$\frac{1}{2}$

CHOQUET REGULARIZATION FOR CONTINUOUS-TIME REINFORCEMENT LEARNING

3

41

42

XIA HAN*, RUODU WANG[†], AND XUN YU ZHOU[‡]

Abstract. We propose *Choquet regularizers* to measure and manage the level of exploration for reinforcement learning (RL), and reformulate the continuous-time entropy-regularized RL problem of [47] in which we replace the differential entropy used for regularization with a Choquet regularizer. We derive the Hamilton–Jacobi–Bellman equation of the problem, and solve it explicitly in the linear–quadratic (LQ) case via maximizing statically a mean–variance constrained Choquet regularizer. Under the LQ setting, we derive explicit optimal distributions for several specific Choquet regularizers, and conversely identify the Choquet regularizers that generate a number of broadly used exploratory samplers such as ε -greedy, exponential, uniform and Gaussian.

12 **Key words.** Reinforcement learning, Choquet integrals, continuous time, exploration, regular-13 izers, quantile, HJB equations, linear-quadratic control

14 **MSC codes.** 93E20, 93E35

15 1. Introduction. Reinforcement learning (RL) is one of the most active and fast developing subareas in machine learning. The foundation of RL is "trial and error" -16to *strategically* explore different action plans in order to find the best plan as efficiently 17 and economically as possible. A key question to this inherent exploratory approach 18for RL is to seek a proper tradeoff between exploration and exploitation, for which 1920 one needs to first quantify the level of exploration. Because exploration is typically 21 captured by randomization in the RL study, entropy has been employed to measure the magnitude of the randomness and hence that of the exploration – a uniform dis-22 tribution representing a completely blind search has the maximum entropy while a 23 Dirac mass signifying no exploration at all has the minimum entropy. Discrete-time 24entropy-regularized (or "softmax") RL formulation has been proposed which intro-2526 duces a weighted entropy value of the exploration as a regularization term into the objective function ([23, 33, 54]). For continuous-time RL, [47] formulate an entropy-27regularized, distribution-valued stochastic control problem for diffusion processes, and 28 derive theoretically the Gibbs (or Boltzmann) measure as the optimal distribution for 29exploration which specializes to Gaussian when the problem is linear-quadratic (LQ). 30 [18] and [48] apply the results of [47] to a Langevin diffusion for simulated annealing and a continuous-time entropy-regularized Markowitz's mean-variance portfolio 32 selection problem, respectively. [22] analyze both quantitatively and qualitatively the 33 impact of entropy regularization for mean-field games with learning in a finite time 34 horizon. There have been recently many other developments along this direction of RL in continuous time; see [25, 26, 27, 32, 43] and the references therein. 36

While the entropy is a reasonable metric to quantify the information gain of exploring the environment and entropy regularization can indeed explain some broadly used exploration distributions such as Gaussian, there are two closely related open questions:

1. Distributions other than Gaussian, such as exponential or uniform, are also widely used for exploration in RL. What regularizer(s) can theoretically jus-

^{*}School of Mathematical Sciences and LPMC, Nankai University, China. Email: xia-han@nankai.edu.cn

 $^{^\}dagger Department$ of Statistics and Actuarial Science, University of Waterloo, Canada. Email: wang@uwaterloo.ca

 $^{^{\}ddagger}\mathrm{D}\mathrm{e}\mathrm{partment}$ of IEOR, Columbia University, USA. Email: xz2574@columbia.edu

43 tify the use of a given class of exploratory distributions?

2. What are the optimal exploratory distributions for regularizers other than the entropy?

In this paper, we study these two questions in the setting of continuous-time 46 diffusion processes, by introducing a new class of regularizers borrowing from the 47 literature of risk metrics. Risk metrics, roughly speaking, include risk measures and 48 variability measures, which are two separate and active research streams in the general 49area of risk management. Value-at-risk (VaR), expected shortfall (ES) and various 50coherent or convex risk measures, introduced by [3, 11, 15], are popular examples of risk measures. Variance, the Gini deviation, interquantile range and deviation measures of [39] are instances of variability measures. There has been a rich body of 53 study on risk metrics in the past two decades; see [16] and the references therein. 54

We introduce what we call *Choquet regularizers*, which belong to the class of the signed Choquet integrals recently studied by [50] in the context of risk management. A signed Choquet integral in general gives rise to a nonlinear and non-monotone expectation in which the state of nature is weighted by a probability weighting or distortion function in calculating the expectation. It includes as special cases Yaari's dual utility ([53]) and distortion risk measures ([2, 29]), which are commonly used monotone functionals, and appears in rank-dependent utility (RDU) theory; see [10, 19, 37, 46] in the related literature of behavioral economic theory.

There are several reasons to use Choquet regularizers for RL due to a number 63 of theoretical and practical advantages. First, they satisfy several "good" properties 64 65 such as quantile additivity, normalization, concavity, and consistency with convex order (mean-preserving spreads) that facilitate analysis as regularizers. Second, Cho-66 quet regularizers are non-monotone. In order to measure exploration, monotonicity is 67 irrelevant, in contrast to assessing risk or wealth. For instance, a degenerate distribu-68 tion should be associated with no-exploration regardless of its position, in which case 69 non-monotone mappings should be used. Moreover, the use of Choquet regularizers is 7071closely connected to distributionally robust optimization (DRO) where a Wasserstein distance naturally induces a special class of Choquet regularizers, whereas DRO itself 72is an important approach for learning and for correcting the inherent flaws suffered by 73 classical model-based estimation and optimization. Finally, as we will see later in the 74paper, for any given location–scale class of distributions, there exists a common Cho-75quet regularizer such that the corresponding regularized continuous-time LQ control 76 77 for RL has optimal distributions in that class.

We take the same continuous-time exploratory stochastic control problem as in 78[47], except that the entropy regularizer is replaced with a Choquet regularizer. In 79 the general case we derive the Hamilton–Jacobi–Bellman (HJB) equation. However, 80 in sharp contrast to [47] in which the optimal control distributions are proved to be 81 82 Gibbs measures, obtaining the class of optimal distributional policies via verification theorem remains a significant open question. To obtain explicit solutions, we focus 83 on the LQ case. The form of the LQ-specialized HJB equation suggests that the 84 problem boils down to a static optimization in which the given Choquet regularizer 85 86 is to be maximized over distributions with given mean and variance. It turns out this last problem has been solved explicitly by [31]. The optimal distributions form a 87 88 location-scale family, whose shape depends on the choices of the Choquet regularizer. The solutions to the static problem are then employed to solve the original LQ-89 based exploratory HJB equation explicitly and to derive the optimal samplers for 90 exploration under the given Choquet regularizer. As expected, optimal distributions 91 are no longer necessarily Gaussian as in [47], and are now dictated by the choice 92

 $\mathbf{2}$

44

CHOQUET REGULARIZATION

93 of Choquet regularizers. However, the following feature of the entropy-regularized 94 solutions revealed in [47] remains intact: the means of the optimal distributions are 95 linear in the current state and independent of the exploration, whereas the variances 96 are determined by the exploration but irrespective of the current state. Along an 97 opposite line of inquiry, we are able to identify a proper Choquet regularizer in order 98 to interpret a given exploratory distribution. Specifically, we derive the regularizers 99 that generate some common exploration measures including ε -greedy, three-point, 90 exponential, uniform and Gaussian.

The rest of the paper is organized as follows. We introduce Choquet regularizers 101 in Section 2, and present their basic properties as well as an axiomatic characterization 102based on existing results of [49, 50]. In Section 3, we formulate the continuous-time 103 Choquet-regularized RL control problem and derive the HJB equation. We then 104 motivate a mean-variance constrained Choquet regularizer maximization problem for 105LQ control. This problem is studied in details in Section 4, including discussions 106 on some special regularizers arising from problems in finance, optimization, and risk 107 management. In Section 5, we return to the exploratory LQ control problem and solve 108 109 it completely. We also present examples linking specific exploratory distributions with 110 the corresponding Choquet regularizers. In Section 6, we discuss the connections between the exploratory LQ problem and the classical LQ problem. Finally, Section 111 7 concludes the paper. 112

113 **2.** Choquet regularizers. Throughout the paper, we assume that $(\Omega, \mathcal{F}, \mathbb{P})$ is an atomless probability space. With a slight abuse of notation, let \mathcal{M} denote both 114 115 the set of (probability) distribution functions of real random variables and the set of Borel probability measures on \mathbb{R} , with the obvious identity $\Pi(x) \equiv \Pi((-\infty, x])$ for 116 $x \in \mathbb{R}$ and $\Pi \in \mathcal{M}$. We denote by $\mathcal{M}^p \subset \mathcal{M}, p \in [1,\infty)$, the set of distribution 117functions or probability measures with finite p-th moment. For a random variable 118X and a distribution Π , we write $X \sim \Pi$ if the distribution of X is Π under \mathbb{P} , and 119 $X \stackrel{\mathrm{d}}{=} Y$ if two random variables X and Y have the same distribution. We denote by 120 μ and σ^2 the mean and variance functionals on \mathcal{M}^2 , respectively; that is, $\mu(\Pi)$ is the 121mean of Π and $\sigma^2(\Pi)$ the variance of Π for $\Pi \in \mathcal{M}^2$. 122

123 Given a function $h: [0,1] \to \mathbb{R}$ of bounded variation with h(0) = 0 and $\Pi \in \mathcal{M}$, 124 the functional I_h on \mathcal{M} is defined as

(2.1)
125
$$I_h(\Pi) \equiv \int h \circ \Pi([x,\infty)) \mathrm{d}x := \int_{-\infty}^0 \left[h \circ \Pi([x,\infty)) - h(1)\right] \mathrm{d}x + \int_0^\infty h \circ \Pi([x,\infty)) \mathrm{d}x,$$

assuming that Equation (2.1) is well defined (which could take the value ∞). The function h is called a distortion function, and the functional I_h is called a signed *Choquet integral* by [50]. If $h(x) \equiv x$ then I_h reduces to the mean functional; thus, I_h is a nonlinear generalization of the mean/expectation. If h is increasing and satisfies h(0) = 1 - h(1) = 0, then I_h is called an *increasing Choquet integral*, which include as special cases the two most important risk measures used in current banking and insurance regulation, VaR and ES.¹

Next, we define the *Choquet regularizer*, a main object of this paper. We are particularly interested in a subclass of signed Choquet integrals, where h satisfies the properties: (i) h is concave, and (ii) h(1) = h(0) = 0.

¹This functional I_h is termed differently in different fields. For example, it is known as Yaari's dual utility ([53]) in decision theory, distorted premium principles ([12, 52]) in insurance and distortion risk measures ([2, 29]) in finance.

136 Let us briefly explain the interpretations and implications of the above two condi-137tions. Condition (i) is equivalent to several other properties, and in particular, to that I_h is a concave mapping and to that I_h is consistent with convex order;² see Theorem 138 3 of [50] for this claim and several other equivalent properties. Here, concavity of I_h 139means $I_h(\lambda \Pi_1 + (1-\lambda)\Pi_2) \ge \lambda I_h(\Pi_1) + (1-\lambda)I_h(\Pi_2)$ for all $\Pi_1, \Pi_2 \in \mathcal{M}$ and $\lambda \in [0,1]$, 140 and consistency with convex order means $I_h(\Pi_1) \leq I_h(\Pi_2)$ for all $\Pi_1, \Pi_2 \in \mathcal{M}$ with 141 $\Pi_1 \preceq_{cx} \Pi_2$. If $\Pi_1 \preceq_{cx} \Pi_2$, then Π_2 is also called a *mean-preserving spread* of Π_1 ([40]), 142 which intuitively means that Π_2 is more spread-out (and hence "more random") than 143 Π_1 . The above two properties do indeed suggest that $I_h(\Pi)$ serves as a measure of 144randomness for Π , since both a mixture and a mean-preserving spread introduce extra 145randomness; see e.g., [1] for related discussions. Condition (ii), on the other hand, is 146 147equivalent to $I_h(\delta_c) = 0 \ \forall c \in \mathbb{R}$, where δ_c is the Dirac mass at c. That is, degenerate distributions do not have any randomness measured by I_h . 148

149 DEFINITION 2.1. Let \mathcal{H} be the set of $h : [0,1] \to \mathbb{R}$ satisfying (i)-(ii). A functional 150 $\Phi : \mathcal{M} \to (-\infty, \infty]$ is a Choquet regularizer if there exists $h \in \mathcal{H}$ such that $\Phi = I_h$, 151 that is,

152 (2.2)
$$\Phi(\Pi) = \int_{\mathbb{R}} h \circ \Pi([x, \infty)) dx,$$

153 and this Φ will henceforth be denoted by Φ_h .

To clarify on notation, we require $h \in \mathcal{H}$ for Φ_h , while there is no such requirement for I_h . Moreover, we call Φ_h to be location invariant and scale homogeneous if $\Phi_h(\Pi') = \lambda \Phi_h(\Pi)$ where Π' is the distribution of $\lambda X + c$ for $\lambda > 0, c \in \mathbb{R}$ and $X \sim \Pi$. We summarize some useful properties of Φ_h in the following lemma.

158 LEMMA 2.2. For $h \in \mathcal{H}$, Φ_h is well defined, non-negative, and location invariant 159 and scale homogeneous.

160 Proof. First, a concave h with h(0) = h(1) has to be first increasing and then 161 decreasing on [0, 1]. Hence h has bounded variation, and the two integrals in Equa-162 tion (2.1) are well defined. Moreover, (i) and (ii) imply $h \ge 0$, which further yields 163 that both terms in Equation (2.1) are non-negative. So Φ_h is well defined and non-164 negative. Location invariance and scale homogeneity follow from Proposition 2 (iii) 165 and (iv) of [49].

Each property in Lemma 2.2 has a simple interpretation for a regularizer that measures the level of randomness, or the level of exploration in the RL context of this paper.

- (a) Well-posedness: Any distribution for exploration can be measured.³
- 169 (b) Non-negativity: Randomness is measured in non-negative values.
- (c) Location invariance: The measurement of randomness does not depend on thelocation.
- 172 (d) Scale homogeneity: The measurement of randomness is linear in its scale.

For a distribution $\Pi \in \mathcal{M}$, let its left-quantile for $p \in (0, 1]$ be defined as, recalling that $\Pi(x) = \Pi((-\infty, x])$ for $x \in \mathbb{R}$,

$$Q_{\Pi}(p) = \inf \left\{ x \in \mathbb{R} : \Pi(x) \ge p \right\},\$$

²Let Π_1 and Π_2 be two distribution functions with finite means. Then, Π_1 is smaller than Π_2 in *convex order*, denoted by $\Pi_1 \preceq_{cx} \Pi_2$, if $\mathbb{E}[f(\Pi_1)] \leq \mathbb{E}[f(\Pi_2)]$ for all convex functions f, provided that the two expectations exist. It is immediate that $\Pi_1 \preceq_{cx} \Pi_2$ implies $\mathbb{E}[\Pi_1] \leq \mathbb{E}[\Pi_2]$.

³This property is technically important since functional properties of I_h could be very difficult to analyze if one faces a quantity such as $\infty - \infty$. As an example, consider h(x) = x leading to I_h being the mean functional. In this case, I_h is only well defined on some subsets of \mathcal{M} .

whereas its right-quantile function for $p \in [0, 1)$ be defined as

$$Q_{\Pi}^{+}(p) = \inf \{ x \in \mathbb{R} : \Pi(x) > p \}.$$

It is useful to note that Φ_h admits a quantile representation as follows; see Lemma 1 173of [49]. 174

LEMMA 2.3. For $h \in \mathcal{H}$ and $\Pi \in \mathcal{M}$, 175

(i) if h is right-continuous, then $\Phi_h(\Pi) = \int_0^1 Q_\Pi^+(1-p) dh(p);$ (ii) if h is left-continuous, then $\Phi_h(\Pi) = \int_0^1 Q_\Pi(1-p) dh(p);$ (iii) if Q_Π is continuous, then $\Phi_h(\Pi) = \int_0^1 Q_\Pi(1-p) dh(p).$ 176

177

178

Choquet regularizers include, for instance, range, mean-median deviation, the 179180Gini deviation, and inter-ES differences. Moreover, standard deviation can be written as the supremum of Choquet regularizers; see Examples 1, 3 and 4 of [50]. 181 Variance also has a related representation (Example 2.2 of [31]) given as $\sigma^2(\Pi) =$ 182 $\sup_{h \in \mathcal{H}} \{\Phi_h(\Pi) - \frac{1}{4} ||h'||_2^2\} \text{ for any } \Pi \in \mathcal{M}, \text{ where } ||h'||_2^2 = \int_0^1 (h'(p))^2 dp \text{ if } h \text{ is continuous with a.e. right-derivative } h', \text{ and } ||h'||_2^2 := \infty \text{ if } h \text{ is not continuous.}$ 183 184

Concave signed Choquet integrals are characterized by, e.g., [50], which is essen-185 tially a consequence of the seminal works of [41, 53]; see also Theorem 2.4 below. 186 In what follows, we say that $\Phi = \Phi_h$ is quantile additive if for all $\Pi_1, \Pi_2 \in \mathcal{M}$, 187 $\Phi(\Pi_1 \oplus \Pi_2) = \Phi(\Pi_1) + \Phi(\Pi_2)$ where the quantile function of $\Pi_1 \oplus \Pi_2$ is the sum of 188 those of Π_1 and Π_2 . In other words, $Q_{\Pi_1 \oplus \Pi_2} = Q_{\Pi_1} + Q_{\Pi_2}$. Moreover, we say that Φ is continuous at infinity if $\lim_{M \to 1} \Phi((\Pi \land M) \lor (1 - M)) = \Phi(\Pi)$, and Φ is uniform 189 190 sup-continuity if for any $\varepsilon > 0$, there exists $\delta > 0$, such that $|\Phi(\Pi_1) - \Phi(\Pi_2)| < \varepsilon$ 191 whenever ess-sup $|\Pi_1 - \Pi_2| < \delta$, where ess-sup is the essential supremum defined by 192 $\Pi^{-1}(1)$. 193

We give the following simple characterization for our Choquet regularizers based 194on Theorems 1 and 3 of [49]. 195

THEOREM 2.4. A functional Φ_h is a Choquet regularizer in Equation (2.2) if and 196197 only if it satisfies all of the following properties

- (i) Φ_h is quantile additive; 198
- (ii) Φ_h is concave or \leq_{cx} -consistent; 199

(iii) $\Phi_h \ge 0$ and $\Phi_h(\delta_c) = 0$ for all $c \in \mathbb{R}$; 200

(iv) Φ_h is continuous at infinity and uniformly sup-continuous. 201

Note that Theorems 1 and 3 of [49] are stated in terms of a risk measure defined 202 on the space of real random variables, say \mathcal{X} , while here Φ_h is defined on \mathcal{M} . To 203 use these results, we can define $\rho: \mathcal{X} \to \mathbb{R}$ by $\rho(X) = \Phi_h(\Pi)$ where $X \sim \Pi$, which 204is automatically law-invariant.⁴ On the other hand, Theorem 1 in [49] requires an 205extra continuity condition to imply that h has bounded variation on [0, 1], which is 206guaranteed here by condition (iii). In fact, condition (i) is equivalent to comonotonic 207additivity of ρ .⁵ Continuity at infinity and uniform sup-continuity of ρ can be defined 208in parallel to those of Φ_h . Finally, h(1) = h(0) = 0 is equivalent to $\Phi_h(\delta_c) = 0$ for all 209 $c \in \mathbb{R}$. Theorem 2.4 hence follows directly from Theorems 1 and 3 of [49]. 210

Remark 2.5. If h is not constantly 0, Choquet regularizers belong to the class of 211 212 generalized deviation measures in [21] and [39]. Moreover, Choquet regularizers can

⁴Law-invariance means that $\rho(X) = \rho(Y)$ for $X \stackrel{d}{=} Y$.

⁵A random vector (X_1, \ldots, X_n) is called *comonotonic* if there exists a random variable $Z \in \mathcal{X}$ and increasing functions f_1, \ldots, f_n on \mathbb{R} such that $X_i = f_i(Z)$ almost surely for all $i = 1, \ldots, n$. Comonotoic-additivity means that $\rho(X+Y) = \rho(X) + \rho(Y)$ if X and Y are comonotonic.

be used to construct law-invariant generalized deviation measures. Indeed, combining characterization of generalized deviation measures in Proposition 2.2 of [21] and the quantile representation of signed Choquet integrals in Lemma 2.3, all law-invariant generalized deviation measures can be represented as a supremum of some Choquet regularizers of the type Equation (2.2). This includes standard deviation and mean absolute deviation as special cases.

We conclude this section by comparing the Choquet regularization with the differential entropy regularization, the latter having been used for exploration–exploitation balance in RL; see [22, 47, 48]. For an absolutely continuous Π, we define DE, Shannon's differential entropy, as

223 (2.3)
$$DE(\Pi) := -\int_{\mathbb{R}} \Pi'(x) \log(\Pi'(x)) dx$$

224 [42] show that Equation (2.3) admits a different quantile representation

225 (2.4)
$$DE(\Pi) = \int_0^1 \log(Q'_{\Pi}(p)) dp.$$

It is clear that DE is location invariant, but not scale homogeneous. It is not quantile additive either. Therefore, DE is *not* a Choquet regularizer.

3. Exploratory control with Choquet regularizers. In this section, we first introduce an exploratory stochastic control problem for RL in continuous time and spaces which was originally proposed in [47], and then reformulate it with Choquet regularizers.

Let $\mathbb{F} = \{\mathcal{F}_t\}_{t \ge 0}$ be a filtration defined on $(\Omega, \mathcal{F}, \mathbb{P})$ along with an $\{\mathcal{F}_t\}_{t \ge 0}$ adapted Brownian motion $W = \{W_t\}_{t \ge 0}$, the filtered probability space satisfying the usual assumptions of completeness and right continuity. All stochastic processes introduced below are supposed to be adapted processes in this space.

The classical stochastic control problem is to control the state dynamic described by a stochastic differential equation (SDE)

238 (3.1)
$$dX_t^u = b(X_t^u, u_t) dt + \xi(X_t^u, u_t) dW_t, \ t > 0; \quad X_0^u = x \in \mathbb{R},$$

where $u = \{u_t\}_{t \ge 0}$ is the control process taking value in a given action space U. The aim of the problem is to achieve the maximum expected total discounted reward represented by the value function

242 (3.2)
$$V^{\mathrm{cl}}(x) := \sup_{u \in \mathcal{A}^{\mathrm{cl}}(x)} \mathbb{E}_x \left[\int_0^\infty e^{-\rho t} r\left(X_t^u, u_t \right) dt \right],$$

where r is the reward function, $\rho > 0$ is the discount rate, and $\mathcal{A}^{cl}(x)$ denotes the set of all admissible controls which in general may depend on x. Throughout this paper, for ease of notation we assume that the state and Brownian motion are scalar-valued processes. Moreover, we suppose that the control is also one-dimensional, which is however an essential assumption because the Choquet regularizer to be involved is defined only for distributions on $\mathbb{R}^{.6}$

With the complete knowledge of the model parameters, the theory for solving the classical, model-based problem (3.1)–(3.2) has been developed and established

 $^{^{6}\}mathrm{See}$ Section 7 for a discussion about how we may extend the notion of Choquet regularizer to multi-dimensions.

CHOQUET REGULARIZATION

7

thoroughly. In the RL setting, where those parameters are partly or completely unknown and therefore dynamic learning is needed, the agent employs exploration to interact with and learn the unknown environment through trial and error. The key idea is to model exploration by a distribution of controls $\Pi = {\Pi_t}_{t\geq 0}$ over the control space U from which each "trial" is sampled. Thus, the notion of controls is extended to distributions. The agent executes controls for N rounds over the same time horizon, while at each round, a classical control is sampled from the distribution II. The reward of such a policy becomes accurate enough when N is large.

Thus, similarly to [47], we give the "exploratory" version of the state dynamic (3.1) motivated by repetitive learning in RL. The control process is now randomized, leading to a distributional or exploratory control process $\Pi = {\Pi_t}_{t \ge 0}$, where $\Pi_t \in$ $\mathcal{M}(U)$ is the probability distribution function for control at time t, with $\mathcal{M}(U)$ being the set of distribution functions on U. For a given such distributional control Π , the exploratory version of the state dynamics is

265 (3.3)
$$dX_t^{\Pi} = \tilde{b}\left(X_t^{\Pi}, \Pi_t\right) dt + \tilde{\xi}\left(X_t^{\Pi}, \Pi_t\right) dW_t, \ t > 0; \quad X_0^{\Pi} = x \in \mathbb{R},$$

where the coefficients $\tilde{b}(\cdot, \cdot)$ and $\tilde{\xi}(\cdot, \cdot)$ are defined as

267 (3.4)
$$\tilde{b}(y,\Pi) = \int_U b(y,u) \mathrm{d}\Pi(u), \quad y \in \mathbb{R}, \ \Pi \in \mathcal{M}(U),$$

268 and

269 (3.5)
$$\tilde{\xi}(y,\Pi) = \sqrt{\int_U \xi^2(y,u) \mathrm{d}\Pi(u)}, \quad y \in \mathbb{R}, \ \Pi \in \mathcal{M}(U).$$

The "exploratory state process" $\{X_t^{\Pi}\}_{t \ge 0}$ describes the average of the state processes under (infinitely) many different classical control processes sampled from the exploratory control $\Pi = \{\Pi_t\}_{t \ge 0}$. Further, the reward function r in (3.2) needs also to be modified to the exploratory reward

274 (3.6)
$$\tilde{r}(y,\Pi) = \int_U r(y,u) d\Pi(u), \quad y \in \mathbb{R}, \ \Pi \in \mathcal{M}(U).$$

A detailed explanation of where this exploratory formulation comes from is provided in [47, pp. 6–8]. We reiterate that the exploratory state process $\{X_t^{\Pi}\}_{t\geq 0}$ is the *average* of the sample state trajectories under infinitely many actions generated from the same distribution Π and is in itself *not* a sample state trajectory nor observable. The exploratory formulation above just provides a framework for *theoretical* analysis. See [26, p. 9] for more discussion on this point.

Next, we use a Choquet regularizer Φ_h to measure the level of exploration, and the aim of the exploratory control is to achieve the maximum expected total discounted and regularized exploratory reward represented by the optimal value function

284 (3.7)
$$V(x) = \sup_{\Pi \in \mathcal{A}(x)} \mathbb{E}_x \left[\int_0^\infty e^{-\rho t} \left(\hat{r}(X_t^{\Pi}, \Pi) + \lambda \Phi_h(\Pi) \right) \mathrm{d}t \right],$$

where $\lambda > 0$ is the *temperature* parameter representing the weight on exploration, $\mathcal{A}(x)$ is the set of admissible distributional controls (which may in general depend on x), and \mathbb{E}_x represents the conditional expectation given $X_0^{\Pi} = x$. The precise definition of $\mathcal{A}(x)$ depends on the specific dynamic model under consideration and the specific problems one wants to solve, which may vary from case to case. We will define $\mathcal{A}(x)$ precisely later for the linear-quadratic (LQ) control case, which will be the main focus of the paper. Note that (3.7) is mathematically a socalled relaxed stochastic control problem; see [47, Footnote 7] for a detailed discussion about the connection between the exploratory formulation and relaxed control.

Controls in $\mathcal{A}(x)$ are measure (distribution function)-valued stochastic adapted processes, which are open-loop controls in the control terminology. A more important notion in RL is the feedback (control) *policy*. Specifically, a deterministic mapping $\Pi(\cdot; \cdot)$ is called a feedback policy if i) $\Pi(\cdot; x)$ is a distribution function for each $x \in \mathbb{R}$; ii) the following SDE (which is the system dynamic after the feedback law $\Pi(\cdot; \cdot)$ is applied)

$$dX_t = \tilde{b}\left(X_t, \Pi\left(\cdot; X_t\right)\right) dt + \tilde{\xi}\left(X_t^{\Pi}, \Pi\left(\cdot; X_t\right)\right) dW_t, \quad t > 0; \quad X_0 = x \in \mathbb{R}$$

has a unique strong solution $\{X_t\}_{t\geq 0}$; and iii) the open-loop control $\Pi = \{\Pi_t\}_{t\geq 0} \in \mathcal{A}(x)$ where $\Pi_t := \Pi(\cdot; X_t)$. In this case, the resulting open-loop control Π is said to be generated from the feedback policy $\Pi(\cdot; \cdot)$ with respect to the initial state x. On the other hand, for a continuous $h \in \mathcal{H}$, we have $\Phi_h(\Pi) = \int_0^1 Q_{\Pi}(1-p) dh(p) = \int_{U}^1 uh'(1-\Pi(u)) d\Pi(u)$.

We present the general procedure for solving the problem (3.7), following [47]. Applying the classical Bellman principle of optimality, we deduce that the optimal value function V satisfies the Hamilton-Jacobi-Bellman (HJB) equation

(3.8)

303
$$\rho v(x) = \max_{\Pi \in \mathcal{M}(U)} \left(\tilde{r}(x,\Pi) + \lambda \int_{U} uh'(1-\Pi(u)) d\Pi(u) + \frac{1}{2} \tilde{\xi}^{2}(x,\Pi) v''(x) + \tilde{b}(x,\Pi) v'(x) \right),$$

or equivalently,

$$\rho v(x) = \max_{\Pi \in \mathcal{M}(U)} \int_{U} \left(r(x, u) + \lambda u h'(1 - \Pi(u)) + \frac{1}{2} \xi^{2}(x, u) v''(x) + b(x, u) v'(x) \right) d\Pi(u),$$

where v denotes the generic unknown solution of the equation. The verification theorem then yields that the feedback policy Π^* defined as

(3.9)

306
$$\Pi^*(x) := \underset{\Pi \in \mathcal{M}(U)}{\arg \max} \int_U \left(r(x,u) + \lambda u h'(1 - \Pi(u)) + \frac{1}{2} \xi^2(x,u) v''(x) + b(x,u) v'(x) \right) d\Pi(u)$$

is an optimal policy if it generates an admissible open-loop control for any x.

When the regularizer is the entropy, [47] applied the corresponding verification 308 theorem to conclude that the Gibbs (or Boltzmann) measures are generally optimal 309 samplers for exploration, which specialize to Gaussian in the LQ case. However, no 310 311 general study on the entropy-regularized exploratory HJB equation was available until [43] established the well-posedness and regularity of its viscosity solution. With the 312 current Choquet regularizers, studying (3.8) and solving the maximization problem 313 in (3.9) generally remain (significant) open questions because (3.8) is very different 314 from its entropy counterpart and it is unclear whether the analyses in [43, 47] carry 315 316over.

In this paper, we focus on the LQ setting, in which the exploratory HJB equation (3.8) can be explicitly solved, to study how different Choquet regularizers may generate the optimal policy distributions. Specifically, we consider

320 (3.10) b(x, u) = Ax + Bu and $\xi(x, u) = Cx + Du, x, u \in \mathbb{R}$,

where $A, B, C, D \in \mathbb{R}$, and 321

322 (3.11)
$$r(x,u) = -\left(\frac{M}{2}x^2 + Rxu + \frac{N}{2}u^2 + Px + Lu\right), \quad x,u \in \mathbb{R}$$

where $M \ge 0$, N > 0, and $R, P, L \in \mathbb{R}$. Moreover, as in standard LQ theory we 323 assume henceforth that $U = \mathbb{R}$ and thus write $\mathcal{M} = \mathcal{M}(U)$ and $\mathcal{M}^2 = \mathcal{M}^2(U)$. 324

Remark 3.1. LQ control plays a vitally important role in the classical control 325 literature, not only because it usually admits elegant and simple solutions, but also 326 because more complex, nonlinear problems can be approximated by LQ problems. 327 Indeed, one can simply apply a second-order Taylor approximation to the reward 328 function and a first-order Taylor approximation to the dynamics coefficient functions 329 to define an approximate LQ problem; see [6, 7, 28, 30, 44] and the reference therein 330 for more details. 331

Fix an initial state $x \in \mathbb{R}$. For each open-loop control $\Pi \in \mathcal{A}(x)$, denote its mean 332 and variance processes $\{\mu_t\}_{t \ge 0}$ and $\{\sigma_t^2\}_{t \ge 0}$ by $\mu_t \equiv \mu(\Pi_t) = \int_U u d\Pi_t(u)$ and $\sigma_t^2 \equiv \sigma^2(\Pi_t) = \int_U u^2 d\Pi_t(u) - \mu_t^2$. By (3.4) and (3.5), we have 333 334 (3.12)

335
$$\tilde{b}(x,\Pi) = Ax + B\mu(\Pi), \quad \tilde{\xi}(x,\Pi) = \sqrt{C^2 x^2 + 2CDx\mu(\Pi) + D^2[\mu^2(\Pi) + \sigma^2(\Pi)]}.$$

Thus, the state dynamic X^{Π} in (3.3) is given by 336

337 (3.13)
$$\mathrm{d}X_t^{\Pi} = (AX_t^{\Pi} + B\mu_t)\mathrm{d}t + \sqrt{(CX_t^{\Pi} + D\mu_t)^2 + D^2\sigma_t^2}\,\mathrm{d}W_t, \quad X_0^{\Pi} = x \in \mathbb{R},$$

which implies that the state process only depends on the mean process $\{\mu_t\}_{t\geq 0}$ and 338 the variance process $\{\sigma_t^2\}_{t\geq 0}$ of the given distributional control $\{\Pi_t\}_{t\geq 0}$. Let \mathcal{B} be 339 the Borel algebra on \mathbb{R} . A control process Π is said to be admissible, denoted by 340 $\Pi \in \mathcal{A}(x)$, if (i) for each $t \ge 0$, $\Pi_t \in \mathcal{M}$ a.s.; (ii) for each $A \in \mathcal{B}$, $\{\Pi_t(A), t \ge 0\}$ 341 is \mathcal{F}_t -progressively measurable; (iii) for each $t \ge 0$, $\mathbb{E}[\int_0^t (\mu_s^2 + \sigma_s^2) \mathrm{d}s] < \infty$; (iv) with $\{X_t^{\Pi}\}_{t\ge 0}$ solving (3.3), $\liminf_{T\to\infty} e^{-\rho T} \mathbb{E}[(X_T^{\Pi})^2] = 0$; (v) with $\{X_t^{\Pi}\}_{t\ge 0}$ solving (3.3), $\mathbb{E}[\int_0^\infty e^{-\rho t} |\tilde{r}(X_t^{\Pi}, \Pi_t) + \lambda \Phi_h(\Pi_t)| \mathrm{d}t] < \infty$. 342 343 344

In the above, condition (iii) is to ensure that for any $\Pi \in \mathcal{A}(x)$, both the drift 345 and volatility terms of (3.3) satisfy a global Lipschitz condition and a linear growth 346 condition in the state variable and, hence, the SDE (3.3) admits a unique strong solu-347 tion X^{Π} . Condition (iv) is used to ensure that dynamic programming and verification 348 theorem are applicable, as will be evident in the sequel. Finally, the reward is finite 349 under condition (v). 350

By (3.6) and (3.11), we have 351

352 (3.14)
$$\tilde{r}(x,\Pi) = -\frac{M}{2}x^2 - Rx\mu(\Pi) - \frac{N}{2}[\mu^2(\Pi) + \sigma^2(\Pi)] - Px - L\mu(\Pi).$$

Thus, plugging (3.12) and (3.14) back into (3.8), we can derive the HJB equation for 353 LQ control as 354

$$\rho v(x) = \max_{\Pi \in \mathcal{M}^2} \left\{ -Rx\mu(\Pi) - \frac{N}{2} \left[\mu^2(\Pi) + \sigma^2(\Pi) \right] - L\mu(\Pi) + \lambda \Phi_h(\Pi) + CDx\mu(\Pi)v''(x) + \frac{1}{2}D^2 \left[\mu^2(\Pi) + \sigma^2(\Pi) \right] v''(x) + B\mu(\Pi)v'(x) \right\} + Axv'(x) - \frac{M}{2}x^2 - Px + \frac{1}{2}C^2x^2v''(x).$$

To analyze and solve this equation, we need to study the maximization problem therein. Denote by $\varphi(x,\Pi)$ the term inside the max operator above. Observe that $\varphi(x,\Pi)$ depends on Π via only its mean $\mu(\Pi)$ and variance $\sigma^2(\Pi)$, except for the term $\Phi_h(\Pi)$, which motivates us to write

$$\underset{361}{\overset{360}{=}} (3.16) \qquad \max_{\Pi \in \mathcal{M}^2} \varphi(x, \Pi) = \max_{m \in \mathbb{R}, s > 0} \max_{\Pi \in \mathcal{M}^2, \mu(\Pi) = m, \sigma^2(\Pi) = s^2} \varphi(x, \Pi).$$

362 The inner maximization problem is in turn equivalent to

363 (3.17) $\max_{\Pi \in \mathcal{M}^2} \Phi_h(\Pi) \quad \text{subject to } \mu(\Pi) = m \text{ and } \sigma^2(\Pi) = s^2.$

This is a *static* optimization problem, which holds the key to solve the HJB 365 366 equation (3.15) and thus to our exploratory problem with Choquet regularizers. It is interesting to note that when the regularizer is the entropy, the optimal solution to the 367 above problem is Gaussian, which is indeed the essential reason behind the Gaussian 368 exploration derived in [47]. More specifically, for LQ control any regularized payoff 369 function depends only on the mean and variance processes of the distributional control, 370 and the Gaussian distribution maximizes the entropy when the mean and variance 371 are fixed. The natural question in our setting is what distribution with given mean 372 and variance maximizes a Choque regularizer, which is exactly the problem (3.17). 373 The next section is devoted to solving explicitly this maximization problem (3.17) of 374 "mean-variance constrained Choquet regularizers" with a variety of specific Choquet 375 regularizers.

4. Maximizing mean–variance constrained Choquet regularizers.

4.1. General results. For given $h \in \mathcal{H}$, $m \in \mathbb{R}$ and s > 0, we consider the problem (3.17), which has been motivated by the exploratory control for RL as discussed in the previous section. Note that since Φ_h is location-invariant and scalable, (3.17) is equivalent to the following problem

$$\underset{\Pi \in \mathcal{M}^2}{\overset{382}{\text{smax}}} \Phi_h(\Pi) \quad \text{subject to } \mu(\Pi) = 0 \text{ and } \sigma^2(\Pi) = 1.$$

In what follows, h' represents the right-derivative of h, which exists on [0, 1) since his concave on [0, 1]. It turns out that a general solution to (3.17) has been given by Theorem 3.1 of [31].

LEMMA 4.1. If h is continuous and not constantly zero, then a maximizer Π^* to (3.17) has the following quantile function

389 (4.1)
$$Q_{\Pi^*}(p) = m + s \frac{h'(1-p)}{||h'||_2}, \quad a.e. \ p \in (0,1),$$

390 and the maximum value of (3.17) is $\Phi_h(\Pi^*) = s ||h'||_2$.

In the context of RL, an interesting question arises: Given a distribution used for exploration, what is the regularizer that leads to that distribution? This is a practically important question that can provide interpretability to some widely used samplers for exploration in practice. Theoretically, answering this question is in some sense a converse of Lemma 4.1 at least in the LQ setting.

In what follows, we denote by $\mathcal{M}^2(m, s^2)$ the set of $\Pi \in \mathcal{M}^2$ satisfying $\mu(\Pi) = m \in \mathbb{R}$ and $\sigma^2(\Pi) = s^2 > 0$. Also, recall that given a distribution Π the *location-scale* family of Π is the set of all distributions $\Pi_{a,b}$ parameterized by $a \in \mathbb{R}$ and b > 0 such that $\Pi_{a,b}(x) = \Pi((x-a)/b)$ for all $x \in \mathbb{R}$.

PROPOSITION 4.2. Let $\Pi \in \mathcal{M}^2(m, s^2)$ be given, where $m \in \mathbb{R}$ and s > 0. Then 400 Π maximizes Φ_h as well as $\Phi_{\lambda h}$ for any $\lambda > 0$ over $\mathcal{M}^2(m, s^2)$ for a continuous $h \in \mathcal{H}$ 401 specified by 402

(1 2)

403 (4.2)
$$h'(p) = Q_{\Pi}(1-p) - m, \quad a.e. \ p \in (0,1).$$

Moreover, for any Π in the location-scale family of Π , Π also maximizes Φ_h over 404 $\mathcal{M}^2(\mu(\hat{\Pi}), \sigma^2(\hat{\Pi})).$ 405

406 *Proof.* By Lemma 4.1, given a continuous $h \in \mathcal{H}$, we have $h'(p) = ||h'||_2 (Q_{\Pi}(1 - p))$ (p) - m)/s for $p \in (0, 1)$ a.e., where Π maximizes Φ_h over $\mathcal{M}^2(m, s^2)$. Since $\Phi_{\lambda h}(\Pi) = 0$ 407 $\lambda \Phi_h(\Pi)$ for any $\lambda > 0$, Π that maximizes Φ_h also maximizes $\Phi_{\lambda h}$, which means that a 408 positive constant multiplier in Φ_h does not affect problem (3.17). Hence, Π maximizes 409 Φ_h over $\mathcal{M}^2(m, s^2)$ with $h'(p) = Q_{\Pi}(1-p) - m$ for $p \in (0, 1)$ a.e. Moreover, if Π 410 is in the location-scale family of Π , then we have $\hat{\Pi}(x) = \Pi((x-a)/b)$ for some 411 412 $a \in \mathbb{R}$ and b > 0 for all $x \in \mathbb{R}$, which implies that $h'(p) = Q_{\Pi}(1-p) - m =$ $(Q_{\Pi}(1-p)-a)/b-m$ for $p \in (0,1)$ a.e. Since $\mu(\Pi) = a+bm$, it follows that Π 413 maximizes Φ_h over $\mathcal{M}^2(\mu(\hat{\Pi}), \sigma^2(\hat{\Pi}))$. 414 Π

A simple but important implication from Proposition 4.2 is that every non-degenerate 415distribution with finite first and second moments is the optimizer of some Φ_h in (3.17) 416 over $\mathcal{M}^2(m, s^2)$ for some $m \in \mathbb{R}$ and s > 0. Therefore, any distribution used for static 417 exploration can be interpreted by certain suitable Choquet regularizer Φ_h . Moreover, 418 there is a common distortion function h, which is explicitly specified by Proposition 419 4.2, for any given location-scale family, in the sense that any distribution function Π 420 belonging to this location-scale family maximizes Φ_h over $\mathcal{M}^2(\mu(\Pi), \sigma^2(\Pi))$. In other 421 words, a single Φ_h can serve as the same regularizer for a whole location-scale family 422 of distributions. We remark that optimization of a general functional I_h may also be 423 feasible where h is not necessarily concave (see [34] for inverse S-shaped distortion 424 functions); however, this is not desirable for an exploration regularizer. 425

In the following subsections, we present specific examples applying the above 426 general results, involving several samplers commonly used in RL for exploration, as 427 well as measures commonly used in finance and operations research for evaluating 428 distribution variability. 429

4.2. Some common exploratory distributions. We first present some simple 430 distributions which have been widely used for exploration in the RL literature. 431

EXAMPLE 4.3 (Bang-bang exploration). Let Π be a Bernoulli distribution with 432 $\Pi(\{0\}) = 1 - \varepsilon \in (0,1)$ and $\Pi(\{1\}) = \varepsilon$. In this case, the RL agent explores only 433 two states 0 and 1, which is called a bang-bang exploration. In particular, in the 434 classical two-armed bandit problem, 0 is the currently more promising arm and 1 is 435the other arm. Proposition 4.2 gives $h'(p) = \mathbb{1}_{\{p < \varepsilon\}} - \varepsilon$ for $p \in (0,1)$ a.e., and thus 436 $h(p) = p \wedge \varepsilon - \varepsilon p$. The corresponding regularizer Φ_h is given by, using the quantile 437 representation in Lemma 2.3, $\Phi_h(\Pi) = \int_0^{\varepsilon} Q_{\Pi}(1-p) dp - \varepsilon \int_0^1 Q_{\Pi}(1-p) dp = \varepsilon(\mu_{\varepsilon}(\Pi) - \mu(\Pi))$, where $\mu_{\varepsilon}(\Pi)$ is the ε -tail mean defined by $\mu_{\varepsilon}(\Pi) := \frac{1}{\varepsilon} \int_0^{\varepsilon} Q_{\Pi}(1-p) dp$. Since 438 439a constant multiplier in Φ_h does not affect problem (3.17), a Bernoulli distribution 440 441 with parameter ε maximizes $\Phi_h = \mu_{\varepsilon} - \mu$. Note that the tail mean corresponds to ES in risk management with an axiomatic foundation laid out in [51]. The difference 442 between an ES and the mean, $\mu_{\varepsilon} - \mu$, is an example of generalized deviation measures 443 in Example 3 of [39], which has an axiomatic characterization similar to ES. 444

EXAMPLE 4.4 (ε -greedy exploration). Let Π be a discrete distribution satisfying 445446 $\Pi(\{0\}) = 1 - \varepsilon \in (0,1)$ and $\Pi(\{j\}) = \varepsilon/(2n)$ for $j \in \{-n, \ldots, -1, 1, \ldots, n\}$. In this 447 case, the RL agent explores 2n + 1 states where 0 is the currently most "exploita-448 tive" state and $\{-n, \ldots, -1, 1, \ldots, n\}$ represent the other states surrounding 0. From 449 Proposition 4.2, we have

$$450 \quad (4.3) \quad h'(p) = \sum_{i=1}^{n} (n-i+1) \mathbb{1}_{\left\{\frac{(i-1)\varepsilon}{2n} \leqslant p < \frac{i\varepsilon}{2n}\right\}} - \sum_{i=n+1}^{2n} (i-n) \mathbb{1}_{\left\{\frac{(i-1)\varepsilon}{2n} + 1 - \varepsilon \leqslant p < \frac{i\varepsilon}{2n} + 1 - \varepsilon\right\}}$$

451 for $p \in (0,1)$ a.e.; and thus h is a piece-wise linear function. An example of h 452 in (4.3) is plotted in FIG. 1. Using the quantile representation in Lemma 2.3, the 453 corresponding regularizer Φ_h is given by $\Phi_h(\Pi) = \varepsilon(\sum_{i=1}^n \mu_{\varepsilon}^+(i,\Pi) - \sum_{i=n+1}^{2n} \mu_{\varepsilon}^-(i,\Pi)),$ 454 where $\mu_{\varepsilon}^+(i,\Pi)$ and $\mu_{\varepsilon}^-(i,\Pi)$ are defined by

455 (4.4)
$$\mu_{\varepsilon}^{+}(i,\Pi) := \frac{n-i+1}{\varepsilon} \int_{\frac{(i-1)\varepsilon}{2n}}^{\frac{i\varepsilon}{2n}} Q_{\Pi}(1-p) \mathrm{d}p \quad for \quad i=1,\ldots,n,$$

456 and

458

457 (4.5)
$$\mu_{\varepsilon}^{-}(i,\Pi) := \frac{i-n}{\varepsilon} \int_{\frac{(i-1)\varepsilon}{2n}+(1-\varepsilon)}^{\frac{i\varepsilon}{2n}+(1-\varepsilon)} Q_{\Pi}(1-p) \mathrm{d}p \quad for \quad i=n+1,\ldots,2n.$$

This example is related to the ε -greedy strategy in multi-armed bandit problem, where



FIG. 1. The plots of h (left panel) and h' (right panel) in Example 4.4 corresponding to a discrete distribution Π where n = 5 and $\varepsilon = 0.4$.

459 ε signifies the probability of exploring. To be specific, the ε -greedy exploration is to 460 select the current best arm with probability $1 - \varepsilon$, and the other 2n arms uniformly 461 with probability $\varepsilon/(2n)$. It is worth noting that ES is also used as a criterion in the 462 multi-armed bandit problem with exploration; see [5, 9].

463 EXAMPLE 4.5 (Exponential exploration). Let Π be an exponential distribution 464 with mean 1. It follows from Proposition 4.2 that $h'(p) = -\log(p) - 1$ for $p \in$ 465 (0,1) a.e., and thus $h(p) = -p\log(p)$. The corresponding Choquet regularizer Φ_h 466 is given by $\Phi_h(\Pi) = -\int_0^1 Q_{\Pi}(1-p)(\log(p)+1)dp =: CRE(\Pi)$ for $\Pi \in \mathcal{M}$, where 467 CRE(Π) := $-\int_0^\infty \Pi([x,\infty))\log(\Pi([x,\infty)))dx$, which is called the cumulative residual 468 entropy (CRE) and studied by [24] and [38]. [45] argue that CRE can be viewed 469 as a measure of dispersion or variability. Thus, the exponential exploration can be 470 interpreted by the CRE regularizer.

471 EXAMPLE 4.6 (Gaussian exploration). If Π is a Gaussian distribution, then 472 Proposition 4.2 gives h'(p) = z(1-p) for $p \in (0,1)$ a.e., where z is the quantile 473 function of a standard normal distribution.⁷ This gives $h(p) = \int_0^p z(1-s) ds$, which 474 is plotted in FIG. 2. The corresponding regularizer Φ_h is given by

475 (4.6)
$$\Phi_h(\Pi) = \int_0^1 Q_{\Pi}(1-p)z(1-p)\mathrm{d}p = \int_0^1 Q_{\Pi}(p)z(p)\mathrm{d}p, \ \Pi \in \mathcal{M}.$$

476 Thus, any Gaussian distribution maximizes the regularize Φ_h given by $\Phi_h(\Pi) =$

477 $\int_0^1 Q_{\Pi}(p) z(p) dp$. This example also indicates that there are multiple regularizers (in-478 cluding the above regularizer and differential entropy) that induce Gaussian exploration.



FIG. 2. The plots of h (left panel) and h' (right panel) in Example 4.6 corresponding to a Gaussian distribution.

479

4.3. The inter-ES difference as a Choquet regularizer. We look at a regularizer based on ES. For $\Pi \in \mathcal{M}$, ES at level p is defined as

$$\mathrm{ES}_{p}(\Pi) := \frac{1}{1-p} \int_{p}^{1} Q_{\Pi}(r) \mathrm{d}r, \quad p \in (0,1),$$

and the left-ES is defined as

$$\mathrm{ES}_{p}^{-}(\Pi) := \frac{1}{p} \int_{0}^{p} Q_{\Pi}(r) \mathrm{d}r, \quad p \in (0, 1).$$

480 For $\alpha \in (0, 1)$, let

481 (4.7)
$$h_{\alpha}(p) := p/(1-\alpha) \wedge 1 + (\alpha - p)/(1-\alpha) \wedge 0, \ p \in [0,1].$$

Define $\Phi_{h_{\alpha}} = \text{IER}_{\alpha}$ by $\text{IER}_{\alpha}(\Pi) := \text{ES}_{\alpha}(\Pi) - \text{ES}_{1-\alpha}^{-}(\Pi)$, which is known as the inter-ES difference. Here, we assume $\alpha \in [1/2, 1)$. The inter-ES difference is a relatively new notion: it appears in Example 4 of [50] as a signed Choquet integral. In a recent work by [4], various properties are studied to underline the special role the inter-ES difference plays among other variability measures.

487 PROPOSITION 4.7. Suppose that $\alpha \in [1/2, 1)$. For $m \in \mathbb{R}$ and $s^2 > 0$, the opti-488 mization problem

$$\max_{\Pi \in \mathcal{M}^2} \operatorname{IER}_{\alpha}(\Pi) \quad subject \ to \ \mu(\Pi) = m \ and \ \sigma^2(\Pi) = s^2$$

⁷In statistics, the quantile of a standard normal distribution corresponding to a test statistic is often referred to as a z-score – hence the notation z.

491 is solved by a three-point distribution Π^* with its quantile function uniquely specified 492 as

493 (4.8)
$$Q_{\Pi^*}(p) = m + \frac{s}{\sqrt{2(1-\alpha)}} \left[\mathbb{1}_{\{p > \alpha\}} - \mathbb{1}_{\{p \le 1-\alpha\}} \right], \quad a.e. \ p \in (0,1).$$

494 Proof. Note that for $\Phi_h = \text{IER}_{\alpha}$, we have $h'(p) = \frac{1}{1-\alpha} \mathbb{1}_{\{p < 1-\alpha\}} - \frac{1}{1-\alpha} \mathbb{1}_{\{p \ge \alpha\}}$ for 495 $\alpha \in [1/2, 1)$, By (4.1), we can show that a maximizer Π^* satisfies (4.8), which is a 496 three-point distribution.

497 So the inter-ES difference regularizer encourages exploration at three points. One 498 of them is the mean m corresponding to the best single-point exploitation without 499 exploration, while the other two spots are symmetric to m capturing the exploration 500 part.

Remark 4.8. For $\alpha \in [1/2, 1)$, if we take the function $h_{\alpha}(p) = \mathbb{1}_{[1-\alpha,\alpha]}(p), p \in$ 501[0,1], the inter-quantile difference $\Phi_{h_{\alpha}} := IQR_{\alpha}$ is given by $IQR_{\alpha}(\Pi) := Q_{\Pi}^{+}(\alpha) - Q_{\Pi}^{+}(\alpha)$ 502 $Q_{\Pi}(1-\alpha)$, which is a classical measure of statistical dispersion widely used in e.g., 503504box plots. Unlike the inter-ES difference, the distortion function h_{α} for IQR_{α} is not concave. However, the concave envelopes of h is give by $h^*(p) = p/(1-\alpha) \wedge 1 + (\alpha - \alpha)$ 505 $p/(1-\alpha) \wedge 0, p \in [0,1]$, which is exactly (4.7). According to Theorem 1 in [34], we 506have $\sup_{\Pi \in \mathcal{M}^2} IQR_{\alpha}(\Pi) = \sup_{\Pi \in \mathcal{M}^2} IER_{\alpha}(\Pi)$ and the maximizer is obtained by Π^* 507 which satisfies (4.8). Thus, the optimization problem is still solvable even if h is not 508 509 concave.

4.4. The L^1 -Wasserstein distance to Dirac measures as a Choquet reg-510**ularizer.** Let $W : \mathcal{M} \times \mathcal{M} \to \mathbb{R}_+$ be a statistical distance between two distributions, 511such as a Wasserstein distance. Since an exploration is essentially to move away from 512 degenerate (Dirac) distributions, a natural way to encourage exploration is to use 513 $W(\Pi, \delta_x)$, where δ_x is the Dirac measure at $x \in \mathbb{R}$, as a regularizer. Moreover, to re-514move the location dependence, we modify the regularizer to be $\min_{x \in \mathbb{R}} W(\Pi, \delta_x)$. For 515any statistical distance satisfying $W(\Pi, \hat{\Pi}) = 0$ if and only if $\Pi = \hat{\Pi}$, it is clear that 516517 $\min_{x \in \mathbb{R}} W(\Pi, \delta_x) = 0$ if and only if Π itself is a Dirac measure (hence deterministic). The use of Wasserstein distance to model distributional uncertainty in other set-518

tings naturally gives rise to a regularization term, yielding a theoretical justification for its use in practice; see for example [8, 14, 35] that formulate different models with distributional robustness based on Wasserstein distances.

We focus on the case where W is the Wasserstein L^1 distance, defined as

$$W_1(\Pi, \hat{\Pi}) := \int_0^1 |Q_\Pi(p) - Q_{\hat{\Pi}}(p)| \mathrm{d}p.$$

In this case, $W_1(\Pi, \delta_x)$ is the L^1 distance between x and $X \sim \Pi$, and it is well known via L^1 loss minimization that the minimizers of $\min_{x \in \mathbb{R}} W_1(\Pi, \delta_x)$ are the medians of Π (unique if Q_{Π} is continuous) given as $\arg \min_{x \in \mathbb{R}} W_1(\Pi, \delta_x) = [Q_{\Pi}(1/2), Q_{\Pi}^+(1/2)]$. Moreover, for a median of Π , $x^* \in [Q_{\Pi}(1/2), Q_{\Pi}^+(1/2)]$, we have that $W_1(\Pi, \delta_{x^*})$ is the mean-median deviation; namely

527
$$\min_{x \in \mathbb{R}} W_1(\Pi, \delta_x) = W_1(\Pi, \delta_{x^*}) = \int_0^{1/2} (x^* - Q_{\Pi}(p)) dp + \int_{1/2}^1 (Q_{\Pi}(p) - x^*) dp$$

528
529
$$= \int_{1/2} Q_{\Pi}(p) dp - \int_{0}^{+} Q_{\Pi}(p) dp.$$

This manuscript is for review purposes only.

This in turn shows that $\arg \min_{x \in \mathbb{R}} W_1(\Pi, \delta_x)$ belongs to the class of Choquet regularizers.

PROPOSITION 4.9. For $m \in \mathbb{R}$ and $s^2 > 0$, the optimization problem

$$\max_{\Pi \in \mathcal{M}^2} \min_{x \in \mathbb{R}} W_1(\Pi, \delta_x) \qquad subject \ to \ \mu(\Pi) = m \ and \ \sigma^2(\Pi) = s^2,$$

is solved by a unique Π^* with the quantile function specified as

533 (4.9)
$$Q_{\Pi^*}(p) = m + s \mathbb{1}_{\{p > 1/2\}} - s \mathbb{1}_{\{p \le 1/2\}}, \quad a.e. \ p \in (0, 1).$$

534 Proof. Applying Lemma 2.3 to get $\min_{x \in \mathbb{R}} W_1(\Pi, \delta_x) = \Phi_h(\Pi)$ with h'(p) = 1 for 535 p < 1/2 and h'(p) = -1 for $p \ge 1/2$. Using (4.1) in Lemma 4.1 yields (4.9), which 536 implies a symmetric two-point distribution.

As $\Phi_h(\Pi) = \min_{x \in \mathbb{R}} W_1(\Pi, \delta_x)$ induces a symmetric exploration around the mean, we call it a symmetric Wasserstein regularizer with $h(p) = p \mathbb{1}_{\{p < 1/2\}} + (1-p) \mathbb{1}_{\{p \ge 1/2\}}$. Next, let us discuss two-point asymmetric exploration. Suppose that two directions are not symmetric, and we would like to regularize in a way to encourage more exploration in a certain direction. Take a constant $\alpha \in (0, 1)$, and choose W as an asymmetric Wasserstein distance

$$W_1^{\alpha}(\Pi, \hat{\Pi}) = \int_0^1 \left(\alpha (Q_{\Pi}(p) - Q_{\hat{\Pi}}(p))_+ + (1 - \alpha) (Q_{\Pi}(p) - Q_{\hat{\Pi}}(p))_- \right) \mathrm{d}p.$$

537 The minimizers are the α -quantiles $\arg\min_{x\in\mathbb{R}} W_1^{\alpha}(\Pi, \delta_x) = [Q_{\Pi}(\alpha), Q_{\Pi}^+(\alpha)]$, and for 538 $x^* \in [Q_{\Pi}(\alpha), Q_{\Pi}(\alpha)]$, we have

539
$$\min_{x \in \mathbb{R}} W_1^{\alpha}(\Pi, \delta_x) = W_1^{\alpha}(\Pi, \delta_{x^*}) = \int_0^{\alpha} (1 - \alpha)(x^* - Q_{\Pi}(p)) dp + \int_{\alpha}^1 \alpha (Q_{\Pi}(p) - x^*) dp$$

540
541
$$= \alpha \int_{\alpha}^{1} Q_{\Pi}(p) \mathrm{d}p - (1-\alpha) \int_{0}^{\alpha} Q_{\Pi}(p) \mathrm{d}p.$$

- 542 We call $\Phi_h(\Pi) = \min_{x \in \mathbb{R}} W_1^{\alpha}(\Pi, \delta_x)$ an asymmetric Wasserstein regularizer with 543 $h(p) = \alpha p \mathbb{1}_{\{p < 1-\alpha\}} + (1-\alpha)(1-p) \mathbb{1}_{\{p \ge 1-\alpha\}}.$
- 544 PROPOSITION 4.10. For $m \in \mathbb{R}$ and $s^2 > 0$, the optimization problem

545
$$\max_{\Pi \in \mathcal{M}^2} \min_{x \in \mathbb{R}} W_1^{\alpha}(\Pi, \delta_x) \quad subject \ to \ \mu(\Pi) = m \ and \ \sigma^2(\Pi) = s^2$$

547 has a unique maximizer Π^* with the quantile function uniquely specified as

548 (4.10)
$$Q_{\Pi^*}(p) = m + s \left(\frac{\alpha}{1-\alpha}\right)^{1/2} \mathbb{1}_{\{p > \alpha\}} - s \left(\frac{1-\alpha}{\alpha}\right)^{1/2} \mathbb{1}_{\{p \leqslant \alpha\}}, \quad a.e. \ p \in (0,1).$$

549 Proof. For $\Phi_h(\Pi) = \min_{x \in \mathbb{R}} W_1^{\alpha}(\Pi, \delta_x)$, we have $h'(p) = \alpha$ for $p < 1 - \alpha$, and 550 $h'(p) = -1 + \alpha$ for $p \ge 1 - \alpha$. Using (4.1), the optimization problem has a solution 551 Π^* satisfying (4.10), which is an asymmetric two-point distribution. \square

To recap, the Wasserstein L^1 regularization encourages possibly asymmetric (with respect to the mean) two-point exploration, which is an instance of the bang-bang exploration in Example 4.3. 4.5. The Gini mean difference or maxiance as a Choquet regularizer. By letting $h(p) = p - p^2$, $p \in [0, 1]$, we consider the regularizer $\Phi_{\sigma} := \Phi_h$ given by

$$\Phi_{\sigma}(\Pi) = \int_{\mathbb{R}} \left(\Pi([x,\infty)) - \Pi^2([x,\infty)) \right) dx$$

There are two ways to represent $\Phi_{\sigma}(\Pi)$ in terms of two iid copies X_1 and X_2 from the distribution Π . First, Φ_{σ} can be rewritten as $\Phi_{\sigma}(\Pi) = \frac{1}{2}\mathbb{E}[|X_1 - X_2|]$, which is the *Gini mean difference* (e.g., [17]; sometimes without the factor 1/2). Alternatively, Φ_{σ} can be represented as $\Phi_{\sigma}(\Pi) = \mathbb{E}[\max\{X_1, X_2\}] - \mu(\Pi)$, which is called the *maxiance* by [13]. The two representations are identical as seen from the following equality

560
$$\mathbb{E}[\max\{X_1, X_2\}] - \mu(\Pi) = \mathbb{E}\left[\max\{X_1, X_2\} - \frac{1}{2}(X_1 + X_2)\right]$$

561
$$= \mathbb{E}\left[\max\{X_1, X_2\} - \frac{1}{2}(\max\{X_1, X_2\} + \min\{X_1, X_2\})\right]$$

$$= \frac{1}{2} \mathbb{E} \left[\max\{X_1, X_2\} - \min\{X_1, X_2\} \right] = \frac{1}{2} \mathbb{E} \left[|X_1 - X_2| \right]$$

As argued by [13], the maxiance can be seen as the dual version of the variance, due to the identities $\sigma^2(\Pi) = \int_{\mathbb{R}} (x - \mu(\Pi))^2 d\Pi$ and $\Phi_{\sigma}(\Pi) = \int_{\mathbb{R}} (x - \mu(\Pi)) d\Pi^2$. Moreover, the maxiance can be used to approximate a local index of absolute risk aversion in [53]'s dual theory of choice under risk, which is similar to the role of variance in the classic expected utility theory.

We now show that the maxiance regularizer Φ_{σ} leads to a uniform distribution for exploration.

571 PROPOSITION 4.11. For $m \in \mathbb{R}$ and $s^2 > 0$, the optimization problem

572 (4.11)
$$\max_{\Pi \in \mathcal{M}^2} \Phi_{\sigma}(\Pi) \quad subject \ to \ \mu(\Pi) = m \ and \ \sigma^2(\Pi) = s^2$$

has a unique maximizer $\Pi^* = U[m - \sqrt{3}s, m + \sqrt{3}s].$

575 Proof. Note that for $\Phi_h = \Phi_\sigma$, we have h'(p) = 1 - 2p. It follows from (4.1) that 576 a maximizer Π^* is a uniform distribution. By matching the moments in (4.11), we 577 obtain $\Pi^* = U[m - \sqrt{3}s, m + \sqrt{3}s]$. The uniqueness statement is guaranteed by e.g. 578 Theorem 2 of [34].

Proposition 4.11 provides a foundation for a uniformly distributed exploration strategy on \mathbb{R} . Note that this is different from the result of uniform distributions maximizing entropy on a fixed, given bounded region: here in our setting the region is *not* fixed, since we allow II to be chosen from arbitrary distributions on \mathbb{R} , and thus the bounded region $[m - \sqrt{3}s, m + \sqrt{3}s]$ is endogenously derived rather than exogenously given.

Remark 4.12. The inequality $\sigma(\Pi) \ge \sqrt{3}\Phi_{\sigma}(\Pi)$ for all $\Pi \in \mathcal{M}^2$ is known as Glasser's inequality ([20]). For the uniform distribution Π^* in Proposition 4.11 with $\sigma(\Pi^*) = s$, we have $\Phi_{\sigma}(\Pi^*) = \sqrt{3}s/3$ by Lemma 4.1. Thus, Π^* attains the sharp bound of Glasser's inequality, which holds naturally since Π^* maximizes Φ_{σ} for a fixed σ^2 .

590 **5.** Solving the exploratory stochastic LQ control problem. We are now 591 ready to solve the exploratory stochastic LQ control problem presented in Section 3. 592 Let

593 (5.1)
$$W(x,\Pi) = \mathbb{E}_x \left[\int_0^\infty e^{-\rho t} \left(\tilde{r}(X_t^\Pi,\Pi_t) + \lambda \Phi_h(\Pi_t) \right) \mathrm{d}t \right], \ x \in \mathbb{R}, \ \Pi \in \mathcal{A}(x).$$

594 We have the following result based on Lemma 4.1.

595 PROPOSITION 5.1. Let a continuous $h \in \mathcal{H}$ be given. For any $\Pi = {\Pi_t}_{t\geq 0} \in \mathcal{A}(x)$ with mean process ${\mu_t}_{t\geq 0}$ and variance process ${\sigma_t^2}_{t\geq 0}$, there exists $\Pi^* = {\Pi_t^*}_{t\geq 0} \in \mathcal{A}(x)$ given by

598 (5.2)
$$Q_{\Pi_t^*}(p) = \mu_t + \sigma_t \frac{h'(1-p)}{||h'||_2}, \quad a.e. \ p \in (0,1), \ t \ge 0,$$

599 which has the same mean and variance processes satisfying $W(x, \Pi^*) \ge W(x, \Pi)$.

600 Proof. It follows from (3.13) and (3.14) that the term $\mathbb{E}_x \left[\int_0^\infty e^{-\rho t} \tilde{r}(X_t^{\Pi}, \Pi_t) dt \right]$ 601 in (5.1) only depends on the mean process $\{\mu_t\}_{t\geq 0}$ and variance process $\{\sigma_t^2\}_{t\geq 0}$ 602 of $\{\Pi_t\}_{t\geq 0}$. Thus, for any fixed $t \geq 0$, choose Π_t^* with mean μ_t and variance σ_t^2 603 that maximizes $\Phi_h(\Pi)$. From Lemma 4.1, it follows that Π_t^* satisfies (5.2) and the 604 maximum value is $\Phi_h(\Pi_t) = \sigma_t ||h'||_2$. Clearly, the strategy $\Pi^* = \{\Pi_t^*\}_{t\geq 0} \in \mathcal{A}(x)$ is 605 the desired one.

Proposition 5.1 indicates that the control problem (3.7) in the LQ setting is maximized within a location–scale family of distributions, which is determined only by h.

We go back to the HJB equation (3.15). It follows from (3.16)–(3.17) along with Lemma 4.1 that (3.15) is equivalent to (5.3)

611

$$\rho v(x) = \max_{\mu \in \mathbb{R}, \sigma > 0} \left[-Rx\mu - \frac{N}{2} \left(\mu^2 + \sigma^2 \right) - L\mu + \lambda \sigma \|h'\|_2 + CDx\mu v''(x) + \frac{1}{2}D^2 \left(\mu^2 + \sigma^2 \right) v''(x) + B\mu v'(x) \right] + Axv'(x) - \frac{M}{2}x^2 - Px + \frac{1}{2}C^2x^2v''(x)$$

By the first-order conditions, we get the maximizers $\mu^*(x) = \frac{CDxv''(x) + Bv'(x) - Rx - L}{N - D^2v''(x)}$

and $(\sigma^*(x))^2 = \frac{\lambda^2 \|h'\|_2^2}{(N-D^2v''(x))^2}$ of the max operator in (5.3), which in turn leads to the optimal distributional policy $\Pi^*(\cdot; x)$ prescribed by Lemma 4.1.

Bringing the above expressions of $\mu^*(x)$ and $\sigma^*(x)$ back into (5.3), we can further write the HJB equation as

617 (5.4)

$$\rho v(x) = \frac{[CDxv''(x) + Bv'(x) - Rx - L]^2 + \lambda^2 \|h'\|_2^2}{2[N - D^2v''(x)]} + \frac{1}{2} [C^2v''(x) - M] x^2 + [Av'(x) - P]x.$$

618 We now solve this equation explicitly. Denote $\Delta = [\rho - (2A + C^2)]N + 2(B + CD)R - D^2M$. Under the assumptions that $\rho > 2A + C^2$ and $MN > R^2$, a smooth solution 620 to (5.4) is given by $v(x) = \frac{1}{2}k_2x^2 + k_1x + k_0$, where⁸

621 (5.5)
$$k_2 = \frac{\Delta - \sqrt{\Delta^2 - 4[(B + CD)^2 + (\rho - (2A + C^2))D^2](R^2 - MN)}}{2[(B + CD)^2 + D^2(\rho - (2A + C^2))]}$$

⁸Values of k_2 , k_1 and k_0 are obtained by solving the system of equations $\rho k_2 = \frac{(k_2(B+CD)-R)^2}{N-k_2D^2} + k_2 \left(2A+C^2\right) - M$, $\rho k_1 = \frac{(k_1B-L)(k_2(B+CD)-R)}{N-k_2D^2} + k_1A - P$, and $\rho k_0 = \frac{(k_1B-L)^2 + \lambda^2 ||h'||_2^2}{2(N-k_2D^2)}$.

This manuscript is for review purposes only.

622

623 (5.6)
$$k_1 = \frac{P(N - k_2 D^2) - LR}{k_2 B(B + CD) + (A - \rho)(N - k_2 D^2) - BR},$$

624 and

625 (5.7)
$$k_0 = \frac{(k_1 B - L)^2 + \lambda^2 ||h'||_2^2}{2\rho(N - D^2 k_2)}.$$

We can verify easily that $k_2 < 0$. Hence, v is concave, a property that is essential for v to be actually the value function. Next, we state the main result of this section, whose proof follows essentially the same lines of that of Theorem 4 in [47], thanks to the analysis above and the results obtained. We omit the details here.

THEOREM 5.2. Consider the LQ control specified by (3.10)-(3.11), where we assume $M \ge 0$, N > 0, $MN > R^2$ and $\rho > 2A + C^2 + \max\left(\frac{D^2R^2 - 2NR(B+CD)}{N}, 0\right)$. Then the value function in (3.7) is given by $V(x) = \frac{1}{2}k_2x^2 + k_1x + k_0$ for each $x \in \mathbb{R}$, where k_2 , k_1 and k_0 are as in (5.5)-(5.7), respectively. The optimal feedback policy has the distribution function $\Pi^*(\cdot; x)$ whose quantile function is (5.8)

635
$$Q_{\Pi^*(\cdot;x)}(p) = \frac{(k_2(B+CD)-R)x + k_1B - L}{N - k_2D^2} + \frac{\lambda h'(1-p)}{N - k_2D^2}, \quad a.e. \ p \in (0,1), \quad x \in \mathbb{R},$$

636 with the mean and variance given by (5.9)

637
$$\mu^*(x) = \frac{(k_2(B+CD)-R)x+k_1B-L}{N-k_2D^2}$$
 and $(\sigma^*(x))^2 = \frac{\lambda^2 \|h'\|_2^2}{(N-k_2D^2)^2}, \quad x \in \mathbb{R}.$

Finally, the associated optimal state process $\{X_t^*\}_{t\geq 0}$ with $X_0^* = x$ under $\Pi^*(\cdot; \cdot)$ is the unique solution of the SDE

$$\begin{aligned} & 640 \qquad \mathrm{d}X_t^* = \left[\left(A + \frac{B\left(k_2(B+CD)-R\right)}{N-k_2D^2} \right) X_t^* + \frac{B\left(k_1B-L\right)}{N-k_2D^2} \right] \mathrm{d}t \\ & 641 \qquad + \sqrt{\left[\left(C + \frac{D\left(k_2(B+CD)-R\right)}{N-k_2D^2} \right) X_t^* + \frac{D\left(k_1B-L\right)}{N-k_2D^2} \right]^2 + \frac{D^2\lambda^2 \left\| h' \right\|_2^2}{(N-k_2D^2)^2} \mathrm{d}W_t. \end{aligned}$$

Some remarks are in order. First of all, (5.8) implies that for any Choquet regu-643 larizer, the optimal exploratory distribution in the regularized LQ problem is uniquely 644 determined by h'. Note that h'(x) is the "probability weight" put on x when calcu-645 lating the (nonlinear) Choquet expectation; see e.g. [19, 37]. Second, we can see from 646 (5.9) that the mean of the optimal distribution does not depend on the exploration 647 represented by h and λ , and only the variance does. In particular, the mean is exactly 648 649 the same as the one in [47] when the differential entropy is used as a regularizer, which is also identical to the optimal control of the classical, non-exploratory LQ problem. 650 Third, the mean of the exploration distributions is a linear function of the state, while 651 its variance is independent of the state. 652

These observations are intuitive in the context of RL. Different h's correspond to different Choquet regularizers; hence they will certainly affect the way and the

⁹The constraint on ρ is used not only to ensure $k_2 < 0$ but also to show $\liminf_{T \to \infty} e^{-\rho T} \mathbb{E}[(X_T^T)^2] = 0$; see the proof of Theorem 4 in [47] for more details.

level of exploration. Also, the more weight put on the level of exploration, the more 655 656 spreaded out the exploration becomes around the current position. Furthermore, the second and third observations above show a perfect separation between exploitation 657 and exploration, as the former is captured by the mean and the latter by the variance 658 of the optimal exploration distributions. This property is also consistent with the LQ 659 case studied in [47, 48] even though a different type of regularizer is applied therein. 660

Next, we investigate optimal exploration samplers under the LQ framework for some concrete choices of h studied in Section 4. For convenience, we denote

$$\tilde{\sigma}^*(x) := \frac{\sigma^*(x)}{\|h'\|_2} \equiv \frac{\lambda}{N - k_2 D^2}$$

Theorem 5.2 yields that the mean of the optimal distribution is independent of h; 661 so we will specify only its quantile function and variance for each h discussed below. 662 663

Recall that the expressions of $\mu^*(x)$ and $(\sigma^*(x))^2$ for a general h are given by (5.9). (i) Let $h(p) = (p \wedge \varepsilon - \varepsilon p)$, leading to $\Phi_h(\Pi) = \varepsilon(\mu_{\varepsilon}(\Pi) - \mu(\Pi))$; see Example 4.3. The optimal policy is ε -greedy, given as

$$\Pi^* \left(\{ \mu^*(x) + (1-\varepsilon)\tilde{\sigma}^*(x) \} \right) \equiv \Pi^* \left(\left\{ \frac{(k_2(B+CD)-R)x + k_1B - L + (1-\varepsilon)\lambda}{N - k_2D^2} \right\} \right)$$
$$= \varepsilon,$$

and

$$\Pi^*\left(\{\mu^*(x) - \varepsilon \tilde{\sigma}^*(x)\}\right) \equiv \Pi^*\left(\left\{\frac{(k_2(B+CD) - R)x + k_1B - L - \varepsilon \lambda}{N - k_2D^2}\right\}\right) = 1 - \varepsilon.$$

664 At each state x, the control policy takes a more "promising" action at $\mu^*(x) - \varepsilon \tilde{\sigma}^*(x)$

with a large probability $1 - \varepsilon$, and tries an alternative action $\mu^*(x) + (1 - \varepsilon)\tilde{\sigma}^*(x)$ with 665 666

probability ε .¹⁰ Since $||h'||_2^2 = \varepsilon(1-\varepsilon)$, the variance of Π^* is $(\sigma^*(x))^2 = \frac{\varepsilon(1-\varepsilon)\lambda^2}{(N-k_2D^2)^2}$. (ii) Let h(p) be specified by the discrete exploration in (4.3), leading to $\Phi_h(\Pi) = \varepsilon(\sum_{i=1}^n \mu_{\varepsilon}^+(i,\Pi) - \sum_{i=n+1}^{2n} \mu_{\varepsilon}^-(i,\Pi))$, where $\mu_{\varepsilon}^+(i,\Pi)$ and $\mu_{\varepsilon}^-(i,\Pi)$ are defined by (4.4) and (4.5); see Example 4.4. The optimal policy is a (2n+1)-point distribution given as

$$\Pi^* \left(\{ \mu^*(x) + j \tilde{\sigma}^*(x) \} \right) \equiv \Pi^* \left(\left\{ \frac{(k_2(B + CD) - R)x + k_1B - L + j\lambda}{N - k_2D^2} \right\} \right) = \frac{\varepsilon}{2n},$$

for $j \in \{-n, \dots, -1, 1, \dots, n\}$, and

$$\Pi^* \left(\{ \mu^*(x) \} \right) \equiv \Pi^* \left(\left\{ \frac{(k_2(B+CD)-R)x + k_1B - L}{N - k_2D^2} \right\} \right) = 1 - \varepsilon.$$

Similarly, at each state x, the control policy takes a more "exploitative" action at $\mu^*(x)$ 667

with a large probability $1 - \varepsilon$, and tries 2n alternative actions $\mu^*(x) + j\tilde{\sigma}^*(x)$ for $j \in$ 668

 $\{-n, \ldots, -1, 1, \ldots, n\}$, each with probability $\varepsilon/(2n)$. Since $\|h'\|_2^2 = \varepsilon(n+1)(2n+1)/6$, the variance of Π^* is given by $(\sigma^*(x))^2 = \frac{\varepsilon(n+1)(2n+1)\lambda^2}{6(N-k_2D^2)^2}$. 669

¹⁰Precisely speaking, the policy presented here is not exactly the ε -greedy strategy in the classical two-arm bandit problem because the two "arms" in our setting depend on the current state x and hence are dynamically changing over time. However, at any point of time one needs to explore only two action points.

(iii) Let $h(p) = -p \log(p)$, leading to $\Phi_h(\Pi) = \int_0^\infty \Pi([x,\infty)) \log(\Pi([x,\infty))) dx$; see Example 4.5. The optimal policy is a shifted-exponential distribution given as

$$\Pi^*(u;x) = 1 - \exp\left\{\frac{[(k_2(B+CD) - R)x + k_1B - L]}{\lambda} - 1\right\} \exp\left\{-\frac{(N - D^2k_2)u}{\lambda}\right\}.$$

Since $||h'||_2^2 = 1$, the variance of Π^* is given by $(\sigma^*(x))^2 = \frac{\lambda^2}{(N-k_2D^2)^2}$. 671

(iv) Let $h(p) = \int_0^p z(1-s) ds$ where z is the standard normal quantile function. We have $\Phi_h(\Pi) = \int_0^1 Q_{\Pi}(p) z(p) dp$; see Example 4.6. The optimal policy is a normal distribution given by

$$\Pi^*(\cdot; x) = \mathcal{N}\left(\frac{(k_2(B+CD)-R)x + k_1B - L}{N - k_2D^2}, \frac{\lambda^2}{(N - k_2D^2)^2}\right),$$

owing to the fact that $||h'||_2^2 = 1$. Recall that the optimal distribution is also Gaussian 672 in [47] using the entropy regularizer. This is an example of different regularizers 673 leading to the same class of exploration samplers. On the other hand, examining more 674 closely the Gaussian policy derived above and the one in [47, eq. (40)], we observe 675 676 that the means of the two are identical but the variance of the former is the square of that of the latter. The reason of the discrepency in variance is because the maximized 677 mean-variance constrained Choquet regularizer $\Phi_h(\Pi)$ is always linear in the given 678 standard deviation σ whereas the corresponding maximized entropy regularizer DE(Π) 679 is logorithmic in σ . 680

(v) Let $h(p) = p/(1-\alpha) \wedge 1 + (\alpha - p)/(1-\alpha) \wedge 0$ with $\alpha \in [1/2, 1)$. Then $\Phi_h(\Pi) = \mathrm{ES}_{\alpha}(\Pi) - \mathrm{ES}_{1-\alpha}^{-}(\Pi)$; see Section 4.3. The optimal policy is a three-point distribution given as

$$\Pi^* \left(\left\{ \frac{(1-\alpha)[(k_2(B+CD)-R)x+k_1B-L]+\lambda}{(1-\alpha)(N-k_2D^2)} \right\} \right) = 1-\alpha,$$
$$\Pi^* \left(\left\{ \frac{(k_2(B+CD)-R)x+k_1B-L}{N-k_2D^2} \right\} \right) = 2\alpha - 1,$$

and

$$\Pi^*\left(\left\{\frac{(1-\alpha)[(k_2(B+CD)-R)x+k_1B-L]-\lambda}{(1-\alpha)(N-k_2D^2)}\right\}\right) = 1-\alpha.$$

681

Since $||h'||_2^2 = 2a/(1-\alpha)^2$, the variance of Π^* is given by $(\sigma^*(x))^2 = \frac{2\alpha\lambda^2}{(1-\alpha)^2(N-k_2D^2)^2}$. (vi) Let $h(p) = \alpha p \mathbb{1}_{\{p < 1-\alpha\}} + (1-\alpha)(1-p)\mathbb{1}_{\{p \ge 1-\alpha\}}$ with $\alpha \in (0,1)$. Then $\Phi_h(\Pi) = \min_{x \in \mathbb{R}} W_1(\Pi, \delta_x)$; see Section 4.4. The optimal feedback policy is an asymmetric two-point distribution given as

$$\Pi^*\left(\left\{\frac{(k_2(B+CD)-R)x+k_1B-L+\alpha\lambda}{N-k_2D^2}\right\}\right)=1-\alpha,$$

and

$$\Pi^*\left(\left\{\frac{(k_2(B+CD)-R)x+k_1B-L-(1-\alpha)\lambda}{N-k_2D^2}\right\}\right)=\alpha.$$

Since $||h'||_2^2 = \alpha(1-\alpha)$, the variance of Π^* is given by $(\sigma^*(x))^2 = \frac{\alpha(1-\alpha)\lambda^2}{(N-k_2D^2)^2}$. 682

(vii) Let $h(p) = p - p^2$. Then $\Phi_h(\Pi) = \mathbb{E}[|X_1 - X_2|]/2$; see Section 4.5. The optimal policy $\Pi^*(\cdot; x)$ is a uniform distribution given as

$$U\left[\frac{(k_2(B+CD)-R)x+k_1B-L-\lambda}{N-k_2D^2},\frac{(k_2(B+CD)-R)x+k_1B-L+\lambda}{N-k_2D^2}\right]$$

This manuscript is for review purposes only.

CHOQUET REGULARIZATION

Since $||h'||_2^2 = 1/3$, the variance of Π^* is given by $(\sigma^*(x))^2 = \frac{\lambda^2}{3(N-k_2D^2)^2}$. 683 Note here the uniform distribution is on a state-dependent bounded region cen-684 tering around the mean $\mu^*(x)$, rather than on a pre-specified bounded region. 685

6. Relationship between classical and exploratory problems. In this sec-686 tion, similarly to the discussions in [47, 48], we study the relationship between the 687 classical (unregularized and non-exploratory) and exploratory stochastic LQ prob-688 lems. Since most results are parallel, we will make the exposition brief. 689

Recall the classical LQ problem (3.2) where the reward function is given by (3.11). 690 The explicit forms of optimal control and value function, denoted respectively by 691 u^* and V^{cl} , were given by Theorem 9-(b) of [47]. We now provide the solvability 692 equivalence between the problems (3.2) and (3.7). 693

694

THEOREM 6.1. The following two statements (a) and (b) are equivalent. (a) The function $V(x) = \frac{1}{2}\alpha_2 x^2 + \alpha_1 x + \alpha_0 + \frac{\lambda^2 \|h'\|_2^2}{2\rho(N-\alpha_2 D^2)}, x \in \mathbb{R}$, with $\alpha_0, \alpha_1 \in \mathbb{R}$ and $\alpha_2 < 0$, is the value function of the exploratory problem (3.7) and the 695 696 corresponding optimal feedback policy has the distribution function $\Pi^*(\cdot;x)$ whose quantile function is $Q_{\Pi^*(\cdot;x)}(p) = \frac{(\alpha_2(B+CD)-R)x+\alpha_1B-L}{N-\alpha_2D^2} + \frac{\lambda h'(1-p)}{N-\alpha_2D^2}$ with the mean and variance given by $\mu^*(x) = \frac{(\alpha_2(B+CD)-R)x+\alpha_1B-L}{\alpha_2D^2}$ and 697 698 699

700
$$(\sigma^*(x))^2 = \frac{\lambda \|h\|_2}{(N - \alpha_2 D^2)^2}.$$

7

(b) The function $w(x) = \frac{1}{2}\alpha_2 x^2 + \alpha_1 x + \alpha_0, x \in \mathbb{R}$, with $\alpha_0, \alpha_1 \in \mathbb{R}$ and $\alpha_2 < 0$, is 701 the value function of the classical problem (3.2) and the corresponding optimal 702 feedback control is $u^*(x) = \frac{(\alpha_2(B+CD)-R)x + \alpha_1B-L}{N}$ 703

Proof. We rewrite the exploratory dynamics of X^* under Π^* as 704

$$dX_{t}^{*} = \left(AX_{t}^{*} + B\frac{\left(\alpha_{2}(B+CD)-R\right)X_{t}^{*} + \alpha_{1}B - L}{N-\alpha_{2}D^{2}}\right)dt$$

$$+ \sqrt{\left(CX_{t}^{*} + D\frac{\left(\alpha_{2}(B+CD)-R\right)X_{t}^{*} + \alpha_{1}B - L}{N-\alpha_{2}D^{2}}\right)^{2} + \frac{D^{2}\lambda^{2} \|h'\|_{2}^{2}}{(N-\alpha_{2}D^{2})^{2}}}dW_{t}$$

$$\equiv (A_{1}X_{t}^{*} + A_{2})dt + \sqrt{(B_{1}X_{t}^{*} + B_{2})^{2} + C_{1}}dW_{t},$$

706

where $A_1 := A + \frac{B(\alpha_2(B+CD)-R)}{N-\alpha_2D^2}, A_2 := \frac{B(\alpha_1B-L)}{N-\alpha_2D^2}, B_1 := C + \frac{D(\alpha_2(B+CD)-R)}{N-\alpha_2D^2}, B_2 := \frac{D(\alpha_1B-L)}{N-\alpha_2D^2}$ and $C_1 := \frac{D^2\lambda^2 \|h'\|_2^2}{(N-\alpha_2D^2)^2}$. This has exactly the same form as that appearing in the proof of Theorem 9 in Appendix C of [47], except that the values of 707 708 C_1 are different.¹¹ Thus, the rest of the proof is the same as in [47]. Π 709

Note that, although the value function V of the exploratory problem (3.7) has 710 711been explicitly given by Theorem 5.2, the above theorem focuses on the equivalence of *solvability* of the two problems without having to know the explicit expression of 712 the value function of either problem. Hence we use generic letters $(\alpha_0, \alpha_1, \alpha_2)$ instead 713of (k_0, k_1, k_2) to express the value functions. 714

The following result shows that the Choquet-regularized LQ problem converges 715to its classical counterpart if the exploration weight λ goes to zero. 716

PROPOSITION 6.2. Assume that statement (a) (or equivalently, (b)) of Theorem 717 6.1 holds. Then, for each $x \in \mathbb{R}$, $\lim_{\lambda \to 0} \Pi^*(\cdot; x) = \delta_{u^*(x)}(\cdot)$ weakly. Moreover, for 718

¹¹There is a typo in the title of Appendix C of [47]: it should be the proof of Theorem 9, instead of Theorem 7.

719 each $x \in \mathbb{R}$, $\lim_{\lambda \to 0} |V(x) - V^{cl}(x)| = 0$.

720 *Proof.* Noting that $\lim_{\lambda\to 0} \frac{\lambda^2 ||h'||_2^2}{2\rho(N-\alpha_2 D^2)} = 0$, the proof is the same as that of 721 Theorem 11 in [47].

Finally, we examine the "cost of exploration" – the loss in the original (i.e., nonregularized) objective due to exploration, which was originally defined and derived in [47] for problems with entropy regularization. The notion can be extended readily to the current Choquet setting, namely, it is the difference between the two optimal value functions, adjusting for the additional contribution coming from the Choquet regularizer of the optimal exploratory strategy:

728 (6.2)
$$\mathcal{C}^{u^*,\Pi^*}(x) := V^{\mathrm{cl}}(x) - \left(V(x) - \lambda \mathbb{E}_x \left[\int_0^\infty e^{-\rho t} \left(\int_U uh'(1 - \Pi_t^*(u)) \mathrm{d}\Pi_t^*(u)\right) \mathrm{d}t\right]\right),$$

for $x \in \mathbb{R}$.

THEOREM 6.3. Assume that statement (a) (or equivalently, (b)) of Theorem 6.1 holds. Then, the exploration cost for the stochastic LQ problem is

732 (6.3)
$$\mathcal{C}^{u^*,\Pi^*}(x) = \frac{\lambda^2 \|h'\|_2^2}{2\rho(N - \alpha_2 D^2)}, \text{ for } x \in \mathbb{R}$$

733 Proof. Let $\{\Pi_t^*\}_{\{t \ge 0\}}$ be the open-loop control generated by the feedback control 734 $\Pi^*(\cdot; x)$ given in statement (a) with respect to the initial state x whose quantile 735 function is $Q_{\Pi^*(\cdot;x)}(p) = \frac{(\alpha_2(B+CD)-R)x+\alpha_1B-L}{N-\alpha_2D^2} + \frac{\lambda h'(1-p)}{N-\alpha_2D^2}$ with the mean and variance 736 given by $\mu^*(x) = \frac{(\alpha_2(B+CD)-R)x+\alpha_1B-L}{N-\alpha_2D^2}$ and $(\sigma^*(x))^2 = \frac{\lambda^2 ||h'||_2^2}{(N-\alpha_2D^2)^2}$. By Lemma 4.1, 737 it is straightforward to calculte $\int_U uh'(1-\Pi_t^*(u))d\Pi_t^*(u) = \frac{\lambda ||h'||_2^2}{N-\alpha_2D^2}$. The desired 738 result now follows immediately from the definition (6.2) and the expressions of $V(\cdot)$ 739 in (a) and $V^{cl}(\cdot)$ in (b).

The costs of exploration derived in [47, 48] for the entropy setting depend on 740 only the temperature parameter and the discounting rate or time horizon which are 741 742chosen by the agents, but not on the state dynamics or the reward coefficients which the agents generally do not know about. In contrast, the derived exploration cost 743 (6.3) for the Choquet setting does depend on the unknown model parameters in a 744 complicated way (mainly through α_2), which seems to be a disadvantage from the 745 learning perspective. However, a bit of reflection reveals that it is more important to 746 know what impact the cost than to know the *precise* value of the cost. For example, 747 748 (6.3) suggests a way to strategically select the regularizers: other things being equal, to reduce the exploration cost one should choose regularizers with smaller values of $||h'||_2$. Moreover, $C^{u^*,\Pi^*}(x) \leq \frac{\lambda^2 ||h'||_2^2}{2\rho N}$ noting $\alpha_2 < 0$; so the cost is bounded above by a constant that is inversely proportional to the unknown parameter N, the 749 750 751 control weight in the reward function. As a result, when executing controls becomes 752increasingly costly, the exploration cost diminishes because the agent is less motivated 753 to do exploration. Furthermore, $C^{u^*,\Pi^*}(x) = \frac{\lambda \|h'\|_2}{2\rho} \sigma^*(x)$, meaning that the cost is proportional to the standardized deviation of the exploratory control, a feature that is 754755 not presented in the entropy setting [47]. Finally, the exploration cost (6.3) depends 756 on λ and ρ in a rather intuitive way: it increases as λ increases, due to more emphasis 757 placed on exploration, or as ρ decreases, indicating an effectively longer horizon for 758exploration. 759

CHOQUET REGULARIZATION

760 7. Conclusion. This paper develops a framework for continuous-time RL that 761can generate or indeed interpret/explain many broadly practiced distributions for exploration. The main contributions are conceptual/theoretical rather than algorith-762 mic: Theorem 5.2 does not lead directly to an algorithm to compute optimal policies, 763 because the expression (5.8) involves the model parameters which are unknown in the 764 RL context. That said, our results do provide important guidance for devising RL 765 algorithms. First, Theorem 5.2 may imply a provable policy improvement theorem 766 and hence result in a q-learning theory analogous to that in the entropy-regularized 767 setting recently established by [27]. Second, the explicit form (5.8) can suggest special 768 structure of function approximators for learning optimal distributions, thereby greatly 769 reduce the number of parameters needed for function approximation and improve the 770 771 efficiency of the resulting learning algorithms. Finally, the availability of a large class 772 of Choquet regularizers makes it possible to compare and choose specific regularizers to achieve certain objectives specific to each learning problem. 773

Another conceptual contribution of the paper is that it establishes a link between risk metrics and RL. This paper is the first to do so, and the attempt is by no means comprehensive. The rich literature on decision theory and risk metrics is expected to further bring in new insights and directions into the RL study, not only related to regularization, but also in terms of motivating new objective functions and axiomatic approaches for learning.

The theory developed in this paper opens up several research directions. Here we comment on some. One is to develop the corresponding q-learning theory mentioned earlier. Another is to find the "best Choquet regularizer" in terms of efficiency of the resulting RL algorithms. Yet another problem is in financial application: to formulate a continuous-time mean-variance portfolio selection problem with a Choquet regularizer and compare the performance with its entropy counterpart solved in [48].

Last but not least, the Choquet regularizers proposed in this paper are defined for distributions on \mathbb{R} , while many RL applications involve multi-dimensional action spaces. Because Choquet regularizers are characterized by quantile additivity as in Theorem 2.4 while quantile functions are not well defined for distributions on \mathbb{R}^d with d > 1, it is very challenging to study Choquet regularizers in high dimensions. To overcome the difficulty, the first possible attempt is to minic (2.2) by defining, for distributions Π on \mathbb{R}^d , the functional $\Phi_h^{\text{joint}}(\Pi) = \int_{\mathbb{R}^d} h \circ \Pi([\mathbf{x}, \infty)) d\mathbf{x}$. This formulation requires some further conditions on $h \in \mathcal{H}$ to guarantee desirable properties, and it is unclear whether we can derive the corresponding optimizers in a form similar to Proposition 4.2. Another possible idea is to use

$$\Phi_h^{\text{sum}}(\Pi) = \sum_{i=1}^d \int_{\mathbb{R}} h \circ \Pi_i([x,\infty)) dx \text{ or } \Phi_h^{\text{prod}}(\Pi) = \prod_{i=1}^d \int_{\mathbb{R}} h \circ \Pi_i([x,\infty)) dx$$

where Π_i is the *i*-th marginal distribution of Π . This formulation relies only on the marginal distributions of Π , allowing us to utilize the existing results for Choquet regularizers on \mathbb{R} . Either formulation mentioned above requires a thorough analysis in a future study.

Acknowledgements. Han is supported by the Fundamental Research Funds for the Central Universities, Nankai University (Grant No. 63231138). Wang is supported by the Natural Sciences and Engineering Research Council of Canada (RGPIN-2018-03823, RGPAS-2018-522590). Zhou is supported by a start-up grant and the Nie Center for Intelligent Asset Management at Columbia University. His work is also part of a Columbia-CityU/HK collaborative project that is supported by the InnotHK Initiative, The Government of the HKSAR, and the AIFT Lab. The authors are grateful to the two anonymous referees for their constructive comments that have led to an improved version of the paper.

799

REFERENCES

- [1] B. ACCIAIO AND G. SVINDLAND, Are law-invariant risk functions concave on distributions?,
 Dependence Modeling, 1 (2013), pp. (2013):54–64.
- [2] C. ACERBI, Spectral measures of risk: A coherent representation of subjective risk aversion,
 Journal of Banking and Finance, 26 (2002), pp. (7):1505–1518.
- [3] P. ARTZNER, F. DELBAEN, J.-M. EBER, AND D. HEATH, Coherent measures of risk, Mathemat ical Finance, 9 (1999), pp. (3):203–228.
- [4] F. BELLINI, T. FADINA, R. WANG, AND Y. WEI, Parametric measures of variability induced by risk measures, Insurance: Mathematics and Economics, 106 (2022), pp. 270–284.
- [5] L. BENAC AND F. GODIN, Risk averse non-stationary multi-armed bandits, arXiv preprint
 arXiv:2109.13977, (2021).
- [6] P. BENIGNO AND M. WOODFORD, Optimal monetary and fiscal policy: A linear-quadratic approach, NBER macroeconomics annual, 18 (2003), pp. 271–333.
- [7] P. BENIGNO AND M. WOODFORD, Linear-quadratic approximation of optimal policy problems,
 Journal of Economic Theory, 147 (2012), pp. (1):1–42.
- [8] J. BLANCHET, L. CHEN, AND X. ZHOU, Distributionally robust mean-variance portfolio selection
 with wasserstein distances, Management Science, 68 (2022), pp. (1):6382–6410.
- [9] J. Q. CHANG, Q. ZHU, AND V. Y. TAN, Risk-constrained thompson sampling for cvar bandits,
 arXiv preprint arXiv:2011.08046, (2020).
- [10] A. DE WAEGENAERE AND P. P. WAKKER, Nonmonotonic choquet integrals, Journal of Mathe matical Economics, 36 (2001), pp. (1):45–60.
- [11] F. DELBAEN, Coherent risk measures on general probability spaces. In Advances in Finance
 and Stochastics (pp. 1–37), Springer, Berlin, Heidelberg, 2002.
- 822 [12] D. DENNEBERG, Non-additive Measure and Integral, Springer Science & Business Media, 1994.
- [13] L. EECKHOUDT AND R. LAEVEN, Dual moments and risk attitudes, Operations Research, 70
 (2022), pp. 1330–1341.
- [14] P. ESFAHANI AND D. KUHN, Data-driven distributionally robust optimization using the wasserstein meric: Performance guarantees and tractable reformulations, Mathematical Programming, 171 (2018), pp. (1):115–166.
- [15] H. FÖLLMER AND A. SCHIED, Convex measures of risk and trading constraints, Finance and
 Stochastics, 6 (2002), pp. (4):429–447.
- [16] H. FÖLLMER AND A. SCHIED, Stochastic Finance. An Introduction in Discrete Time. Fourth
 Edition, Walter de Gruyter, Berlin, 2016.
- [17] E. FURMAN, R. WANG, AND R. ZITIKIS, Gini-type measures of risk and variability: Gini shortfall, capital allocation and heavy-tailed risks, Journal of Banking and Finance, 83 (2017), pp. 70–84.
- [18] X. GAO, Z. XU, AND X. ZHOU, State-dependent temperature control for langevin diffusions,
 SIAM Journal on Control and Optimization, 60 (2022), pp. (3):1250–1268.
- [19] I. GILBOA AND D. SCHMEIDLER, Maxmin expected utility with non-unique prior, Journal of Mathematical Economics, 18 (1989), pp. (2):141–153.
- [20] G. GLASSER, Variance formulas for the mean difference and coefficient of concentration, Jour nal of the American Statistical Association, 57 (1962), pp. (299):648–654.
- [21] B. GRECHUK, A. MOLYBOHA, AND M. ZABARANKIN, Maximum entropy principle with general deviation measures, Mathematics of Operations Research, 34 (2009), pp. 445–467.
- [22] X. GUO, R. XU, AND T. ZARIPHOPOULOU, Entropy regularization for mean field games with
 learning, Mathematics of Operations Research, 47 (2022), pp. (4):3239–3260.
- [23] T. HAARNOJA, T. ZHOU, P. ABBEEL, AND S. LEVINE, Soft actor-critic: Off-policy maximum
 entropy deep reinforcement learning with a stochastic actor, in International Conference
 on Machine Learning, 2018, pp. 1856–1865.
- [24] T. HU AND O. CHEN, On a family of coherent measures of variability, Insurance: Mathematics
 and Economics, 95 (2020), pp. 173–182.
- [25] Y. JIA AND X. ZHOU, Policy evaluation and temporal-difference learning in continuous time
 and space: A martingale approach., Journal of Machine Learning Research, 23 (2022a),
 pp. (154):1–55.
- [26] Y. JIA AND X. ZHOU, Policy gradient and actor-critic learning in continuous time and space:
 Theory and algorithms, Journal of Machine Learning Research, (2022b), pp. (275):1–50.

[27] Y. JIA AND X. ZHOU, q-learning in continuous time, Journal of Machine Learning Research,

K. JUDD, Numerical Methods in Economomics, MIT Press, Cambridge, MA, 1998.

855

856

857

[28]

forthcoming, (2022c).

- 858 [29]S. KUSUOKA, On law invariant coherent risk measures, Advances in Mathematical Economics, 859 3 (2001), pp. 83–95. [30] W. LI AND E. TODOROV, Iterative linearization methods for approximately optimal control and 860 estimation of non-linear stochastic system, International Journal of Control, 80 (2007), 861 pp. (9):1439-1453. 862 863 [31] F. LIU, J. CAI, C. LEMIEUX, AND R. WANG, Convex risk functionals: Representation and 864 applications, Insurance: Mathematics and Economics, 90 (2020), pp. 66–79. 865 [32] C. MOU, W. ZHANG, AND C. ZHOU, Robust exploratory mean-variance problem with drift 866 uncertainty, arXiv preprint arXiv:2108.04100, (2021). [33] O. NACHUM, M. NOROUZI, K. XU, AND D. SCHUURMANS, Bridging the gap between value and 867 policy based reinforcement learning, in Advances in Neural Information Processing Systems, 868 869 2017, pp. 2775-2785. 870 [34] S. PESENTI, Q. WANG, AND R. WANG, Optimizing distortion risk metrics with distributional 871 uncertainty, arXiv preprint arXiv:2011.04889, (2020). [35] G. PFLUG AND D. WOZABAL, Ambiguity in portfolio selection, Quantitative Finance, 7 (2007), 872 873 pp. (4):435-442. 874 [36] G. PSARRAKOS AND J. NAVARRO, Generalized cumulative residual entropy and record values, Metrika, 76 (2013), pp. (5):623-640. 875 876 [37] J. QUIGGIN, A theory of anticipated utility, Journal of Economic Behavior and Organization, 3 877 (1982), pp. (4):323-343. 878 [38] M. RAO, Y. CHEN, B. C. VEMURI, AND F. WANG, Cumulative residual entropy: A new measure 879 of information, IEEE Transactions on Information Theory, 50 (2004), pp. 1220–1228. 880 [39] R. T. ROCKAFELLAR, S. URYASEV, AND M. ZABARANKIN, Generalized deviation in risk analysis, 881 Finance and Stochastics, 10 (2006), pp. 51–74. 882 [40] M. ROTHSCHILD AND J. E. STIGLITZ, Increasing risk: I. A definition, Journal of Economic 883 Theory, 2 (1978), pp. 99–121. 884 [41] D. SCHMEIDLER, Subjective probability and expected utility without additivity, Econometrica, 57 885 (1989), pp. (3):571-587. 886 [42] S. M. SUNOJ AND P. G. SANKARAN, Quantile based entropy function, Statistics and Probability 887 Letters, 82 (2012), pp. (6):1049-1053. 888 [43] W. TANG, Y. ZHANG, AND X. ZHOU, Exploratory HJB equations and their convergence, SIAM Journal on Control and Optimization, 60 (2022), pp. (6):3191-3216. 889 890 [44] E. TODOROV AND W. LI, A generalized iterative lqg method for locally-optimal feedback con-891 trol of constrained nonlinear stochastic systems, in American Control Conference, 2005. 892 Proceedings of the 2005, IEEE, 2005, pp. 300-306. 893 [45] A. TOOMAJ, S. M. SUNOJ, AND J. NAVARRO, Some properties of the cumulative residual entropy 894 of coherent and mixed systems, Journal of Applied Probability, 54 (2017), pp. (2):379–393. 895 [46] A. TVERSKY AND D. KAHNEMAN, Advances in prospect theory: Cumulative representation of
- uncertainty, Journal of Risk and Uncertainty, 5 (1992), pp. 297–323.
 [47] H. WANG, T. ZARIPHOPOULOU, AND X. ZHOU, Exploration versus exploitation in reinforcement
- [11] In third, T. Zhan horocoloo, Alb A. Enov, Exploration versus exploration in reinforcement
 learning: A stochastic control approach, Journal of Machine Learning Research, 21 (2020a),
 pp. (198):1–34.
- [48] H. WANG AND X. ZHOU, Continuous-time mean-variance portfolio selection: A reinforcement learning framework., Mathematical Finance, 30 (2020), pp. 1273–1308.
- [49] Q. WANG, R. WANG, AND Y. WEI, Distortion risk metrics on general spaces, ASTIN Bulletin,
 50 (2020b), pp. 827–851.
- [50] R. WANG, Y. WEI, AND G. E. WILLMOT, Characterization, robustness and aggregation of signed choquet integrals, Mathematics of Operations Research, 45 (2020c), pp. 993–1015.
- [51] R. WANG AND R. ZITIKIS, An axiomatic foundation for the expected shortfall, Management
 Science, 67 (2021), pp. (3):1413–1429.
- [52] S. WANG, V. R. YOUNG, AND H. H. PANJER, Axiomatic characterization of insurance prices, Insurance: Mathematics and Economics, 21 (1997), pp. (2):173–183.
- 910 [53] M. E. YAARI, The dual theory of choice under risk, Econometrica, 55 (1987), pp. (1):95–115.
- [54] B. D. ZIEBART, A. L. MAAS, J. A. BAGNELL, AND A. K. DEY, Maximum entropy inverse reinforcement learning, in AAAI Conference on Artificial Intelligence, vol. 8, 2008, pp. 1433– 1438.