_c(\mathbf{x}_t), u_a(\mathbf{x}_t)) + \ell_f(\mathbf{x}_T) \right]$$

$$\text{First Exit (FE): } \mathcal{J}(\mathbf{x}, u_c, u_a) = \mathbb{E}_{u_c, u_a, \mathbf{x}_0 = \mathbf{x}} \left[\sum_{t=0}^{T-1} \ell(\mathbf{x}_t, u_c(\mathbf{x}_t), u_a(\mathbf{x}_t)) + \ell_f(\mathbf{x}_T) \right]$$

$$T = \min\{t : s_t \in \mathcal{T}\} - 1$$

$$\text{Infinite Horizon (IH): } \mathcal{J}(\mathbf{x}, u_c, u_a) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{u_c, u_a, \mathbf{x}_0 = \mathbf{x}} \left[\sum_{t=0}^T \ell(\mathbf{x}_t, u_c(\mathbf{x}_t), u_a(\mathbf{x}_t)) \right]$$

where the expectations are under the stochastic dynamics of the system defined by P^0, u_c, u_a .

THEOREM 4. If $\alpha <$, the LC always has a saddle point equilibrium provided (4.1) has a continuous solution $z : \mathcal{X} \rightarrow (\epsilon, 1)$ for some $\epsilon > 0$. The saddle point equilibrium is given by

$$u_a^*(\mathbf{x}; t) = P^0(\mathbf{x}) \otimes (z_{t+1})^{\frac{\alpha}{\alpha-1}}, u_c^*(\mathbf{x}; t) = u_a^*(\mathbf{x}; t) \otimes (z_{t+1})^{\frac{1}{1-\alpha}}$$

for FH and

$$u_a^*(\mathbf{x}) = P^0(\mathbf{x}) \otimes (z)^{\frac{\alpha}{\alpha-1}}, u_c^*(\mathbf{x}) = u_a^*(\mathbf{x}) \otimes (z)^{\frac{1}{1-\alpha}}$$

for IH, FE.

Proof. We do the proof for the IH case, the proof for the other cases is similar. The verification theorem in [17] (theorem 3.4) says that if v is a continuous bounded solution to the Shapeley equation [28]

$$\begin{aligned} \lambda + v(\mathbf{x}) &= \min_{u_c \ll P^0(\mathbf{x}) \ll u_c} \max_{u_a \ll P^0(\mathbf{x}) \ll u_a} \left[\ell(\mathbf{x}) - \mathbb{D}_{1/\alpha}(P^0(\mathbf{x}) \parallel u_a) + \text{KL}(u_c \parallel u_a) + \mathbb{E}_{u_c}[v] \right] \\ &= \max_{u_a \ll P^0(\mathbf{x}) \ll u_a} \min_{u_c \ll P^0(\mathbf{x}) \ll u_c} \left[\ell(\mathbf{x}) - \mathbb{D}_{1/\alpha}(P^0(\mathbf{x}) \parallel u_a) + \text{KL}(u_c \parallel u_a) + \mathbb{E}_{u_c}[v] \right] \end{aligned}$$

then the saddle point equilibrium is attained for the feedback policies setting u_c, u_a to the values attaining the extremum in the Shapely equation above. We are given that the linear BE (4.1) has a solution, say z . Consider the function $v = \frac{1}{\alpha-1} \log(z)$. This is bounded and continuous function since z is continuous, bounded above and below away from 0. Using theorem 9, we know that

$$\begin{aligned} & \max_{u_a \ll P^0(\mathbf{x}) \ll u_a} \min_{u_c \ll P^0(\mathbf{x}) \ll u_c} \left[\ell(\mathbf{x}) - \mathbb{D}_{1/\alpha}(P^0(\mathbf{x}) \parallel u_a) + \text{KL}(u_c \parallel u_a) + \mathbb{E}_{u_c}[v] \right] \\ &= \ell(\mathbf{x}) - \Psi_{P^0(\mathbf{x})}^{\frac{1}{\alpha}-1}[-v] = \ell(\mathbf{x}) + \Psi_{P^0(\mathbf{x})}^{1-\alpha}[v] \end{aligned}$$

and the same result holds if the min, max are interchanged. Since z satisfies the linear BE, we know that v satisfies

$$v(\mathbf{x}) = \ell(\mathbf{x}) + \frac{1}{\alpha-1} \log \left(\mathbb{E}_{P^0(\mathbf{x})} [\exp((\alpha-1)v)] \right) = \ell(\mathbf{x}) + \Psi_{P^0(\mathbf{x})}^{1-\alpha}[v].$$

Combining the results above, we have that v is a bounded continuous solution to the Shapely equation and hence by the verification theorem (theorem 3.4 in [17]), we have the result. \square

4.3. LCs: Properties and Efficient Algorithms. For control problems with discrete (finite) state spaces, the linearity implies that the optimal control problem reduces to solving a linear equation (in the first exit setting), a linear eigenvalue problem (in the infinite horizon average cost setting) and a set of matrix-vector multiplications (in the finite horizon setting). Also, for most control problems of practical interest, the linear systems involved are very sparse and this can be exploited further to get very efficient algorithms. We can easily solve problems with millions of discrete states in time of the order of minutes. For continuous state spaces, the linearity makes the Bellman equation amenable to cost-to-go function approximation techniques. It also leads to other interesting properties that we describe in the following subsections.

4.3.1. Compositionality of Optimal Control Laws. **THEOREM 5.** *Consider a set of FH or FE LC problems with the different terminal costs $\ell_f^{(i)}(\mathbf{x}), i = 1, 2, \dots, k$ but the same running cost $\ell(\mathbf{x})$ (or $\ell_t(\mathbf{x})$), passive dynamics $P^0(\mathbf{x})$ and risk parameter α . Suppose that for each LC, the optimal desirability function $z^{(i)}$ satisfying the linear BE exists and is bounded away from 0. Then, the optimal desirability for a the LC with terminal cost given by*

$$\ell_f(\mathbf{x}) = \frac{1}{\alpha-1} \log \left(\sum_{i=1}^k w_i \exp \left((\alpha-1) \ell_f^{(i)}(\mathbf{x}) \right) \right), w_i > 0$$

is

$$z(\mathbf{x}) = \sum_{i=1}^k w_i z^{(i)}(\mathbf{x}).$$

z, ℓ are additionally time dependent in the FH case.

Proof. This is a direct consequence of linearity. Let $\beta = \alpha - 1$. For the FE problem, the linear BE is given by

$$\begin{aligned} z^{(i)}(\mathbf{x}) &= \exp(\beta \ell(\mathbf{x})) \left(\int (1 - \mathbb{I}[\mathcal{T}]) z^{(i)} dP^0(\mathbf{x}) + \int \mathbb{I}[\mathcal{T}] \exp(\beta \ell_f^{(i)}(\mathbf{x})) dP^0(\mathbf{x}) \right) \mathbf{x} \notin \mathcal{T} \\ z^{(i)}(\mathbf{x}) &= \exp(\beta \ell_f^{(i)}(\mathbf{x})) \mathbf{x} \in \mathcal{T} \end{aligned}$$

Thus,

$$\begin{aligned} \sum_i w_i z^{(i)}(\mathbf{x}) &= \sum_i w_i \exp(\beta \ell_f^{(i)}(\mathbf{x})) \mathbf{x} \in \mathcal{T} \\ \sum_i w_i z^{(i)}(\mathbf{x}) &= \sum_i w_i \exp(\beta \ell(\mathbf{x})) \left(\int (1 - \mathbb{I}[\mathcal{T}]) z^{(i)} dP^0(\mathbf{x}) + \int \mathbb{I}[\mathcal{T}] \exp(\beta \ell_f^{(i)}(\mathbf{x})) dP^0(\mathbf{x}) \right) \\ &= \exp(\beta \ell(\mathbf{x})) \int (1 - \mathbb{I}[\mathcal{T}]) \left(\sum_i w_i z^{(i)} \right) dP^0(\mathbf{x}) \\ &\quad + \exp(\beta \ell(\mathbf{x})) \int \mathbb{I}[\mathcal{T}] \left(\sum_i w_i \exp(\beta \ell_f^{(i)}(\mathbf{x})) \right) dP^0(\mathbf{x}) \end{aligned}$$

Thus $z = \sum_i w_i z^{(i)}$ satisfies the linear BE for the new terminal cost. Also, since $w_i > 0$ and k is finite, z is also bounded away from 0.

□

NOTE 4. *This result was presented for the risk neutral case in [33]. It was stated in the conference paper [12] for the risk sensitive case without formal proof. Here, we derive it for arbitrary Borel spaces.*

NOTE 5. *This gives a set of non-LQG problems for which we can compute analytical solutions. If we have a FH problem with linear dynamics and quadratic costs, the optimal desirability function can be computed analytically. Now, we can compose the terminal costs $\ell_f(\mathbf{x}) = \frac{1}{\alpha-1} \log \left(\sum_i w_i \exp \left((\alpha-1) \frac{\mathbf{x}^T Q_f^{(i)} \mathbf{x}}{2} \right) \right)$ to get a non quadratic terminal cost. By compositionality, the optimal desirability function for this problem is available analytically by linearly combining those for the constituent problems. The corresponding control policy can be seen as interpolating between the optimal linear feedback policies for each of the constituent problems.*

4.3.2. Path Integral Representation. THEOREM 6. *If there exists a solution z to the linear BE for an FH or FE LC problem, then it has the path integral representation*

$$z(\mathbf{x}) = \mathbb{E}_{\substack{\mathbf{x}_{\tau+1} \sim P^0(\mathbf{x}_\tau), \mathbf{x}_0 = \mathbf{x} \\ T_f = \min \tau: \mathbf{x}_\tau \in \mathcal{T}}} \left[\exp \left((\alpha - 1) \left(\left(\sum_{\tau=0}^{\min(T_f-1, 0)} \ell(\mathbf{x}_\tau) \right) + \ell_f(\mathbf{x}_{T_f}) \right) \right) \right]$$

for FE and

$$z_t(\mathbf{x}) = \mathbb{E}_{\mathbf{x}_{\tau+1} \sim P^0(\mathbf{x}_\tau), \mathbf{x}_0 = \mathbf{x}} \left[\exp \left((\alpha - 1) \left(\left(\sum_{\tau=t}^{T-1} \ell_\tau(\mathbf{x}_\tau) \right) + \ell_f(s_T) \right) \right) \right]$$

for FH.

Proof. For FH the result is easily proven by induction. For FE if $\mathbf{x} \in \mathcal{T}$, $T_f = 0$, and from the formula above, $z(\mathbf{x}) = \exp((\alpha - 1)\ell_f(\mathbf{x}))$. If $\mathbf{x} \notin \mathcal{T}$, then

$$\begin{aligned} \mathbb{E}_{\mathbf{x}' \sim P^0(\mathbf{x})} [z(\mathbf{x}')] &= \mathbb{E}_{\mathbf{x}' \sim P^0(\mathbf{x})} \left[\mathbb{E}_{\substack{\mathbf{x}_{\tau+1} \sim P^0(\mathbf{x}_\tau), \mathbf{x}_0 = \mathbf{x}' \\ T_f = \min \tau: \mathbf{x}_\tau \in \mathcal{T}}} \left[\exp \left((\alpha - 1) \sum_{\tau=0}^{\min(T_f-1, 0)} \ell(\mathbf{x}_\tau) + \ell_f(\mathbf{x}_{T_f}) \right) \right] \right] \\ &= \mathbb{E}_{\substack{\mathbf{x}_{\tau+1} \sim P^0(\mathbf{x}_\tau), \mathbf{x}_0 = \mathbf{x} \\ T_f = \min \tau: \mathbf{x}_\tau \in \mathcal{T}}} \left[\exp \left((\alpha - 1) \sum_{\tau=1}^{\min(T_f-1, 0)} \ell(\mathbf{x}_\tau) + \ell_f(\mathbf{x}_{T_f}) \right) \right] \end{aligned}$$

Multiplying throughout by $\exp((\alpha - 1)\ell(\mathbf{x}))$, we see that the z function defined through the formula above satisfies the linear BE. \square

NOTE 6. *The path integral formulation was developed for continuous time problems by Kappen in [4, 20]. For the LMDP case, it was developed in [34]. It was stated in the conference paper [12] without formal proof for finite state spaces. Here, we derive it for arbitrary Borel spaces.*

4.3.3. Inverse Optimal Control. The inverse optimal control problem is the problem of inferring the state cost $\ell(\mathbf{x})$ of an LC given the risk parameter α , the passive dynamics P^0 and state transitions sampled from the optimal control law $\mathcal{D} = \{\mathbf{x}_i, \mathbf{x}'_i\}_{i=1}^N$ where $\mathbf{x}'_i \sim u^*(\mathbf{x}_i)$.

THEOREM 7. *Consider an LC with a finite state space $\mathcal{X} = \{1, 2, \dots, n\}$ and let $\mathcal{F}[\mathcal{X}]$ be the power set of \mathcal{X} . Let $v(\mathbf{x}; \theta)$ parameterized family of candidate cost-to-gofunctions. The log-likelihood of observing a set of transitions $\{\mathbf{x}_i, \mathbf{x}'_i\}$ under the optimally controlled system (assuming each transition is independent) given θ is*

$$\begin{aligned} L(\theta) &= \sum_i \log(u^*(\mathbf{x}'_i | \mathbf{x}_i; \theta)) \\ &= \sum_i \log(P^0(\mathbf{x}'_i | \mathbf{x}_i) \exp(-v(\mathbf{x}'_i; \theta))) - \log \left(\mathbb{E}_{\mathbf{x}' \sim P^0(\mathbf{x}_i)} [\exp(-v(\mathbf{x}'; \theta))] \right) \\ &= c + \sum_i -v(\mathbf{x}'_i; \theta) - \log \left(\mathbb{E}_{\mathbf{x}' \sim P^0(\mathbf{x}_i)} [\exp(-v(\mathbf{x}'; \theta))] \right) \end{aligned}$$

where c is a constant independent of θ . The maximum likelihood estimate of the state cost is then

$$\ell^*(\mathbf{x}) = v(\mathbf{x}; \theta^*) - \frac{1}{\alpha - 1} \log \left(\mathbb{E}_{P^0(\mathbf{x})} [\exp((\alpha - 1)v(\theta^*))] \right).$$

Further, the maximum likelihood estimate $\theta^* = \max_{\theta} L(\theta)$ can be found efficiently by solving a convex optimization problem if the parametrization is linear $v(\mathbf{x}; \theta) = f(\mathbf{x})^T \theta$.

NOTE 7. *This result was first presented in [11] for the LMDP case and extended to the risk sensitive case in [12]. We restate the result in this section for completeness. This method is computationally more efficient than previous methods for inverse optimal control [1] [43], where one needs to solve the forward optimal control problem repeatedly, which is computationally demanding for large state spaces.*

4.3.4. Characterizing the Most Likely Trajectory under Optimal Control. THEOREM 8. *Suppose the state space \mathcal{X} is finite. Consider the game theoretic interpretation of LCs. The most likely state trajectory of the optimally controlled stochastic system (that includes both the controller and the adversary) starting at \mathbf{x}_0 is the one that minimizes the following objective:*

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_T} \sum_{t=0}^{T-1} -\log(P^0(\mathbf{x}_{t+1} | \mathbf{x}_t)) + (\alpha - 1) \left(\sum_{t=0}^{T-1} \ell_t(\mathbf{x}_t) + \ell_f(\mathbf{x}_T) \right)$$

Proof. The optimally controlled system has the dynamics of the adversary (interpretation 2) and is given by $u^*(\mathbf{x}; t) = P^0(\mathbf{x}) \otimes z_{t+1}$. So the probability of a trajectory $\mathbf{x}_0, \dots, \mathbf{x}_T$ is given by

$$\prod_{t=0}^{T-1} \frac{P^0(\mathbf{x}_t) z_{t+1}(\mathbf{x}_{t+1})}{E_{P^0(\mathbf{x}_t)}[z_{t+1}]}$$

From the linear BE we have $E_{P^0(\mathbf{x})}[z_{t+1}] = \exp((\alpha - 1) \ell_t(\mathbf{x})) z_t(\mathbf{x})$. The above expression then becomes

$$\begin{aligned} & \prod_{t=0}^{T-1} \frac{P^0(\mathbf{x}_{t+1} | \mathbf{x}_t) z_{t+1}(\mathbf{x}_{t+1})}{\exp((\alpha - 1) \ell_t(\mathbf{x}_t)) z_t(\mathbf{x}_t)} = \\ & \frac{1}{z_0(\mathbf{x}_0)} \left(\prod_{t=0}^{T-1} P^0(\mathbf{x}_{t+1} | \mathbf{x}_t) \exp((1 - \alpha) \ell_t(\mathbf{x}_t)) \right) z_T(\mathbf{x}_T) = \\ & \frac{1}{z_0(\mathbf{x}_0)} \left(\prod_{t=0}^{T-1} P^0(\mathbf{x}_{t+1} | \mathbf{x}_t) \exp((1 - \alpha) \ell_t(\mathbf{x}_t)) \right) \exp((\alpha - 1) \ell_f(\mathbf{x}_T)) \end{aligned}$$

The most likely trajectory maximizes the RHS above, or equivalently, minimizes its negative logarithm:

$$-\log(z_0(\mathbf{x}_0)) + \sum_{t=0}^{T-1} -\log(P^0(\mathbf{x}_{t+1} | \mathbf{x}_t)) + (\alpha - 1) \left(\sum_{t=0}^{T-1} \ell_t(\mathbf{x}_t) \right) + \ell_f(\mathbf{x}_T)$$

Since \mathbf{x}_0 is fixed, the first term can be dropped from the objective giving the required result. \square

NOTE 8. *For the risk neutral case, this result was derived in [36] and the risk sensitive extension was stated without formal proof in [12]. Here we prove it formally for finite state spaces.*

NOTE 9. *The result, appropriately phrased, should extend to non-finite state spaces under suitable assumptions. We do not consider those technicalities here. This result gives us intuition about the behavior of the system under the optimal policies in a game theoretic setting: When $\alpha < 1$, the most likely trajectory minimizes state costs while when $\alpha > 1$, the most likely trajectory maximizes state costs. This corresponds to the fact that when $\alpha < 1$, the controller “wins” the game with high probability while when $\alpha > 1$, the adversary “wins” with high probability. This is reflected in the control policy, for $\alpha < 1$, the controller plans so as to minimize long term costs being slightly conservative to counter effects of the adversary. When $\alpha > 1$, the controller has to give up and settle for minimizing the damage caused by the adversary.*

5. Numerical Examples. In this section, we present numerical examples that illustrate the effect of the parameter α . We demonstrate that changing α can have nontrivial effects on the optimal control policy, and the effects changes depending on the type of control problem being solved.

5.1. Crossing a narrow bridge. Consider the problem of driving across a narrow bridge over a pond on a windy day where there is a risk of being blown off of the bridge. The state space consists of the position of the car $\mathbf{x} = (x, y)$ and the controller specifies the velocity (u_x, u_y) . The continuous time dynamics of the system is given by

$$\dot{x} = u_x + \sigma d\omega, \dot{y} = u_y + \sigma d\omega.$$

The state cost reflects a high cost for falling into the pond and a constant small cost elsewhere, that encourages getting to the goal as soon as possible. The situation is depicted in figure 5.1.

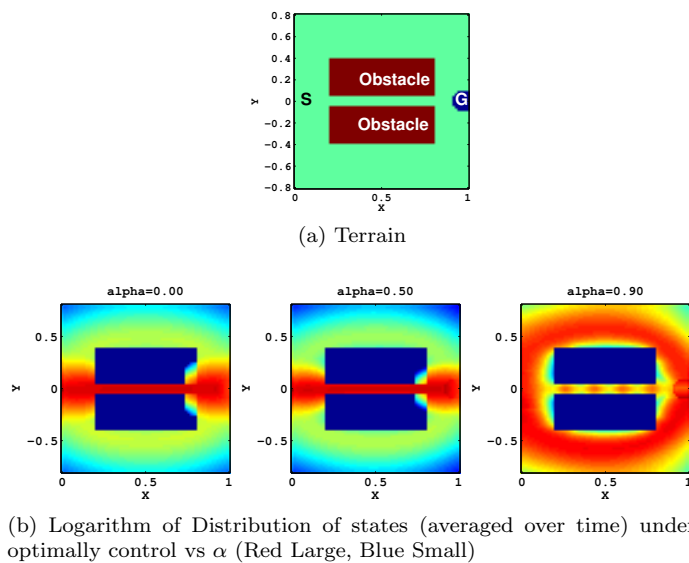


FIG. 5.1. *Driving across a narrow bridge*

We solve this problem by discretizing time (with time step Δ) to get an LRD. We define a passive dynamics $P^0(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \Delta\sigma^2 I)$. This says that the state remains unchanged on average under 0 control (0 velocity) and is perturbed only by noise. The $\Delta\sigma^2 I$ covariance term reflects the fact that the variance of noise for Brownian motion grows linearly with time. When α is small (risk-neutral), the optimal control strategy is to cross the bridge and get to the other side. However, as α becomes large (risk-seeking), the optimal strategy is to play it safe and go around the pond (figure 5.1b).

5.2. Driving up hills. Consider the problem of driving on a slippery hilly terrain, with a cost function that encourages driving up as high as possible. We consider an example with 2 hills, where one hill is higher but steeper than the other (figure 5.2a). When $\alpha = -.5$ (risk-seeking), the optimally controlled distribution peaks at

the higher/steeper hill. When $\alpha = .5$ (risk-averse), the optimally controlled distribution peaks at the lower/less-steep hill. When $\alpha = 0$ (risk-neutral), the probability is shared between both hills.

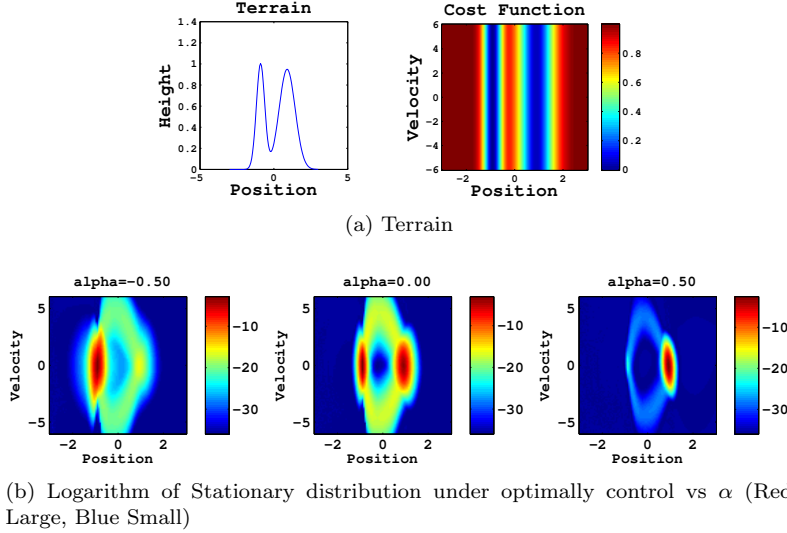


FIG. 5.2. *Driving up hills*

5.3. Stabilizing an unstable system. Consider an unstable system on \mathfrak{R} . The continuous time deterministic dynamics is plotted in 5.3a. We consider a system that has this nominal deterministic dynamics plus zero mean Brownian noise: $d\mathbf{x} = f(\mathbf{x}) + \sigma d\omega$. We can visualize the solution by constructing a “mean” control law: $u_c(\mathbf{x}) = (\mathbb{E}_{\mathbf{x}' \sim u_c^*(\mathbf{x})}[\mathbf{x}'] - \mathbb{E}_{\mathbf{x}' \sim P^0(\mathbf{x})}[\mathbf{x}'])$. This has the interpretation that the control input is the expected change in the system state due to the controller. The control law constructed this way for various values of α is plotted in figure 5.3b. We see that the controls become smaller as α increases. This is reminiscent of several examples in robust control, where high gains could lead to instability when one has imperfect state estimation and robust control policies tend to have lower gains.

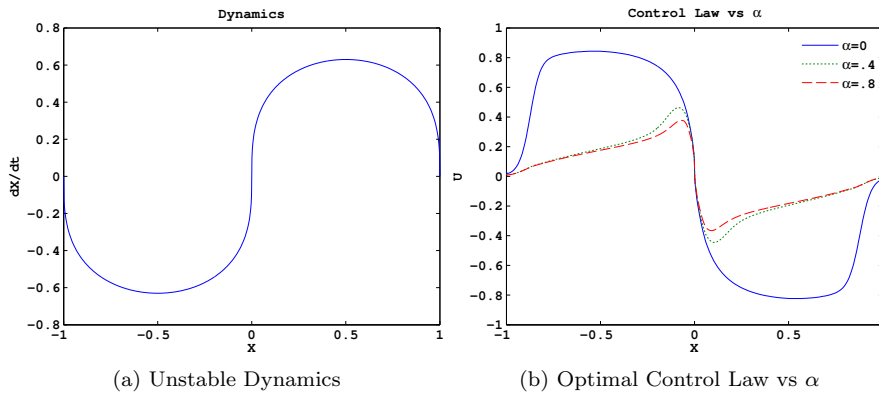


FIG. 5.3. *Controlling an unstable system*

6. Conclusions and future work. Linearly-solvable optimal control is an exciting new development in control theory and has been the subject of many papers over the past few years. In this paper we presented a unified theory of linearly solvable control, with rigorous derivations of all the results for arbitrary compact Borel state spaces. The work so far has been mostly aimed at understanding the framework and its properties. We are now at a stage where the framework is mature and well understood and can lead to the development of algorithms that scale to hard real-world control problems from various application domains. Impressive results in robotics [29] and character animation [8] have recently been obtained. We feel that the surface has barely been scratched in terms of developing more efficient numerical methods for stochastic optimal control.

Appendix A. Proofs.

THEOREM 9. *Let $(\Omega, \mathcal{F}[\Omega], \nu)$ be a probability space. If $\alpha > 1$ and $f : \Omega \mapsto \mathfrak{R}$ is a bounded measurable function, then the problem*

$$\begin{aligned} & \min_{\mu} \max_{\mu'} -\mathbb{D}_{\alpha}(\nu \parallel \mu') + \text{KL}(\mu \parallel \mu') + \mathbb{E}_{\mu}[f] \\ & \text{Subject to } \mu, \mu' \in \mathcal{P}[\Omega] \\ & \mu \ll \nu \ll \mu, \mu' \ll \nu \ll \mu' \end{aligned}$$

has optimum $-\Psi_{\nu}^{\frac{\alpha-1}{\alpha}}[-f]$ with the optimum attained at $\mu^* = \nu \otimes \exp\left(\frac{f}{\alpha}\right), \mu'^* = \mu^* \otimes \exp(-f)$. The results hold even if the order of the min, max are reversed.

Proof. Let μ, μ' be any measures satisfying the constraints of the problem. Then by lemma 3, we have $\mathbb{D}_{\alpha}(\nu \parallel \mu') - \text{KL}(\mu \parallel \mu') \geq \frac{\alpha \text{KL}(\mu \parallel \nu)}{1-\alpha}$. This bound is attained when $\mu' = \mu \otimes \left(\frac{d\mu}{d\nu}\right)^{\frac{\alpha}{1-\alpha}}$ (assuming this measure exists). Also, it is easy to see that this choice for μ' satisfies the constraints. So we have

$$\max_{\mu'} -\mathbb{D}_{\alpha}(\nu \parallel \mu') + \text{KL}(\mu \parallel \mu') = -\min_{\mu'} \mathbb{D}_{\alpha}(\nu \parallel \mu') - \text{KL}(\mu \parallel \mu') = \frac{\alpha}{1-\alpha} \text{KL}(\mu \parallel \nu)$$

We are left with $\frac{\alpha}{1-\alpha} \text{KL}(\mu \parallel \nu) + \mathbb{E}_{\mu}[f]$, which can be lower bounded using lemma 1:

$$\frac{\alpha}{1-\alpha} \left(\text{KL}(\mu \parallel \nu) + \mathbb{E}_{\mu} \left[\frac{(1-\alpha)f}{\alpha} \right] \right) \geq \frac{\alpha}{\alpha-1} \Psi_{\nu} \left[\frac{(1-\alpha)f}{\alpha} \right] = -\Psi_{\nu}^{\frac{\alpha-1}{\alpha}}[-f]$$

Since f is bounded, the probability measure $\mu = \nu \otimes \exp\left(\frac{(1-\alpha)f}{\alpha}\right)$ exists and attains the lower bound above. Thus, the minimum is attained at the measure $\mu^* = \nu \otimes \exp\left(\frac{(1-\alpha)f}{\alpha}\right)$ which satisfies the constraints and the minimum value is $-\Psi_{\nu}^{\frac{\alpha-1}{\alpha}}[f]$.

Given this choice for μ , we can verify that the measure $\mu \otimes \left(\frac{d\mu}{d\nu}\right)^{\frac{\alpha}{1-\alpha}}$ exists (since f is bounded). Thus, the optimal value of the overall minimax problem is $-\Psi_{\nu}^{\frac{\alpha-1}{\alpha}}[-f]$ and the saddle point is given by $\mu^* = \nu \otimes \exp\left(\frac{(1-\alpha)f}{\alpha}\right), \mu'^* = \mu^* \otimes \left(\frac{d\mu^*}{d\nu}\right)^{\frac{\alpha}{1-\alpha}} = \nu \otimes \exp\left(\frac{f}{\alpha}\right)$. This can be rewritten as in the statement of the theorem, $\mu'^* = \nu \otimes \exp\left(\frac{f}{\alpha}\right)$ and $\mu^* = \mu'^* \otimes \exp(-f)$.

If we switch the max and min, we can write $\min_{\mu} \text{KL}(\mu \parallel \mu') + \mathbb{E}_{\mu}[f] = -\Psi_{\mu'}[-f]$

with the minimum attained at $\mu^* = \mu' \otimes \exp(-f)$ (this measure exists since f is bounded) by lemma 1. We are then left with

$$\max_{\mu'} -\mathbb{D}_\alpha(\nu \parallel \mu') - \Psi_{\mu'}[\exp(-f)] = -\min_{\mu'} \mathbb{D}_\alpha(\nu \parallel \mu') + \Psi_{\mu'}[-f]$$

By lemma 2, we know that $\mathbb{D}_\alpha(\nu \parallel \mu') + \Psi_{\mu'}[-f] \geq -\Psi_{\nu^{\frac{1-\alpha}{\alpha}}} [f]$ with the minimum attained at $\mu'^* = \nu \otimes \exp\left(\frac{f}{\alpha}\right)$. \square

LEMMA 1. Let $\mu, \nu \in \mathcal{P}[(\Omega, \mathcal{F}[\Omega])]$, $\mu \ll \nu, \nu \ll \mu$ and $f : \Omega \rightarrow \mathfrak{R}$ be a bounded measurable function. Then, $\text{KL}(\mu \parallel \nu) + \mathbb{E}_\mu[f] \geq -\Psi_\nu[-f]$.

Proof. The objective can be rewritten as

$$\mathbb{E}_\mu \left[\log \left(\frac{d\mu}{d\nu} \right) + f \right] = \mathbb{E}_\mu \left[-\log \left(\frac{d\nu}{d\mu} \right) + f \right] = \mathbb{E}_\mu \left[-\log \left(\frac{d\nu}{d\mu} \exp(-f) \right) \right].$$

By Jensen's inequality, since $-\log$ is convex, the RHS is larger than $-\log \left(\mathbb{E}_\mu \left[\frac{d\nu}{d\mu} \exp(-f) \right] \right) = -\log(\mathbb{E}_\nu[\exp(-f)])$, establishing the result. \square

LEMMA 2. Let $\mu, \nu \in \mathcal{P}[(\Omega, \mathcal{F}[\Omega])]$, $\mu \ll \nu, \nu \ll \mu$ and $f : \Omega \rightarrow \mathfrak{R}$ be a bounded measurable function. Then,

$$\text{sgn}(\alpha) \mathbb{D}_\alpha(\nu \parallel \mu) + \Psi_\mu[f] \geq \Psi_{\nu^{\frac{\alpha-1}{\alpha}}} [f] \text{ if } \alpha > 0$$

$$\text{sgn}(\alpha) \mathbb{D}_\alpha(\nu \parallel \mu) + \Psi_\mu[f] \leq \Psi_{\nu^{\frac{\alpha-1}{\alpha}}} [f] \text{ if } \alpha < 0$$

Proof. Letting $g = \exp(f)$, the LHS can be rewritten as

$$\frac{\log \left(\mathbb{E}_\nu \left[\left(\frac{d\mu}{d\nu} \right)^{1-\alpha} \right] \right)}{\alpha-1} + \log \left(\mathbb{E}_\mu [g] \right) = \log \left(\left(\mathbb{E}_\nu \left[\left(\frac{d\mu}{d\nu} \right)^{1-\alpha} \right] \right)^{\frac{1}{\alpha-1}} \mathbb{E}_\nu \left[\frac{d\mu}{d\nu} g \right] \right).$$

First suppose $\alpha < 0$. Then, using Holder's inequality, we have

$$\begin{aligned} \mathbb{E}_\nu \left[\frac{d\mu}{d\nu} g \right] &\leq \left(\mathbb{E}_\nu \left[\left(\frac{d\mu}{d\nu} \right)^{1-\alpha} \right] \right)^{\frac{1}{1-\alpha}} \left(\mathbb{E}_\nu \left[(g)^{\frac{1-\alpha}{\alpha}} \right] \right)^{\frac{-\alpha}{1-\alpha}} \implies \\ \mathbb{E}_\nu \left[\frac{d\mu}{d\nu} g \right] \left(\mathbb{E}_\nu \left[\left(\frac{d\mu}{d\nu} \right)^{1-\alpha} \right] \right)^{\frac{1}{\alpha-1}} &\leq \left(\mathbb{E}_\nu \left[(g)^{\frac{1-\alpha}{\alpha}} \right] \right)^{\frac{-\alpha}{1-\alpha}} = \left(\mathbb{E}_\nu \left[\exp \left(\frac{(\alpha-1)f}{\alpha} \right) \right] \right)^{\frac{\alpha}{\alpha-1}}. \end{aligned}$$

Taking log on both sides, we have the result. The measure $\nu \otimes (g)^{\frac{\alpha-1}{\alpha}}$ exists since f is bounded. Also note that $h(x) = x^{\frac{1}{1-\alpha}}$ is convex if $\alpha > 0, \alpha \neq 1$. The first term inside the log can be bounded as follows:

$$\begin{aligned} \left(\mathbb{E}_\nu \left[\left(\frac{d\mu}{d\nu} \right)^{1-\alpha} \right] \right)^{\frac{1}{1-\alpha}} &= \left(\mathbb{E}_{\nu \otimes g^{\frac{\alpha-1}{\alpha}}} \left[g^{\frac{1-\alpha}{\alpha}} \left(\frac{d\mu}{d\nu} \right)^{1-\alpha} \right] \right)^{\frac{1}{1-\alpha}} \left(\mathbb{E}_\nu \left[g^{\frac{\alpha-1}{\alpha}} \right] \right)^{\frac{1}{1-\alpha}} \\ &\leq \mathbb{E}_{\nu \otimes g^{\frac{\alpha-1}{\alpha}}} \left[g^{\frac{1}{\alpha}} \frac{d\mu}{d\nu} \right] \left(\mathbb{E}_\nu \left[g^{\frac{\alpha-1}{\alpha}} \right] \right)^{\frac{1}{1-\alpha}} \text{ (Jensen's Inequality)} \\ &= \mathbb{E}_\nu \left[g^{\frac{\alpha-1}{\alpha} + \frac{1}{\alpha}} \frac{d\mu}{d\nu} \right] \left(\mathbb{E}_\nu \left[g^{\frac{\alpha-1}{\alpha}} \right] \right)^{\frac{1}{1-\alpha}-1} = \mathbb{E}_\nu \left[g \frac{d\mu}{d\nu} \right] \left(\mathbb{E}_\nu \left[g^{\frac{\alpha-1}{\alpha}} \right] \right)^{\frac{\alpha}{1-\alpha}} \end{aligned}$$

Rewriting the last inequality, we get $\left(E_\nu \left[g^{\frac{\alpha-1}{\alpha}} \right]\right)^{\frac{\alpha}{\alpha-1}} \leq E_\mu [g] \left(E_\nu \left[\left(\frac{d\mu}{d\nu}\right)^{1-\alpha} \right]\right)^{\frac{1}{1-\alpha}}$.

Taking log on both sides gives the result.

□

LEMMA 3. *Let $\mu, \mu', \nu \in \mathcal{P}[(\Omega, \mathcal{F}[\Omega])]$, $\mu' \ll \nu, \nu \ll \mu', \mu \ll \mu', \mu' \ll \mu$ and $f : \Omega \rightarrow \mathfrak{R}$ be a bounded measurable function, $\alpha > 1$. Then, $\mathbb{D}_\alpha(\nu \parallel \mu') - \text{KL}(\mu \parallel \mu') \geq \frac{\alpha}{1-\alpha} \text{KL}(\mu \parallel \nu)$.*

Proof. By definition,

$$\text{KL}(\mu \parallel \mu') = E_\mu \left[-\log \left(\frac{d\mu'}{d\mu} \right) \right], \mathbb{D}_\alpha(\nu \parallel \mu') = \frac{\log \left(E_\nu \left[\left(\frac{d\mu'}{d\nu} \right)^{1-\alpha} \right] \right)}{\alpha - 1}.$$

$$\text{Now, } E_\nu \left[\left(\frac{d\mu'}{d\nu} \right)^{1-\alpha} \right] = E_\mu \left[\frac{d\nu}{d\mu} \left(\frac{d\mu'}{d\nu} \right)^{1-\alpha} \right] = E_\mu \left[\frac{\frac{d\mu'}{d\mu} \left(\frac{d\mu'}{d\nu} \right)^{1-\alpha}}{\frac{d\mu'}{d\nu}} \right] = E_\mu \left[\frac{d\mu'}{d\mu} \left(\frac{d\mu'}{d\nu} \right)^{-\alpha} \right].$$

This gives us

$$\mathbb{D}_\alpha(\nu \parallel \mu') - \text{KL}(\mu \parallel \mu') = \frac{\log \left(E_\mu \left[\frac{d\mu'}{d\mu} \left(\frac{d\mu'}{d\nu} \right)^{-\alpha} \right] \right)}{\alpha - 1} + E_\mu \left[\log \left(\frac{d\mu'}{d\mu} \right) \right]$$

When $\alpha > 1$, $\frac{\log}{\alpha-1}$ is concave and by Jensen's inequality the first term is $\geq \frac{E_\mu \left[\log \left(\frac{d\mu'}{d\mu} \left(\frac{d\mu'}{d\nu} \right)^{-\alpha} \right) \right]}{\alpha-1} = \frac{E_\mu \left[-\alpha \log \left(\frac{d\mu'}{d\nu} \right) + \log \left(\frac{d\mu'}{d\mu} \right) \right]}{\alpha-1}$. This implies:

$$\mathbb{D}_\alpha(\nu \parallel \mu') - \text{KL}(\mu \parallel \mu') \geq \frac{\alpha}{\alpha-1} E_\mu \left[\log \left(\frac{d\mu'}{d\mu} \right) \right] = \frac{\alpha}{\alpha-1} E_\mu \left[\log \left(\frac{d\nu}{d\mu} \right) \right] = \frac{\alpha \text{KL}(\mu \parallel \nu)}{1-\alpha}.$$

□

REFERENCES

- [1] P. ABBEEL AND A.Y. NG, *Apprenticeship learning via inverse reinforcement learning*, in Proceedings of the twenty-first international conference on Machine learning, ACM, 2004, p. 1.
- [2] T. BAŞAR AND G.J. OLSDER, *Dynamic noncooperative game theory*, Society for Industrial Mathematics, 1999.
- [3] R.E. BELLMAN AND RAND CORPORATION, *Dynamic programming*, Rand Corporation research study, Princeton University Press, 1957.
- [4] B. VAN DEN BROEK, W. WIEGERINCK, AND B. KAPPEN, *Risk sensitive path integral control*, in Uncertainty in AI, 2010. Proceedings of the 2010, 2010.
- [5] J. BROEK, W. WIEGERINCK, AND KAPPEN H., *Stochastic optimal control of state constrained systems*, International Journal of Control, (2011), pp. 1–9.
- [6] J. BROEK, W. WIEGERINCK, AND H. KAPPEN, *Risk sensitive path integral control*, Uncertainty in Artificial Intelligence, (2010).
- [7] R. CAVAZOS-CADENA, *Optimality equations and inequalities in a class of risk-sensitive average cost markov decision chains*, Mathematical Methods of Operations Research, 71 (2010), pp. 47–84.
- [8] M. DA SILVA, F. DURAND, AND J. POPOVIĆ, *Linear Bellman combination for control of character animation*, ACM Transactions on Graphics (TOG), 28 (2009), pp. 1–10.
- [9] G.B. DI MASI AND L. STETTNER, *Risk-sensitive control of discrete-time markov processes with infinite horizon*, SIAM Journal on Control and Optimization, 38 (1999), p. 61.
- [10] P. DUPUIS AND R. S. ELLIS, *A Weak Convergence Approach to the Theory of Large Deviations*, Wiley, New York, NY, 1997.

- [11] K. DVIJOTHAM AND E. TODOROV, *Inverse optimal control with linearly solvable MDPs*, International Conference on Machine Learning, (2010).
- [12] KRISHNAMURTHY DVIJOTHAM AND EMANUEL TODOROV, *A unifying framework for linearly solvable control*, in UAI, Fabio Gagliardi Cozman and Avi Pfeffer, eds., AUAI Press, 2011, pp. 179–186.
- [13] W.H. FLEMING, *Exit probabilities and optimal stochastic control*, Applied Mathematics & Optimization, 4 (1977), pp. 329–346.
- [14] W. FLEMING AND W. MCENEANEY, *Risk sensitive optimal control and differential games*, Stochastic Theory and Adaptive Control, (1992), pp. 185–197.
- [15] W. FLEMING AND S. MITTER, *Optimal control and nonlinear filtering for nondegenerate diffusion processes*, Stochastics, 8 (1982), pp. 226–261.
- [16] W.H. FLEMING AND H.M. SONER, *Controlled Markov processes and viscosity solutions*, vol. 25, Springer Verlag, 2006.
- [17] M. K. GHOSH AND A. BAGCHI, *Stochastic games with average payoff criterion*, Applied Mathematics & Optimization, 38 (1998), pp. 283–301. 10.1007/s002459900092.
- [18] C. HOLLAND, *A new energy characterization of the smallest eigenvalue of the Schrödinger equation*, Comm Pure Appl Math, 30 (1977), pp. 755–765.
- [19] E. HOPF, *The partial differential equation $ut + uux = \mu xx$* , Communications on Pure and Applied mathematics, 3 (1950), pp. 201–230.
- [20] H. KAPPEN, *Linear theory for control of nonlinear stochastic systems*, Physical Review Letters, 95 (2005).
- [21] M.A. KRASNOSELSKI, E.A. LIFSHITS, AND A.V. SOBOLEV, *Positive linear systems: the method of positive operators*, vol. 5, Heldermann, 1989.
- [22] H.J. KUSHNER AND P. DUPUIS, *Numerical methods for stochastic control problems in continuous time*, vol. 24, Springer Verlag, 2001.
- [23] S.I. MARCUS, E. FERNÁNDEZ-GAUCHERAND, D. HERNÁNDEZ-HERNANDEZ, S. CORALUPPI, AND P. FARD, *Risk sensitive Markov decision processes*, Systems and Control in the Twenty-First Century, 29 (1997).
- [24] T. MENSINK, J. VERBEEK, AND H. KAPPEN, *EP for efficient stochastic control with obstacles*, ECAI, (2010).
- [25] S. MITTER AND N. NEWTON, *A variational approach to nonlinear estimation*, SIAM J Control Opt, (2003), pp. 1813–1833.
- [26] B.K. ØKSENDAL, *Stochastic differential equations: an introduction with applications*, Springer Verlag, 2003.
- [27] A. RENYI, *On measures of entropy and inf*, in Proc. 4th Berkley symp. Stat. and Prob, vol. 1, pp. 541–561.
- [28] L.S. SHAPLEY, *Stochastic games*, Proceedings of the National Academy of Sciences of the United States of America, 39 (1953), p. 1095.
- [29] E. THEODOROU, J. BUCHLI, AND S. SCHAAL, *Learning policy improvements with path integrals*, AISTATS, (2010).
- [30] E. A. THEODOROU, *Iterative Path Integral Stochastic Optimal Control: Theory and Applications to Motor Control*, PhD thesis, University of Southern California, 2011.
- [31] E. TODOROV, *Linearly-solvable Markov decision problems*, Advances in Neural Information Processing Systems, (2006).
- [32] EMANUEL TODOROV, *General duality between optimal control and estimation*, in CDC, IEEE, 2008, pp. 4286–4292.
- [33] E. TODOROV, *Compositionality of optimal control laws*, in NIPS, Yoshua Bengio, Dale Schuurmans, John D. Lafferty, Christopher K. I. Williams, and Aron Culotta, eds., Curran Associates, Inc., 2009, pp. 1856–1864.
- [34] ———, *Efficient computation of optimal actions*, PNAS, 106 (2009), pp. 11478–11483.
- [35] EMANUEL TODOROV, *Policy gradients in linearly-solvable mdps*, in NIPS, John D. Lafferty, Christopher K. I. Williams, John Shawe-Taylor, Richard S. Zemel, and Aron Culotta, eds., Curran Associates, Inc., 2010, pp. 2298–2306.
- [36] E. TODOROV, *Finding the most likely trajectories of optimally-controlled stochastic systems*, in World Congress of the International Federation of Automatic Control (IFAC), 2011.
- [37] M. TOUSSAINT, *Robot trajectory optimization using approximate inference*, International Conference on Machine Learning, 26 (2009), pp. 1049–1056.
- [38] TIM VAN ERVEN, *When Data Compression and Statistics Disagree*, PhD thesis, CWI, 2010.
- [39] J. VON NEUMANN, O. MORGENSTERN, A. RUBINSTEIN, AND H.W. KUHN, *Theory of games and economic behavior*, Princeton Univ Pr, 2007.
- [40] W. WIEGERINCK, B. BROEK, AND H. KAPPEN, *Stochastic optimal control in continuous space-time multi-AgentSystems*, 22nd annual conference on Uncertainty in Artificial Intelligence,

- (2006).
- [41] M. ZHONG AND E. TODOROV, *Aggregation methods for linearly-solvable MDPs*, IFAC World Congress, (2011).
 - [42] ———, *Moving least-squares approximations for linearly-solvable stochastic optimal control problems*, Journal of Control Theory and Applications, (2011).
 - [43] B.D. ZIEBART, A. MAAS, J.A. BAGNELL, AND A.K. DEY, *Maximum entropy inverse reinforcement learning*, in Proc. AAAI, 2008, pp. 1433–1438.