

BACKWARD STOCHASTIC CONTROL SYSTEM WITH ENTROPY REGULARIZATION*

ZIYUE CHEN [†] AND QI ZHANG [‡]

Abstract. The entropy regularization is inspired by information entropy from machine learning and the ideas of exploration and exploitation in reinforcement learning, which appears in the control problem to design an approximating algorithm for the optimal control. This paper is concerned with the optimal exploratory control for backward stochastic system, generated by the backward stochastic differential equation and with the entropy regularization in its cost functional. We give the theoretical depict of the optimal relaxed control so as to lay the foundation for the application of such a backward stochastic control system to mathematical finance and algorithm implementation. For this, we first establish the stochastic maximum principle by convex variation method. Then we prove sufficient condition for the optimal control and demonstrate the implicit form of optimal control. Finally, the existence and uniqueness of the optimal control for backward linear-quadratic control problem with entropy regularization is proved by decoupling techniques.

Key words. backward stochastic control system, relaxed control, entropy regularization, maximum principle, linear-quadratic problem.

MSC codes. 93E20, 93C15

1. Introduction. Different from the deterministic system which has only one path, there is much difference between the forward stochastic system and the backward one. It is well known that stochastic differential equation (SDE) and backward stochastic differential equation (BSDE) are much different from the form of equation to the form of solution. In fact, BSDE and forward-backward stochastic differential equation (FBSDE) play a special role in many problems of stochastic analysis and stochastic controls. For example, Peng [23] demonstrates that the solution to FBSDE gives the probabilistic interpretation of nonlinear PDE which is known as nonlinear Feynman-Kac formula, Duffie and Epstein [7] put forward the stochastic differential utility which is actually a BSDE with conditional expectation, Zhang and Zhao [30] constructs the stationary solution to parabolic stochastic partial differential equation (SPDE) based on the idea that infinite horizon backward doubly stochastic differential equation can serves as "elliptic" SPDE to give the pathwise steady statue of parabolic SPDE, to name but a few. Without the exception of control system, the forward stochastic control system and backward stochastic control system are also different. The difference not only lies on the state equation, i.e. the state equation of backward stochastic control system is controlled BSDE rather than SDE, but also on the application. It is well known that the controlled BSDE is widely used in mathematical finance, for example, as the dynamic equation for the value of portfolio to replicate a contingent claim. Also, as shown in Karnam, Ma and Zhang [17], many time-inconsistent optimization problems can be transformed into a stochastic controlled problem with multidimensional BSDE dynamics by the so-called dynamic utility approach. Recently, the controlled BSDE is also applied to the numerical calcu-

*This paper is supported by National Key R&D Program of China (No.2022YFA1006101), National Natural Science Foundation of China (No.12371445) and the Science and Technology Commission of Shanghai Municipality (No.22ZR1407600).

[†]School of Mathematical Sciences, Fudan University, Shanghai 200433, China (chenziyue21@m.fudan.edu.cn).

[‡]Corresponding author. School of Mathematical Sciences, Fudan University, Shanghai 200433, China and Laboratory of Mathematics for Nonlinear Science, Fudan University, Shanghai 200433, China (qzh@fudan.edu.cn).

lation of partial differential equation based on the theory of nonlinear Feynman-Kac formula. In an extended work of the deep BSDE numerical scheme by Takahashi, Tsuchida and Yamada [27], the controlled BSDE is used to make the scheme more efficient and stable.

The state equation of backward stochastic control system is a controlled BSDE as below

$$(1.1) \quad \begin{cases} -dy^\pi(t) = f(t, y^\pi(t), z^\pi(t), \pi_t)dt - z^\pi(t)dW(t), \\ y(T) = \xi, \end{cases}$$

where π is the control variable. The studies on the backward stochastic control system emerged soon after the solvability of nonlinear BSDE. Here we recall some early results. Peng [24] first derived the local stochastic maximum principle in 1993. El Karoui, Peng and Quenez [9] demonstrated the application of controlled BSDE in finance. Dokuchaev and Zhou [6] applied the backward stochastic control system to pricing European contingent claims and derived the global stochastic maximum principle in nonconvex control domain. For the linear-quadratic (LQ) case, it was studied in 2001 by Lim and Zhou [20] with the help of decoupling techniques.

For $\sigma > 0$, the first motivation to write this paper is to study the relaxed control of a type of backward stochastic control system with an entropy-regularized cost functional as below

$$(1.2) \quad J^\sigma(\pi) = \mathbb{E} \left[\int_0^T (l(t, y^\pi(t), z^\pi(t), \pi_t) + \frac{\sigma^2}{2} Ent(\pi_t | e^{-U}))dt + \phi(y^\pi(0)) \right].$$

The relaxed control means that the value of the control could depend on a distribution of value space, which actually enhances the possibility to get an optimal control, especially in a case that the classical control doesn't exist. This concept was put forward by Becker and Mandrekar [1] for deterministic control system in 1969 and extended to stochastic control system by Fleming [11] in 1978, and later El Karoui, Huu Nguyen and Jeanblanc-Picqu  [8] further study the relaxed control for stochastic control system with degenerate diffusion. It is worth noting that the relaxed control has a significant applications in machine learning algorithms, especially with the entropy regulation in the cost functional since the use of relaxed controls along with entropy regulation improves algorithm stability and efficiency. Actually, the entropy regulation has been widely used in numerical calculus for a long time. For example, the entropy regularization was ever applied to iterative numerical scheme for solving PDEs, and one can refer to Jordan, Kinderlehrer and Otto [16], Gomes and Valdinoci [13] for details. To apply this method to reinforcement learning, Wang, Zariphopoulou and Zhou [28] studied the relaxed control with entropy-regularized cost functional to devise an exploratory formulation for the feature dynamics which captures learning under exploration based on the dynamic programming method. Moreover, in LQ case, [28] proved that the optimal control is Gaussian. In the meantime, Wang and Zhou [29] showed that it has potential application in continuous mean-variance optimal portfolio problem. From the point of view of stability, Reisinger and Zhang [25] demonstrated the regularised relaxed control formulation ensures that the optimal controls are stable with respect to model perturbations. A recent version of Šiřka and Szpruch [26] added the priori reference measure into the entropy regularization and studied this general control system by the maximum principle method. Besides, there are more results on this topic for forward stochastic control system emerging recently,

such as Gao, Xu and Zhou [12], Firoozi and Jaimungal [10], Jia and Zhou [15], Dai, Dong, Jia and Zhou [5], etc.

Unlike the control systems or research methods in [28] and [26], we study the backward stochastic control systems (1.1) and (1.2) based on maximum principle method, which allows us to consider the control system with random coefficients. We get the necessary and sufficient condition for the optimal control which establishes a theoretical foundation for future applications on mathematical finance and machine learning. With the depict of the optimal control derived from the stochastic Hamiltonian system, it provides us a chance to give a specific implicit or explicit form of the relaxed optimal control. Actually, in the LQ case, we borrow the idea of decoupling techniques for backward stochastic control system in Lim and Zhou [20] to prove the optimal control exactly exists and obtain its explicit form. As expected, the optimal relaxed control for backward stochastic LQ control systems with entropy regularization appears to be Gaussian. According to the theoretical depict of the optimal control for the backward stochastic control systems, we can also design an algorithm to approximate the optimal control by the method of successive approximation. Due to limited space, we will study it in another paper.

This paper is organized as follows. We introduce the necessary notation and state the backward stochastic control problem with entropy regularization in Section 2. Then we prove the extended maximum principle for our concerned control problem, and get the necessary condition of the optimal relaxed control in Section 3. The sufficient condition of the optimal relaxed control is given in Section 4, and the implicit form of optimal control is also discussed. In Section 5, we prove the existence and uniqueness of the optimal control for stochastic LQ control problem with entropy regularization and give the explicit form of optimal control in LQ case.

2. Notation and Control Problem. Given a separable metric space E , let $\mathcal{M}(E)$ denote the set of all measures on E , $\mathcal{P}_q(E)$ denote the set of probability measures defined on E with finite q -th moment for $q \in \mathbb{N}$ and $\mathcal{P}(E) = \mathcal{P}_0(E)$ denote the set of all probability measures on E . In this paper the entropy of the measure $\pi \in \mathcal{P}(E)$ is defined as below

$$R(\pi) = \begin{cases} -\int_E \frac{d\pi}{du} \ln \frac{d\pi}{du} du, & \text{if } \pi \ll \lambda(du), \\ -\infty, & \text{otherwise,} \end{cases}$$

where λ is the Lebesgue measure on a separable metric space E . In order to make above entropy well defined, we assume that all probability-measure-valued control in this paper are absolutely continuous with respect to λ . Hence by Radon-Nikodym theorem, there exists a measurable function $g : U \rightarrow \mathbb{R}^+$ such that $\mu(du) = g(u)du$ a.e. Throughout the paper we will abuse the notation of probability measure which are absolutely continuous with respect to Lebesgue measure and do not distinguish a probability measure from its density which exists due to Radon-Nikodym theorem.

We begin with a finite time horizon $[0, T]$ for $T > 0$ and a complete filtered probability space $(\Omega, \mathcal{F}, \mathbb{P})$, on which a standard \mathbb{R}^m -valued Brownian motion W is defined. Moreover, $\mathbb{F} = (\mathcal{F}_t)_{0 \leq t \leq T}$ is the natural filtration generated by W and $\mathcal{F}_T = \mathcal{F}$. Next we introduce some useful spaces. For $t \in [0, T]$ and a Hilbert space S with norm $\|\cdot\|_S$ and Borel σ -field \mathcal{S} , we define

- $S_{\mathbb{F}}^2(0, T; S)$: the space of all \mathbb{F} -adapted processes $x : \Omega \times [0, T] \rightarrow S$ satisfying $t \rightarrow x(t)$ is a.s. continuous and $\mathbb{E} [\sup_{0 \leq t \leq T} \|x(t)\|_S^2] < +\infty$;
- $L_{\mathbb{F}}^2(0, T; S)$: the space of all \mathbb{F} -adapted processes $x : \Omega \times [0, T] \rightarrow S$ satisfying $\mathbb{E} \left[\int_0^T \|x(t)\|_S^2 dt \right] < +\infty$;

- $L^2_{\mathcal{F}_t}(\Omega; S)$: the space of all \mathcal{F}_t -measurable random variables $x : \Omega \rightarrow S$ satisfying $\mathbb{E} [\|x(t)\|_S^2] < +\infty$;
- $L^\infty(0, T; S)$: the space of all measurable maps $x : [0, T] \rightarrow S$ satisfying $\|x\|_{L^\infty} = \text{ess sup}_{t \in [0, T]} \|x(t)\|_S < \infty$;
- $L^\infty(S)$: the space of all measurable maps $x : S \rightarrow \mathbb{R}$ satisfying $\|x\|_{L^\infty(S)} = \text{ess sup}_{a \in S} |x(a)| < \infty$;
- C^∞ : the space of all infinitely differentiable functions.

To avoid heavy notations, we omit the transpose symbols in this paper unless necessary.

For the backward stochastic control systems (1.1) and (1.2), the set of admissible controls \mathcal{A} is defined as follows.

DEFINITION 2.1. For $E = [0, T] \times \mathbb{R}^p$, we define a subset of $\mathcal{M}(E)$:

$$\mathcal{M}_2 := \left\{ \pi \in \mathcal{M}(E) : \text{for a.a. } t \in [0, T], \text{ there exists } \pi_t \in \mathcal{P}(\mathbb{R}^p) \text{ such that} \right. \\ \left. \pi(da, dt) = \pi_t(a)dadt, \int_0^T \int |a|^2 \pi_t(a)dadt < \infty \right\}.$$

Here and in the rest of this paper, the integration without explicit domain is over \mathbb{R}^p unless indicated explicitly. Then the set of admissible controls

$$\mathcal{A} := \left\{ \pi : \Omega \rightarrow \mathcal{M}_2 : \mathbb{E} \left[\int_0^T \text{Ent}(\pi_t | e^{-U})dadt \right] < \infty, \mathbb{E} \left[\int_0^T \int |a|^2 \pi_t(a)dadt \right] < \infty \right. \\ \left. \text{and } \pi_t \text{ is } \mathcal{F}_t\text{-measurable for any } t \in [0, T] \right\}.$$

where $U : \mathbb{R}^p \rightarrow \mathbb{R}$ is a measurable function such that e^{-U} is a density function. Here e^{-U} is regarded as the priori reference measure and $\text{Ent}(\cdot | \cdot)$ is the relative entropy characterized by

$$\text{Ent}(m | m') := \int (\ln m(a) - \ln m'(a)) m(a)da,$$

for $m, m' \in \mathcal{P}(\mathbb{R}^p)$ which are absolutely continuous with respect to the Lebesgue measure (otherwise $\text{Ent}(m | m') = \infty$).

Remark 2.2. By Lemma A.1 in Kerimkulov, Šiška, Szpruch and Zhang [18], the admissible control set \mathcal{A} is convex due to the fact that $\pi \mapsto \text{Ent}(\pi | e^{-U})$ is convex.

For the given measurable functions $f : \Omega \times [0, T] \times \mathbb{R}^n \times \mathbb{R}^{n \times m} \times \mathcal{P}(\mathbb{R}^p) \rightarrow \mathbb{R}^n$, $l : \Omega \times [0, T] \times \mathbb{R}^n \times \mathbb{R}^{n \times m} \times \mathcal{P}(\mathbb{R}^p) \rightarrow \mathbb{R}$, $\xi : \Omega \rightarrow \mathbb{R}^n$, $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and the admissible control $\pi \in \mathcal{A}$, we consider the backward stochastic system (1.1) and (1.2). Then the control problem is

(P1): to find an optimal $\mu \in \mathcal{A}$ such that

$$J^\sigma(\mu) = \inf_{\pi \in \mathcal{A}} J^\sigma(\pi).$$

Since the concerned system involves the measure-valued control, we need define flat derivative on a convex subset $\mathcal{C} \subseteq \mathcal{P}(\mathbb{R}^p)$. The following definition refers to [18] based on the ideas from early literatures (e.g. Lions [21], Carmona and Delarue [3], Buckdahn, Li, Peng and Rainer [2], etc.).

DEFINITION 2.3. A functional $F : \mathcal{C} \rightarrow \mathbb{R}^d$ is said to admit a linear derivative if there is a continuous map $\frac{\delta F}{\delta m} : \mathcal{C} \times \mathbb{R}^p \rightarrow \mathbb{R}^d$ such that for all $m, m' \in \mathcal{C}$, it holds that

$\int |\frac{\delta F}{\delta m}(m)(a)| m'(a) da < \infty$, and

$$(2.1) \quad F(m') - F(m) = \int_0^1 \int \frac{\delta F}{\delta m}(m + \lambda(m' - m))(a) \cdot (m'(a) - m(a)) da d\lambda.$$

In the above definition, $\frac{\delta F}{\delta m}$ is only defined up to a constant according to Remark 5.46 in Carmona and Delarue [4]. The functional $\frac{\delta F}{\delta m}$ is then called the linear (functional) derivative of F on \mathcal{C} . Note that if $\frac{\delta F}{\delta m}$ exists, for any $\nu, \mu \in \mathcal{C}$,

$$(2.2) \quad \lim_{\varepsilon \rightarrow 0^+} \frac{F(\nu + \varepsilon(\mu - \nu)) - F(\nu)}{\varepsilon} = \int \frac{\delta F}{\delta m}(\nu)(a)(\mu(a) - \nu(a)) da.$$

Obviously, (2.1) implies (2.2). To see the implication in the other direction, take $\nu^\lambda := \nu + \lambda(\mu - \nu)$ and $\mu^\lambda := \mu - \nu + \nu^\lambda$ and notice that (2.2) ensures that for all $\lambda \in [0, 1]$,

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0^+} \frac{F(\nu^\lambda + \varepsilon(\mu - \nu)) - F(\nu^\lambda)}{\varepsilon} &= \lim_{\varepsilon \rightarrow 0^+} \frac{F(\nu^\lambda + \varepsilon(\mu^\lambda - \nu^\lambda)) - F(\nu^\lambda)}{\varepsilon} \\ &= \int \frac{\delta F}{\delta m}(\nu^\lambda)(a)(\mu^\lambda - \nu^\lambda) da = \int \frac{\delta F}{\delta m}(\nu^\lambda)(a)(\mu(a) - \nu(a)) da. \end{aligned}$$

By the fundamental theorem of calculus, we can derive that

$$(2.3) \quad F(\mu) - F(\nu) = \int_0^1 \lim_{\varepsilon \rightarrow 0^+} \frac{F(\nu^\lambda + \varepsilon) - F(\nu^\lambda)}{\varepsilon} d\lambda = \int_0^1 \int \frac{\delta F}{\delta m}(\nu^\lambda)(a)(\mu(a) - \nu(a)) da d\lambda.$$

The linear derivative $\frac{\delta F}{\delta \nu}$ is here also defined up to the additive constant as any constant can be added to $\int \frac{\delta F}{\delta \nu}(\nu, t)(a) \nu_t(da)$ without affecting the definition formula. Note that if $\frac{\delta F}{\delta \nu}$ exists, then similarly we have an equivalent form of (2.3) as below

$$\forall \nu, \nu' \in \mathcal{A}, \lim_{\varepsilon \rightarrow 0^+} \frac{F(\nu + \varepsilon(\nu' - \nu)) - F(\nu)}{\varepsilon} = \int \frac{\delta F}{\delta \nu}(\nu, t)(a)(\nu'_t(a) - \nu_t(a)) da.$$

Then we state the assumptions in our concerned control problem.

ASSUMPTION 2.4. (i) $\xi \in L^2_{\mathcal{F}_T}(\Omega; \mathbb{R}^n)$. Denote by \mathcal{P} the predictable sub- σ algebra of $\mathcal{F} \otimes \mathcal{B}([0, T])$, and f is $\mathcal{P} \otimes \mathcal{B}(\mathbb{R}^n) \otimes \mathcal{B}(\mathbb{R}^{n \times m}) \otimes \mathcal{B}(\mathcal{P}(\mathbb{R}^p))$ -measurable and there exists a $m \in \mathcal{P}(\mathbb{R}^p)$ such that $f(\cdot, 0, 0, m) \in L^2_{\mathcal{F}}(0, T; \mathbb{R}^n)$. Also for $(\omega, t) \in \Omega \times [0, T]$, $f(t, y, z, m)$ is continuously differentiable with respect to $(y, z) \in \mathbb{R}^n \times \mathbb{R}^{n \times m}$ and continuously twice differentiable with respect to $m \in \mathcal{P}_2(\mathbb{R}^p)$ (in the sense of flat derivative), and for any $(\omega, t) \in \Omega \times [0, T]$, $(y, z), (y', z') \in \mathbb{R}^n \times \mathbb{R}^{n \times m}$, $m \in \mathcal{P}_2(\mathbb{R}^p)$, there exists a constant $K > 0$ such that $|\nabla_y f| + |\nabla_z f| + |\frac{\delta^2 f}{\delta m^2}| \leq K$ uniformly and

$$\begin{aligned} \left| \frac{\delta f(t, y, z, m)}{\delta m} - \frac{\delta f(t, y', z', m)}{\delta m} \right| &\leq K(|y - y'| + |z - z'|); \\ \left| \frac{\delta f(t, y, z, m)}{\delta m} \right| &\leq K(1 + |a|). \end{aligned}$$

(ii) ϕ is $\mathcal{B}(\mathbb{R}^n)$ -measurable and l is $\mathcal{P} \otimes \mathcal{B}(\mathbb{R}^n) \otimes \mathcal{B}(\mathbb{R}^{n \times m}) \otimes \mathcal{B}(\mathcal{P}(\mathbb{R}^p))$ -measurable. For $(\omega, t) \in \Omega \times [0, T]$, $\phi(y)$ is continuously differentiable with respect to $y \in \mathbb{R}^n$, and $l(t, y, z, m)$ is continuously differentiable with respect to $(y, z) \in \mathbb{R}^n \times \mathbb{R}^{n \times m}$ and

continuously twice differentiable with respect to $m \in \mathcal{P}_2(\mathbb{R}^p)$, and for any $(\omega, t) \in \Omega \times [0, T]$, $(y, z), (y', z') \in \mathbb{R}^n \times \mathbb{R}^{n \times m}$, $m \in \mathcal{P}_2(\mathbb{R}^p)$ and K in (i),

$$\begin{aligned} |\phi| &\leq K(1 + |y|^2), \\ |l| &\leq K(1 + |y|^2 + |z|^2), \\ |\nabla_y \phi| &\leq K(1 + |y|), \\ |\nabla_y l| + |\nabla_z l| &\leq K(1 + |y| + |z|), \\ \left| \frac{\delta l(t, y, z, m)}{\delta m} - \frac{\delta l(t, y', z', m)}{\delta m} \right| &\leq K(|y - y'| + |z - z'|), \\ \left| \frac{\delta l(t, y, z, m)}{\delta m} \right| &\leq K(1 + |a|^2); \\ \left| \frac{\delta^2 l}{\delta m^2} \right| &\leq K. \end{aligned}$$

From Assumption 2.4 (i), we know that f is a uniformly Lipschitz generator. Hence for any $\pi \in \mathcal{A}$, by Pardoux and Peng [22] the state equation (1.1) has a unique solution $(y^\pi, z^\pi) \in S_{\mathcal{F}}^2(0, T; \mathbb{R}^n) \times L_{\mathcal{F}}^2(0, T; \mathbb{R}^{n \times m})$.

From Assumption 2.4 (i) (ii), we know that the optimization of control problem is well-posed since

$$\mathbb{E} [|\phi(y_0^\pi)|] < \infty \quad \text{and} \quad \mathbb{E} \left[\int_0^T |l(t, y_t^\pi, z_t^\pi, \pi_t)| dt \right] < \infty,$$

together with the fact that $\mathbb{E} \left[\int_0^T \int Ent(\pi_t | e^{-U(a)}) da \right] < \infty$.

3. Maximum Principle. In this section, we denote by (Y^π, Z^π) the solution to BSDE (1.1) driven by $\pi \in \mathcal{A}$. Since the admissible control set \mathcal{A} is convex, we will work with an additional control $\pi \in \mathcal{A}$ and define $\mu^\epsilon := \mu + \epsilon(\pi - \mu)$. Denote by (Y^ϵ, Z^ϵ) the solution to BSDE (1.1) driven by $\mu^\epsilon \in \mathcal{A}$. Firstly, we have some prior estimates.

LEMMA 3.1. *Under Assumption 2.4, then*

$$\mathbb{E} \left[\sup_{t \in [0, T]} |Y_t^\epsilon - Y_t^\mu|^2 \right] + \mathbb{E} \left[\int_0^T |Z_t^\epsilon - Z_t^\mu|^2 dt \right] = O(\epsilon^2).$$

Proof. By Assumption 2.4, $\frac{\delta f}{\delta m}$ is linear growth on a . Hence a classical priori estimates of BSDE leads to

$$\begin{aligned} &\mathbb{E} \left[\sup_{t \in [0, T]} |Y_t^\epsilon - Y_t^\mu|^2 \right] + \mathbb{E} \left[\int_0^T |Z_t^\epsilon - Z_t^\mu|^2 dt \right] \\ &\leq C\mathbb{E} \left[\int_0^T |f(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t^\epsilon) - f(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t)|^2 dt \right] \\ &= C\mathbb{E} \left[\int_0^T \left| \int_0^1 \int \frac{\delta f}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) \epsilon(\pi_t(a) - \mu_t(a)) da d\lambda \right|^2 dt \right] \\ &\leq C\mathbb{E} \left[\int_0^T \left| \epsilon K \int_0^1 \int (1 + |a|)(\pi_t(a) + \mu_t(a)) da d\lambda \right|^2 dt \right] \end{aligned}$$

$$\begin{aligned} &\leq C\epsilon^2 K^2 \mathbb{E} \left[\int_0^T \int_0^1 \int (1+|a|)^2 (\pi_t(a) + \mu_t(a)) da d\lambda dt \right] \\ &\leq C_{K,T} \epsilon^2 = O(\epsilon^2). \end{aligned}$$

The first equality is due the Definition 2.3. Here and in the rest of this paper, C is a generic constant whose value may change line by line, and the subscript of C , if it is indicated, is the given parameters C depends on. \square

For $\psi = f, l; x = y, z$, set

$$\begin{cases} \nabla_x \psi(t) = \nabla_x \psi(t, Y_t^\mu, Z_t^\mu, \pi_t), \\ \nabla_x \tilde{\psi}^\epsilon(t) = \int_0^1 \nabla_x \psi(t, Y_t^\mu + \lambda(Y_t^\epsilon - Y_t^\mu), Z_t^\mu + \lambda(Z_t^\epsilon - Z_t^\mu), \pi_t) d\lambda. \end{cases}$$

By the uniform boundeness of $|\nabla_x \tilde{f}^\epsilon|$ and the linear growth of $|\frac{\delta f}{\delta m}|$, BSDE

$$(3.1) \quad \begin{cases} -d(Y_t^\epsilon - Y_t^\mu) = -\left(\nabla_y \tilde{f}^\epsilon(t)(Y_t^\epsilon - Y_t^\mu) + \nabla_z \tilde{f}^\epsilon(t)(Z_t^\epsilon - Z_t^\mu) \right. \\ \quad \left. + \epsilon \int_0^1 \int \frac{\delta f}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a)(\pi_t(a) - \mu_t(a)) da d\lambda \right) dt \\ \quad + (Z_t^\epsilon - Z_t^\mu) dW_t, \\ Y_T^\epsilon - Y_T^\mu = 0, \end{cases}$$

has a unique solution. Next, we introduce the variation equation

$$(3.2) \quad \begin{cases} dV_t = -\left(\nabla_y f(t)V_t + \nabla_z f(t)Z_t^V \right) dt + Z_t^V dW_t \\ \quad + \int \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a)(\pi_t(a) - \mu_t(a)) da dt, \\ V_T = 0. \end{cases}$$

Similarly, by the boundeness of $|\nabla_x f|$ and the linear growth of $|\frac{\delta f}{\delta m}|$, the variation equation (3.2) has a unique solution. Moreover, we can get the following estimate

$$\begin{aligned} &\mathbb{E} \left[\sup_{t \in [0, T]} |V_t|^2 \right] + \mathbb{E} \left[\int_0^T |Z_t^V|^2 dt \right] \\ &\leq C \mathbb{E} \left[\int_0^T \left| \int \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a)(\pi_t(a) - \mu_t(a)) da \right|^2 dt \right] \\ &\leq CK^2 \mathbb{E} \left[\int_0^T \int_0^1 \int (1+|a|)^2 (\pi_t(a) + \mu_t(a)) da d\lambda dt \right] \\ &\leq 2C_{K,T}. \end{aligned}$$

Set $V_t^\epsilon := Y_t^\epsilon - Y_t^\mu - \epsilon V_t$ and $Z_t^{V^\epsilon} := Z_t^\epsilon - Z_t^\mu - \epsilon Z_t^V$. Then we have more estimates.

LEMMA 3.2. *Under Assumption 2.4,*

$$\mathbb{E} \left[\sup_{t \in [0, T]} |V_t^\epsilon|^2 \right] + \mathbb{E} \left[\int_0^T (|Z_t^{V^\epsilon}|^2) dt \right] = o(\epsilon^2).$$

Proof. First note that BSDE

$$(3.3) \quad \begin{cases} d\epsilon V_t = -\epsilon \left(\nabla_y f(t) V_t + \nabla_z f(t) Z_t^V \right) dt + \epsilon Z_t^V dW_t \\ \quad + \epsilon \int_0^1 \int \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da d\lambda dt, \\ \epsilon V_T = 0, \end{cases}$$

has a unique solution. Based on BSDEs (3.1) and (3.3), by the continuity dependence of the solutions to BSDEs, we have

$$\mathbb{E} \left[\sup_{t \in [0, T]} |V_t^\epsilon|^2 \right] + \mathbb{E} \left[\int_0^T |Z_t^{V^\epsilon}|^2 dt \right] \leq C\epsilon^2 h(\epsilon),$$

where

$$\begin{aligned} h(\epsilon) := & \mathbb{E} \left[\int_0^T \left(\left| \left(\nabla_y \tilde{f}^\epsilon(t) - \nabla_y f(t) \right) V_t + \left(\nabla_z \tilde{f}^\epsilon(t) - \nabla_z f(t) \right) Z_t^V \right. \right. \\ & \left. \left. + \int_0^1 \int \left(\frac{\delta f}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) \right. \right. \right. \\ & \left. \left. \left. - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da d\lambda \right|^2 dt \right). \end{aligned}$$

It is sufficient to show that $\lim_{\epsilon \rightarrow 0^+} h(\epsilon) = 0$. Note that

$$(3.4) \quad \begin{aligned} h(\epsilon) & \leq 2\mathbb{E} \left[\int_0^T \left| \left(\nabla_y \tilde{f}^\epsilon(t) - \nabla_y f(t) \right) V_t + \left(\nabla_z \tilde{f}^\epsilon(t) - \nabla_z f(t) \right) Z_t^V \right|^2 dt \right] \\ & \quad + 2\mathbb{E} \left[\int_0^T \left| \int_0^1 \int \left(\frac{\delta f}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) \right. \right. \right. \\ & \quad \left. \left. \left. - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da d\lambda \right|^2 dt \right] \\ & \leq 2\mathbb{E} \left[\int_0^T (h_1(t, \epsilon) + h_2(t, \epsilon)) dt \right], \end{aligned}$$

where

$$\begin{cases} h_1(t, \epsilon) = \left| \left(\nabla_y \tilde{f}^\epsilon(t) - \nabla_y f(t) \right)^2 + \left(\nabla_z \tilde{f}^\epsilon(t) - \nabla_z f(t) \right)^2 \right| (|V_t|^2 + |Z_t^V|^2), \\ h_2(t, \epsilon) = \left| \int_0^1 \int \left(\frac{\delta f}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) \right. \right. \\ \quad \left. \left. - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da d\lambda \right|^2. \end{cases}$$

Set

$$\begin{cases} h_1(t) = 8K^2(|V_t|^2 + |Z_t^V|^2), \\ h_2(t) = 32K^2 + 16K^2(|Y_t^\epsilon - Y_t|^2 + |Z_t^\epsilon - Z_t|^2). \end{cases}$$

By Assumption 2.4 and regularity of probability density, we further have for $i \in \{1, 2\}$ and any $0 < \epsilon < 1$,

$$(3.5) \quad h_i(t, \epsilon) \leq h_i(t) \quad \text{a.e. a.s.}$$

Actually, by a priori estimate (3.3) and Lemma 3.1, for $i = 1, 2$,

$$(3.6) \quad \mathbb{E} \left[\int_0^T h_i(t, \epsilon) dt \right] < \infty.$$

As $\epsilon_n \downarrow 0$, by definition of μ_t^ϵ , we know that $\mu_t^{\epsilon_n}$ converges weakly to μ_t a.e. a.s. Also by Lemma 3.1 it is clear that $(Y_t^{\epsilon_n}, Z_t^{\epsilon_n})$ converges in $S_{\mathcal{F}}^2(0, T; \mathbb{R}^n) \times L_{\mathcal{F}}^2(0, T; \mathbb{R}^{n \times m})$, so there exists a subsequence of $\{\epsilon_n\}$, still denoted by $\{\epsilon_n\}$, such that $(Y_t^{\epsilon_n}, Z_t^{\epsilon_n})$ converges to (Y_t^μ, Z_t^μ) a.e. a.s. By the continuity of $\nabla_y f$ with respect to (y, z) , we have for $0 \leq \lambda \leq 1$,

$$\lim_{n \rightarrow \infty} |\nabla_y f(t, Y_t^\mu + \lambda(Y_t^{\epsilon_n} - Y_t^\mu), Z_t + \lambda(Z_t^{\epsilon_n} - Z_t^\mu), \mu_t) - \nabla_y f(t)| = 0.$$

Hence by the dominated convergence theorem,

$$\begin{aligned} & \lim_{n \rightarrow \infty} |\nabla_y \tilde{f}^{\epsilon_n}(t) - \nabla_y f(t)|^2 \\ & \leq \lim_{n \rightarrow \infty} \int_0^1 |\nabla_y f(t, Y_t^\mu + \lambda(Y_t^{\epsilon_n} - Y_t^\mu), Z_t + \lambda(Z_t^{\epsilon_n} - Z_t^\mu), \mu_t) - \nabla_y f(t)|^2 d\lambda \\ & = 0 \quad \text{a.e. a.s.} \end{aligned}$$

Similarly, $\lim_{n \rightarrow \infty} |\nabla_z \tilde{f}^{\epsilon_n}(t) - \nabla_z f(t)|^2 = 0$.

For the term of control, we notice that

$$\begin{aligned} & \int \left(\frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, (1-\lambda)\mu_t + \lambda\mu_t^{\epsilon_n})(a) - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da \\ & = \int \left(\frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, (1-\lambda)\mu_t + \lambda\mu_t^{\epsilon_n})(a) - \frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da \\ & \quad + \int \left(\frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, \mu_t)(a) - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da. \end{aligned}$$

For $0 \leq \lambda' \leq 1$ and $\mu_t^{\lambda, \lambda'} := (1-\lambda')\mu_t + \lambda'((1-\lambda)\mu_t + \lambda\mu_t^{\epsilon_n}) = \mu_t + \lambda\lambda'(\mu_t^{\epsilon_n} - \mu_t)$,

$$\begin{aligned} & \left| \mathbb{E} \left[\int_0^T \int_0^1 \int \left(\frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, (1-\lambda)\mu_t + \lambda\mu_t^{\epsilon_n})(a) \right. \right. \right. \\ & \quad \left. \left. \left. - \frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, \mu_t)(a) \right) \epsilon_n (\pi_t(a) - \mu_t(a)) da d\lambda dt \right] \right|^2 \\ (3.7) \quad & \leq \mathbb{E} \left[\int_0^T \left| \int_0^1 \int_0^1 \lambda \int \int \frac{\delta^2 f}{\delta m^2}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, \mu_t^{\lambda, \lambda'})(a, a') \epsilon_n (\pi_t(a') - \mu_t(a')) da' \right. \right. \\ & \quad \left. \left. \cdot \epsilon_n (\pi_t(a) - \mu_t(a)) da d\lambda d\lambda' \right|^2 dt \right] \\ & \leq \mathbb{E} \left[\int_0^T \left| \int_0^1 \int_0^1 \lambda K \int (\pi_t(a') + \mu_t(a')) da' \int (\pi_t(a) + \mu_t(a)) da d\lambda d\lambda' \right|^2 dt \right] \epsilon_n^2 \\ & = C_{K, T} \epsilon_n^2 = o(\epsilon_n). \end{aligned}$$

Then, since $\frac{\delta f}{\delta m}$ is uniformly Lipschitz in (y, z) for any fixed $m \in \mathcal{A}$, by Lemma 3.1 we get

$$\left| \mathbb{E} \left[\int_0^T \int_0^1 \int \left(\frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, \mu_t)(a) - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) \epsilon_n (\pi_t(a) - \mu_t(a)) da d\lambda dt \right] \right|^2$$

$$\begin{aligned}
&\leq \left| \mathbb{E} \left[\int_0^T \int_0^1 \int K(|Y_t^{\epsilon_n} - Y_t^\mu| + |Z_t^{\epsilon_n} - Z_t^\mu|) \epsilon_n (\pi_t(a) + \mu_t(a)) da d\lambda dt \right] \right|^2 \\
&\leq C_K \mathbb{E} \left[\int_0^T (|Y_t^{\epsilon_n} - Y_t^\mu|^2 + |Z_t^{\epsilon_n} - Z_t^\mu|^2) dt \right] \epsilon_n^2 = o(\epsilon_n^2),
\end{aligned}$$

which together with (3.7) leads to

$$\begin{aligned}
&\mathbb{E} \left[\int_0^T \int_0^1 \int \left(\frac{\delta f}{\delta m}(t, Y_t^{\epsilon_n}, Z_t^{\epsilon_n}, (1-\lambda)\mu_t + \lambda\mu_t^{\epsilon_n})(a) - \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) \right. \\
&\quad \left. \cdot \epsilon_n (\pi_t(a) - \mu_t(a)) da d\lambda dt \right] \\
&= o(\epsilon_n).
\end{aligned}$$

As a result,

$$(3.8) \quad \lim_{n \rightarrow \infty} |h_1(t, \epsilon_n) + h_2(t, \epsilon_n)| = 0 \quad \text{a.e. a.s.}$$

By (3.5), (3.6), (3.8) and the dominated convergence theorem, we have

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\int_0^T (h_1(t, \epsilon_n) + h_2(t, \epsilon_n)) dt \right] = 0.$$

Hence $\lim_{n \rightarrow \infty} h(\epsilon_n) = 0$ follows from (3.4).

Due to the arbitrariness of ϵ_n and Heine Theorem, it yields that $\lim_{\epsilon \rightarrow 0^+} h(\epsilon) = 0$. \square

The following lemma comes from Lemma 3.2 in [26] and will be used in the variation inequality of cost functional.

LEMMA 3.3. For $\pi, \mu, \gamma \in \mathcal{A}$, set $\mu^\epsilon = \mu + \epsilon(\pi - \mu)$. Then
i) for any $\epsilon \in (0, 1)$,

$$\frac{1}{\epsilon} \int_0^T (Ent(\mu_t^\epsilon | \gamma_t) - Ent(\mu_t | \gamma_t)) dt \geq \int_0^T \int (\ln \mu_t(a) - \ln \gamma_t(a)) (\pi_t(a) - \mu_t(a)) da dt;$$

ii)

$$\begin{aligned}
&\limsup_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \int_0^T (Ent(\mu_t^\epsilon | \gamma_t) - Ent(\mu_t | \gamma_t)) dt \\
&\leq \int_0^T \int (\ln \mu_t(a) - \ln \gamma_t(a)) (\pi_t(a) - \mu_t(a)) da dt.
\end{aligned}$$

Before we prove the variation inequality of cost functional, we give an expression for the Gâteaux derivative of J^0 in terms of the functional derivative of the Hamiltonian.

PROPOSITION 3.4. Under Assumption 2.4,

$$\begin{aligned}
\lim_{\epsilon \rightarrow 0} \frac{J^0(\mu^\epsilon) - J^0(\mu)}{\epsilon} &= \mathbb{E} \left[\nabla_y \phi(Y_0) V_0 + \int_0^T (\nabla_y l(t) V_t + \nabla_z l(t) Z_t^V \right. \\
&\quad \left. + \int \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da dt \right].
\end{aligned}$$

Proof. By the variation equation (3.2), we have

$$J^0(\mu^\epsilon) - J^0(\mu)$$

$$\begin{aligned}
&= \mathbb{E} \left[\int_0^T \left(l(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t^\epsilon) - l(t, Y_t^\mu, Z_t^\mu, \mu_t) \right) dt + \phi(Y_0^\epsilon) - \phi(Y_0^\mu) \right] \\
&= \mathbb{E} \left[\int_0^T \left(\int_0^1 \int \frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) \epsilon (\pi_t(a) - \mu_t(a)) da d\lambda \right. \right. \\
&\quad \left. \left. + \nabla_z \tilde{l}^\epsilon(t)(Z_t^{V^\epsilon} + \epsilon Z_t^V) + \nabla_y \tilde{l}^\epsilon(t)(V_t^\epsilon + \epsilon V_t) \right) dt \right. \\
&\quad \left. + \int_0^1 \nabla_y \phi(Y_0^\mu + \lambda \epsilon (V_0^\epsilon + V_0)) (V_0^\epsilon + \epsilon V_0) d\lambda \right] \\
&= \epsilon \mathbb{E} \left[\int_0^T \left(\nabla_y l(t) V_t + \nabla_z l(t) Z_t^V + \int \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right) dt \right. \\
&\quad \left. + \nabla_y \phi(Y_0^\mu) V_0 \right] + \mathbb{E} \left[\int_0^T (\nabla_y l(t) V_t^\epsilon + \nabla_z l(t) Z_t^{V^\epsilon}) dt \right] \\
&\quad + \mathbb{E} \left[\int_0^T \left((\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t))(V_t^\epsilon + \epsilon V_t) + (\nabla_z \tilde{l}^\epsilon(t) - \nabla_z l(t))(Z_t^{V^\epsilon} + \epsilon Z_t^V) \right) \right. \\
&\quad \left. + \mathbb{E} \left[\int_0^T \int_0^1 \int \left(\frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) - \frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t)(a) \right) \right. \right. \\
&\quad \left. \left. \cdot \epsilon (\pi_t(a) - \mu_t(a)) da d\lambda dt \right] \right. \\
&\quad \left. + \mathbb{E} \left[\int_0^T \int_0^1 \int \left(\frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t)(a) - \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) \right) \epsilon (\pi_t(a) - \mu_t(a)) da d\lambda dt \right] \right. \\
&\quad \left. + \mathbb{E} \left[\nabla_y \phi(Y_0^\mu) V_0^\epsilon + \int_0^1 (\nabla_y \phi(Y_0^\mu + \lambda(V_0^\epsilon + \epsilon V_0)) - \nabla_y \phi(Y_0^\mu)) (V_0^\epsilon + \epsilon V_0) d\lambda \right] \right].
\end{aligned}
\tag{3.9}$$

First by Assumption 2.4 and Lemma 3.2, we know

$$\begin{aligned}
\left| \mathbb{E} \left[\nabla_y \phi(Y_0^\mu) V_0^\epsilon \right] \right|^2 &\leq \mathbb{E} \left[|\nabla_y \phi(Y_0^\mu)|^2 \right] \cdot \mathbb{E} \left[|V_0^\epsilon|^2 \right] \leq 2K^2 \mathbb{E} \left[1 + |Y_0^\mu|^2 \right] \mathbb{E} \left[\sup_{0 \leq t \leq T} |V_t^\epsilon|^2 \right] \\
&\leq 2K^2 \left(1 + \mathbb{E} \left[\sup_{0 \leq t \leq T} |Y_t^\mu|^2 \right] \right) \mathbb{E} \left[\sup_{0 \leq t \leq T} |V_t^\epsilon|^2 \right] = o(\epsilon^2)
\end{aligned}$$

and

$$\begin{aligned}
\left| \mathbb{E} \left[\int_0^T \nabla_y l(t) V_t^\epsilon dt \right] \right|^2 &\leq \mathbb{E} \left[\int_0^T |\nabla_y l(t)|^2 dt \right] \cdot \mathbb{E} \left[\int_0^T |V_t^\epsilon|^2 dt \right] \\
&\leq 3K^2 \mathbb{E} \left[\int_0^T (1 + |Y_t^\mu|^2 + |Z_t^\mu|^2) dt \right] \cdot \mathbb{E} \left[\int_0^T |V_t^\epsilon|^2 dt \right] \\
&\leq C_{\mu, T, K} \mathbb{E} \left[\int_0^T |V_t^\epsilon|^2 dt \right] = o(\epsilon^2).
\end{aligned}$$

Similarly, $\left| \mathbb{E} \left[\int_0^T \nabla_z l(t) Z_t^{V^\epsilon} dt \right] \right|^2 = o(\epsilon^2)$.

Hence

$$\begin{cases} \mathbb{E}[\nabla_y \phi(Y_0^\mu) V_0^\epsilon] = o(\epsilon), \\ \mathbb{E}[\int_0^T \nabla_y l(t) V_t^\epsilon dt] = o(\epsilon), \\ \mathbb{E}[\int_0^T \nabla_z l(t) Z_t^{V^\epsilon} dt] = o(\epsilon). \end{cases}$$

On the other hand, by Lemma 3.1 we have

$$\begin{aligned}
(3.10) \quad & \left| \mathbb{E} \left[\int_0^T \left(\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t) \right) (V_t^\epsilon + \epsilon V_t) dt \right] \right|^2 \\
&= \left| \mathbb{E} \left[\int_0^T \left(\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t) \right) (Y_t^\epsilon - Y_t^\mu) dt \right] \right|^2 \\
&\leq \mathbb{E} \left[\int_0^T |\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t)|^2 dt \right] \cdot \mathbb{E} \left[\int_0^T |Y_t^\epsilon - Y_t^\mu|^2 dt \right] \\
&\leq C\epsilon^2 \cdot \mathbb{E} \left[\int_0^T |\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t)|^2 dt \right].
\end{aligned}$$

Then we prove $\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\int_0^T |\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t)|^2 dt \right] = 0$. For any $0 \leq \lambda \leq 1$, Set

$$r_y(t, \epsilon; \lambda) := |\nabla_y l(t, Y_t^\mu + \lambda(Y_t^\epsilon - Y_t^\mu), Z_t^\mu + \lambda(Z_t^\epsilon - Z_t^\mu), \mu_t) - \nabla_y l(t)|^2.$$

Similar to the proof of Lemma 3.2, we get from Lemma 3.1 that there exists a subsequence of $\{\epsilon_n\}$, still denoted by $\{\epsilon_n\}$, such that $(Y_t^{\epsilon_n}, Z_t^{\epsilon_n})$ converges to (Y_t^μ, Z_t^μ) a.e. a.s. By the continuity of $\nabla_y l$ with respect to (y, z) , we have for any $0 \leq \lambda \leq 1$,

$$\lim_{n \rightarrow \infty} r_y(t, \epsilon_n; \lambda) = 0 \quad \text{a.e. a.s.}$$

Moreover, as n large enough,

$$\begin{aligned}
& \mathbb{E} \left[\int_0^T r_y(t, \epsilon_n; \lambda) \right] \\
&\leq 24K^2 \mathbb{E} \left[\int_0^T \left((1 + |Y_t^\mu|^2 + |Z_t^\mu|^2) + 16K^2(|Y_t^{\epsilon_n} - Y_t^\mu|^2 + |Z_t^{\epsilon_n} - Z_t^\mu|^2) \right) dt \right] \\
&\leq C \mathbb{E} \left[\int_0^T (1 + |Y_t^\mu|^2 + |Z_t^\mu|^2) dt \right] < \infty.
\end{aligned}$$

Notice

$$\begin{aligned}
& \mathbb{E} \left[\int_0^T |\nabla_y \tilde{l}^{\epsilon_n}(t) - \nabla_y l(t)|^2 dt \right] \\
&= \mathbb{E} \left[\int_0^T \left| \int_0^1 \left(\nabla_y l(t, Y_t^\mu + \lambda(Y_t^{\epsilon_n} - Y_t^\mu), Z_t^\mu + \lambda(Z_t^{\epsilon_n} - Z_t^\mu), \mu_t) - \nabla_y l(t) \right) d\lambda \right|^2 dt \right] \\
&\leq \mathbb{E} \left[\int_0^T \int_0^1 r_y(t, \epsilon_n; \lambda) d\lambda \right].
\end{aligned}$$

By the dominated convergence theorem again, we have

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\int_0^T |\nabla_y \tilde{l}^{\epsilon_n}(t) - \nabla_y l(t)|^2 dt \right] = 0.$$

The arbitrariness of $\{\epsilon_n\}$ and Heine theorem leads to

$$\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\int_0^T |\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t)|^2 dt \right] = 0.$$

Back to (3.10), the above deductions imply

$$\mathbb{E}\left[\int_0^T \left(\nabla_y \tilde{l}^\epsilon(t) - \nabla_y l(t)\right) (V_t^\epsilon + \epsilon V_t) dt\right] = o(\epsilon).$$

By similar deductions, we also have

$$(3.11) \quad \begin{cases} \mathbb{E}\left[\int_0^T \left(\nabla_z \tilde{l}^\epsilon(t) - \nabla_z l(t)\right) (Z_t^{V^\epsilon} + \epsilon Z_t^V) dt\right] = o(\epsilon), \\ \mathbb{E}\left[\int_0^1 \left(\nabla_y \phi(Y_0^\mu + \lambda \epsilon (V_0^\epsilon + V_0)) - \nabla_y \phi(Y_0^\mu)\right) (V_0^\epsilon + \epsilon V_0) d\lambda\right] = o(\epsilon). \end{cases}$$

Then we deal with two terms involving $\frac{\delta l}{\delta m}$. For the first one, it turns out that

$$(3.12) \quad \begin{aligned} & \left| \mathbb{E}\left[\int_0^T \int_0^1 \int \left(\frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, (1-\lambda)\mu_t + \lambda\mu_t^\epsilon)(a) \right. \right. \right. \\ & \quad \left. \left. \left. - \frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t)(a)\right) \epsilon (\pi_t(a) - \mu_t(a)) da d\lambda dt\right] \right|^2 \\ & \leq \mathbb{E}\left[\int_0^T \left| \int_0^1 \int_0^1 \lambda \int \int \frac{\delta^2 l}{\delta m^2}(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t^{\lambda, \lambda'})(a, a') \epsilon (\pi_t(a') - \mu_t(a')) da' \right. \right. \\ & \quad \left. \left. \cdot \epsilon (\pi_t(a) - \mu_t(a)) da d\lambda d\lambda' \right|^2 dt\right] \\ & \leq \mathbb{E}\left[\int_0^T \left| \int_0^1 \int_0^1 \lambda K \int (\pi_t(a') + \mu_t(a')) da' \int (\pi_t(a) + \mu_t(a)) da d\lambda d\lambda' \right|^2 dt\right] \epsilon^2 \\ & = C_{K,T} \epsilon^2 = o(\epsilon). \end{aligned}$$

For the other term, since $\frac{\delta l}{\delta m}$ is uniformly Lipschitz in (y, z) , by Lemma 3.1 we have

$$(3.13) \quad \begin{aligned} & \left| \mathbb{E}\left[\int_0^T \int_0^1 \int \left(\frac{\delta l}{\delta m}(t, Y_t^\epsilon, Z_t^\epsilon, \mu_t)(a) - \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a)\right) \epsilon (\pi_t(a) - \mu_t(a)) da d\lambda dt\right] \right|^2 \\ & \leq \left| \mathbb{E}\left[\int_0^T \int_0^1 \int K (|Y_t^\epsilon - Y_t^\mu| + |Z_t^\epsilon - Z_t^\mu|) \epsilon (\pi_t(a) + \mu_t(a)) da d\lambda dt\right] \right|^2 \\ & \leq C_K \mathbb{E}\left[\int_0^T (|Y_t^\epsilon - Y_t^\mu|^2 + |Z_t^\epsilon - Z_t^\mu|^2) dt\right] \epsilon^2 = o(\epsilon). \end{aligned}$$

Therefore, Proposition 3.4 follows from (3.9) and (3.11)–(3.13). \square

Now we are ready to present and prove the variation equality of cost functional.

PROPOSITION 3.5. *Under Assumption 2.4,*

$$\begin{aligned} \mathbb{E}\left[\nabla_y \phi(Y_0) V_0 + \int_0^T \left(\nabla_y l(t) V_t + \nabla_z l(t) Z_t^V + \int \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right. \right. \\ \left. \left. + \frac{\sigma^2}{2} \int ((\ln \mu_t(a) + U(a)) (\pi_t(a) - \mu_t(a))) da\right) dt\right] \geq 0. \end{aligned}$$

where μ is the optimal control for **(P1)**.

Proof. From the optimality of μ , based on Lemma 3.3 and Proposition 3.4 we have

$$0 \leq \limsup_{\epsilon \rightarrow 0^+} \frac{J^\sigma(\mu^\epsilon) - J^\sigma(\mu)}{\epsilon}$$

$$\begin{aligned}
&= \limsup_{\epsilon \rightarrow 0^+} \left(\frac{J^0(\mu^\epsilon) - J^0(\mu)}{\epsilon} + \frac{1}{\epsilon} \mathbb{E} \left[\int_0^T (\text{Ent}(\mu_t^\epsilon | U) - \text{Ent}(\mu_t | U)) dt \right] \right) \\
&\leq \mathbb{E} \left[\nabla_y \phi(Y_0) V_0 + \int_0^T \left(\nabla_y l(t) V_t + \nabla_z l(t) Z_t^V \right. \right. \\
&\quad \left. \left. + \int \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right. \right. \\
&\quad \left. \left. + \frac{\sigma^2}{2} \int ((\ln \mu_t(a) + U(a)) (\pi_t(a) - \mu_t(a))) da \right) dt \right]. \quad \square
\end{aligned}$$

Now we use duality technique to derive the local necessary condition for optimal control. To begin with, we introduce the adjoint equation of variation equation (3.2)

$$(3.14) \quad \begin{cases} dP_t^\mu = -(-\nabla_y f^\top(t) P_t^\mu + \nabla_y l(t)) dt - (-\nabla_z f^\top(t) P_t^\mu + \nabla_z l(t)) dW_t \\ P_0^\mu = -\nabla_y \phi(Y_0^\mu). \end{cases}$$

Actually, SDE (3.14) is a linear equation whose coefficients $\nabla_y f^\top(t)$ and $\nabla_z f^\top(t)$ are duality operators of $\nabla_y f(t)$ and $\nabla_z f(t)$ satisfying Lipschitz conditions, so it is clear that SDE (3.14) has a unique solution in $S_{\mathcal{F}}^2(0, T; \mathbb{R}^n)$. Then introduce the Hamiltonians $H^0 : \Omega \times [0, T] \times \mathbb{R}^n \times \mathbb{R}^{n \times m} \times \mathbb{R}^n \times \mathcal{P}(\mathbb{R}^p) \rightarrow \mathbb{R}$ and $H^\sigma : \Omega \times [0, T] \times \mathbb{R}^n \times \mathbb{R}^{n \times m} \times \mathbb{R}^n \times \mathcal{P}(\mathbb{R}^p) \times \mathcal{P}(\mathbb{R}^p) \rightarrow \mathbb{R}$ as follows

$$\begin{aligned}
H^0(t, y, z, p, m) &:= -pf(t, y, z, m) + l(t, y, z, m), \\
H^\sigma(t, y, z, p, m, m') &:= H^0(t, y, z, p, m) + \frac{\sigma^2}{2} \text{Ent}(m | m').
\end{aligned}$$

It can be seen from Assumption 2.4 that $H^0(t, y, z, p, m)$ and $H^\sigma(t, y, z, p, m, m')$ are continuous with respect to (y, z, p) and differentiable with respect to (y, z) . By Definition 2.3 $H^0(t, y, z, p, m)$ has flat derivative

$$\frac{\delta H^0(t, y, z, p, m)}{\delta m} := \frac{\delta H^0}{\delta m}(t, y, z, p, m)(a).$$

Hence the adjoint equation (3.14) is equivalent to

$$\begin{cases} dP_t^\mu = -\nabla_y H^0(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t) dt - \nabla_z H^0(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t) dW_t \\ P_0 = -\nabla_y \phi(Y_0). \end{cases}$$

Although the term of entropy is lower-semi continuous and does not have flat derivative, we still use the following notation for convenience

$$\frac{\delta H^\sigma}{\delta m}(t, y, z, p, m)(a) := \frac{\delta H^0}{\delta m}(t, y, z, p, m)(a) + \frac{\sigma^2}{2} (U(a) + \ln m(a) + 1).$$

Based on (3.2) and adjoint equation (3.14), the variation inequality of cost functional in Proposition 3.4 can be written in a new form.

PROPOSITION 3.6. *Under Assumption 2.4,*

$$\lim_{\epsilon \rightarrow 0^+} \frac{J^0(\mu^\epsilon) - J^0(\mu)}{\epsilon} = \mathbb{E} \left[\int_0^T \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) \right) (\pi_t(a) - \mu_t(a)) da dt \right].$$

Proof. Applying Itô formula to $P_t^\mu V_t$, where V is the solution to variation equation (3.2), we have

$$\begin{aligned}
& dP_t^\mu V_t \\
&= V_t dP_t^\mu + P_t^\mu dV_t + dP_t^\mu dV_t \\
&= \left(-V_t (-\nabla_y f(t) P_t^\mu + \nabla_y l(t)) - P_t^\mu \left(\nabla_y f(t) V_t + \nabla_z f(t) Z_t^V \right) \right. \\
&\quad \left. - Z_t^V (-\nabla_z f(t) P_t^\mu + \nabla_z l(t)) - P_t^\mu \int \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right) dt \\
(3.15) \quad & + \left(P_t^\mu Z_t^V - V_t (-\nabla_z f(t) P_t^\mu + \nabla_z l(t)) \right) dW_t.
\end{aligned}$$

Taking expectation on both sides of (3.15) and using standard stopping time arguments, we obtain

$$(3.16) \quad \mathbb{E} \left[\nabla_y \phi(Y_0^\mu) V_0 + \int_0^T \left(\nabla_y l(t) V_t + \nabla_z l(t) Z_t^V \right. \right. \\
\left. \left. + P_t^\mu \int \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right) dt \right] = 0.$$

Then, based on Proposition 3.4 and (3.16), we have

$$\begin{aligned}
& \lim_{\epsilon \rightarrow 0^+} \frac{J^0(\mu^\epsilon) - J^0(\mu)}{\epsilon} \\
&= \mathbb{E} \left[\int_0^T \left(\nabla_y l(t) V_t + \nabla_z l(t) Z_t^V \right. \right. \\
&\quad \left. \left. + \int \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right) dt + \nabla_y \phi(Y_0^\mu) V_0 \right] \\
&= \mathbb{E} \left[\int_0^T \left(-P_t^\mu \int \frac{\delta f}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right. \right. \\
&\quad \left. \left. + \int \frac{\delta l}{\delta m}(t, Y_t^\mu, Z_t^\mu, \mu_t)(a) (\pi_t(a) - \mu_t(a)) da \right) dt \right] \\
&= \mathbb{E} \left[\int_0^T \int \frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) \cdot (\pi_t(a) - \mu_t(a)) dadt \right]. \quad \square
\end{aligned}$$

Now we are ready to prove the maximum principle for Problem (P1).

THEOREM 3.7. *Under Assumption 2.4, if $\mu \in \mathcal{A}$ is an optimal control of Problem (P1) with the corresponding optimal state process (Y^μ, Z^μ) , and P^μ is the solution to adjoint equation (3.14), then for any $t \in [0, T]$, $\pi \in \mathcal{A}$,*

$$(3.17) \quad \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) + \frac{\sigma^2}{2} (\ln \mu_t(a) + U(a)) \right) (\pi_t(a) - \mu_t(a)) da \geq 0 \quad \text{a.s.}$$

Proof. Based on Propositions 3.5 and 3.6, we deduce from the optimality of μ that

$$\mathbb{E} \left[\int_0^T \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) + \frac{\sigma^2}{2} (\ln \mu_t(a) + U(a)) \right) (\pi_t(a) - \mu_t(a)) dadt \right] \geq 0.$$

Assume that (3.17) doesn't hold. This means that there is a $\tilde{\pi} \in \mathcal{A}$ and $S_\epsilon \in \mathcal{F} \otimes \mathcal{B}([0, T])$ with a strictly positive measure $\mathbb{P} \otimes \Lambda$, where Λ is the Lebesgue measure on $\mathcal{B}([0, T])$ and

$$S_\epsilon = \left\{ (\omega, t) : \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) + \frac{\sigma^2}{2} (\ln \mu_t(a) + U(a)) \right) \cdot (\tilde{\pi}_t(a) - \mu_t(a)) da \leq -\epsilon < 0 \right\}.$$

Define $\tilde{\mu}_t := \tilde{\pi}_t \mathbb{I}_{S_\epsilon} + \mu_t \mathbb{I}_{S_\epsilon^c}$. We have

$$\begin{aligned} 0 &\leq \mathbb{E} \left[\int_0^T \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) + \frac{\sigma^2}{2} (\ln \mu_t(a) + U(a)) \right) (\tilde{\mu}_t(a) - \mu_t(a)) dadt \right] \\ &= \mathbb{E} \left[\int_0^T \mathbb{I}_{S_\epsilon} \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) + \frac{\sigma^2}{2} (\ln \mu_t(a) + U(a)) \right) \cdot (\tilde{\pi}_t(a) - \mu_t(a)) dadt \right] \\ &\leq -\epsilon \mathbb{E} \int_0^T \mathbb{I}_{S_\epsilon} dt < 0, \end{aligned}$$

which leads to a contradiction. Then the proof follows. \square

4. Further Discussions of Optimal Controls. In this section we present a sufficient condition for the optimal control and give an implicit form of it. For the sufficient condition, the convex conditions for the coefficients are needed.

ASSUMPTION 4.1. *For the coefficients in Problem (P1), $\phi(y)$ is convex in y and $H^\sigma(t, y, z, p, m, m')$ is convex in (y, z, m) .*

With Assumption 4.1 we prove the sufficient condition for the optimal control.

THEOREM 4.2. *Under Assumptions 2.4 and 4.1, the control $\mu \in \mathcal{A}$ is an optimal control of Problem (P1) if for any $t \in [0, T]$, $\pi \in \mathcal{A}$, it satisfies (3.17) with (Y^μ, Z^μ) and P^μ be the solutions to the corresponding BSDE (1.1) and adjoint equation (3.14), respectively.*

Proof. For $\mu \in \mathcal{A}$ satisfying (3.17), we have

$$(4.1) \quad J^\sigma(\pi) - J^\sigma(\mu) = I_1 + I_2,$$

where $I_1 = \mathbb{E} \left[\phi(Y_0^\pi) - \phi(Y_0^\mu) \right]$ and

$$(4.2) \quad I_2 = \mathbb{E} \left[\int_0^T [l(t, Y_t^\pi, Z_t^\pi, \pi_t) + \frac{\sigma^2}{2} Ent(\pi_t | e^{-U}) - l(t, Y_t^\mu, Z_t^\mu, \mu_t) - \frac{\sigma^2}{2} Ent(\mu_t | e^{-U})] dt \right].$$

For I_1 , applying Itô formula to $P_t^\mu(Y_t^\pi - Y_t^\mu)$, we have

$$\begin{aligned} &dP_t^\mu(Y_t^\pi - Y_t^\mu) \\ &= \left(- (Y_t^\pi - Y_t^\mu)(-\nabla_y f(t)P_t^\mu + \nabla_y l(t)) - (Z_t^\pi - Z_t^\mu)(-\nabla_z f(t)P_t^\mu + \nabla_z l(t)) \right. \\ &\quad \left. - P_t^\mu(f(t, Y_t^\mu, Z_t^\mu, \mu_t) - f(t, Y_t^\pi, Z_t^\pi, \pi_t)) \right) dt \\ &\quad + \left(P_t^\mu(Z_t^\pi - Z_t^\mu) - (Y_t^\pi - Y_t^\mu)(-\nabla_z f(t)P_t^\mu + \nabla_z l(t)) \right) dW_t, \end{aligned}$$

so by the convexity of ϕ it yields that

$$\begin{aligned}
(4.3) \quad I_1 &\geq \mathbb{E} \left[\nabla_y \phi(Y_0^\mu)(Y_0^\pi - Y_0^\mu) \right] \\
&= -\mathbb{E} \left[P_0^\mu(Y_0^\pi - Y_0^\mu) \right] \\
&= -\mathbb{E} \left[\int_0^T (Y_t^\pi - Y_t^\mu) \nabla_y H^0(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t) dt \right] \\
&\quad - \mathbb{E} \left[\int_0^T (Z_t^\pi - Z_t^\mu) \nabla_z H^0(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t) dt \right] \\
&\quad - \mathbb{E} \left[\int_0^T P_t^\mu (f(t, Y_t^\mu, Z_t^\mu, \mu_t) - f(t, Y_t^\pi, Z_t^\pi, \pi_t)) dt \right].
\end{aligned}$$

Hence by (4.1)–(4.3) and the convexity of H^σ we have

$$\begin{aligned}
&J^\sigma(\pi) - J^\sigma(\mu) \\
&\geq \mathbb{E} \left[\int_0^T \left(H^\sigma(t, Y_t^\pi, Z_t^\pi, P_t^\pi, \pi_t, e^{-U}) - H^\sigma(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t, e^{-U}) \right) dt \right] \\
&\quad - \mathbb{E} \left[\int_0^T (Y_t^\pi - Y_t^\mu) \nabla_y H^0(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t) dt \right] \\
&\quad - \mathbb{E} \left[\int_0^T (Z_t^\pi - Z_t^\mu) \nabla_z H^0(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t) dt \right] \\
&\geq \mathbb{E} \left[\int_0^T \int \left(\frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a) + \frac{\sigma^2}{2} (\ln \mu_t + U(a)) \right) (\pi_t(a) - \mu_t(a)) da dt \right].
\end{aligned}$$

Therefore, by the condition (3.17), it follows that $J^\sigma(\pi) - J^\sigma(\mu) \geq 0$ for any $\pi \in \mathcal{A}$, which implies that μ is the optimal control. \square

We then give an implicit form of optimal control of BSDE (1.1) with the cost functional (1.2). Assume that an optimal control μ exists. For fixed $t \in [0, T]$ and $\omega \in \Omega$, (3.17) in Theorems 3.7 and 4.2 is equivalent to

$$(4.4) \quad \mu_t \in \arg \min_{m \in \mathcal{P}_2(\mathbb{R}^p)} H^\sigma(t, Y_t^\mu, Z_t^\mu, P_t^\mu, m, e^{-U}),$$

where (Y^μ, Z^μ) and P^μ are the solutions to BSDE (1.1) and the adjoint equation (3.14), respectively, with the control variable μ .

Assume that $U \in C^\infty$, $\nabla_a U$ is Lipschitz-continuous, and there exists constants $C_U > 0$ and $C'_U \in \mathbb{R}$ such that for any $a \in \mathbb{R}^p$, it holds that $\nabla_a U(a) \cdot a \geq C_U |a|^2 + C'_U$. According to Proposition 2.5 in Hu, Ren, Šiška and Szpruch [14], in which the deterministic control system is studied, the admissible control set of the optimization problem (4.4) can be enlarged to $\mathcal{P}(\mathbb{R}^p)$ with above assumptions on U . So (4.4) is equivalent to a constrained optimization problem in these settings, i.e., for $t \in [0, T]$,

$$(4.5) \quad \mu_t \in \arg \min_{m \in \mathcal{M}(\mathbb{R}^p)} -pf(t, Y_t^\mu, Z_t^\mu, P_t^\mu, m) + l(t, Y_t^\mu, Z_t^\mu, P_t^\mu, m) + \frac{\sigma^2}{2} Ent(m | e^{-U}),$$

with a constraint $\int m(a) da = 1$.

By Lagrange multiplier method, we further transform this optimization problem into an equivalent optimization problem without constraint. For this, we introduce

the Lagrange function $L : \mathcal{M}(\mathbb{R}^p) \times \mathbb{R} \rightarrow \mathbb{R}$ with the Lagrange multiplier β as below

$$L(m, \beta) = H^\sigma(t, Y_t^\mu, Z_t^\mu, P_t^\mu, m, e^{-U}) + \beta \left(\int m(a) da - 1 \right).$$

Then we define the Lagrange duality function $G(\beta) = \min_{m \in \mathcal{M}(\mathbb{R}^p)} L(m, \beta)$. By the weak duality theory, we know $G(\beta) \leq \min_{m \in \mathcal{P}(\mathbb{R}^p)} H^\sigma(t, Y_t, Z_t, P_t, m, e^{-U})$. Hence the goal now is to solve

$$(4.6) \quad \beta^* \in \arg \max_{\beta \in \mathbb{R}} G(\beta) = \arg \max_{\beta \in \mathbb{R}} \min_{m \in \mathcal{M}(\mathbb{R}^p)} L(m, \beta).$$

With Assumption 4.1, (4.5) is a convex optimization problem and satisfies Slater's condition of convex optimization theory, which leads to the strong duality of (4.5) and (4.6):

$$\max_{\beta \in \mathbb{R}} \min_{m \in \mathcal{M}(\mathbb{R}^p)} L(m, \beta) = \min_{m \in \mathcal{P}(\mathbb{R}^p)} H^\sigma(t, Y_t^\mu, Z_t^\mu, P_t^\mu, m, e^{-U}).$$

By the first order condition for the flat derivative of H^0 we have

$$\begin{cases} \frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, m)(a) + \frac{\sigma^2}{2}(U(a) + \ln m(a) + 1) + \beta = 0, \\ \int m(a) da = 1. \end{cases}$$

By solving above equation, we know that the control μ_t , $t \in [0, T]$, is a fixed point of the following equation

$$(4.7) \quad \mu_t(a) = \frac{e^{-U(a) - \frac{2}{\sigma^2} \frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a)}}{\int e^{-U(a) - \frac{2}{\sigma^2} \frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a)} da},$$

and the corresponding Lagrange multiplier

$$(4.8) \quad \beta = \frac{\sigma^2}{2} \left(\ln \left(\int e^{-U(a) - \frac{2}{\sigma^2} \frac{\delta H^0}{\delta m}(t, Y_t^\mu, Z_t^\mu, P_t^\mu, \mu_t)(a)} da \right) - 1 \right).$$

Note that here $\mu_t \in \mathcal{P}(\mathbb{R}^p)$, so it is not a solution to Problem (P1) unless $\mu_t \in \mathcal{P}_2(\mathbb{R}^p)$ and $\mathbb{E} \left[\int_0^T Ent(\mu_t | e^{-U}) dadt \right] < \infty$. Actually, according to [14] we further know from the assumption on U that there exist constants C' and C satisfying $0 \leq C' \leq C$ such that for any $a \in \mathbb{R}^p$,

$$C'|a|^2 - C \leq U(a) \leq C(1 + |a|^2).$$

So if μ in (4.7) satisfies $\mathbb{E} \left[\int_0^T Ent(\mu_t | e^{-U}) dadt \right] < \infty$, $Ent(\mu_t | e^{-U}) < \infty$ a.e. a.s., which leads to $\int |a|^2 \mu_t(a) da \leq \int U(a) \mu_t(a) da < \infty$ a.e. a.s., i.e. $\mu_t \in \mathcal{P}_2(\mathbb{R}^p)$.

Without a specific form of H^0 , the above discussion for the existence of an optimal control is based on some assumptions and the optimal control remains implicit. We would give an explicit form of μ in (4.7) and prove that this μ is exactly the optimal control in the LQ case. One can refer to Proposition 2.5 in [14] for more discussions about similar formulations as (4.7) in deterministic cases.

To end this section, let's see a relaxed optimal control problem of BSDE, which is a special case of entropy regularized control problem. For $f : \Omega \times [0, T] \times \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^n$, $\xi \in: \Omega \rightarrow \mathbb{R}^n$ and an admissible control $a \in \mathcal{U}_{ad} = L^2_{\mathcal{F}}(0, T; \mathbb{R}^p)$, consider the controlled BSDE

$$(4.9) \quad \begin{cases} -dy_t^a = f(t, y_t^a, z_t^a, a_t)dt - z_t^a dW_t, \\ y_t = \xi. \end{cases}$$

By law of large numbers we have the exploratory BSDE

$$(4.10) \quad \begin{cases} -d\tilde{y}_t^\pi = \tilde{f}(t, \tilde{y}_t^\pi, \tilde{z}_t^\pi, \boldsymbol{\pi}_t)dt - \tilde{z}_t^\pi dW_t, \\ y_t = \xi, \end{cases}$$

where $\boldsymbol{\pi}$ is the distribution of control, $(\tilde{y}^\pi, \tilde{z}^\pi)$ is the exploratory state variable and $\tilde{f}(t, \tilde{y}_t^\pi, \tilde{z}_t^\pi, \boldsymbol{\pi}_t) = \int f(t, \tilde{y}_t^\pi, \tilde{z}_t^\pi, a) \boldsymbol{\pi}_t(a) da$. To see how to get BSDE (4.10), let's set (y^i, z^i) to be the copy of the path generated from the dynamics (4.9) with the control a^i sampled independently under this policy $\boldsymbol{\pi}$. For any $0 \leq t \leq T$, we have

$$\Delta y_t^i \equiv y_{t+\Delta t}^i - y_t^i \approx -f(t, y_t^i, z_t^i, a_t^i) \Delta t + z_t^i (W_{t+\Delta t}^i - W_t^i).$$

Here each $y^i, i = 1, 2, \dots, N$, can be viewed as a copy of an independent sample from \tilde{y} . It then follows from the law of large numbers that, as $N \rightarrow \infty$,

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \Delta y_t^i &\approx -\frac{1}{N} \sum_{i=1}^N f(t, y_t^i, z_t^i, a_t^i) \Delta t + \frac{1}{N} \sum_{i=1}^N z_t^i (W_{t+\Delta t}^i - W_t^i) \\ &\xrightarrow{\text{a.s.}} \mathbb{E} \left[\int -f(t, \tilde{y}_t, \tilde{z}_t, a) \boldsymbol{\pi}_t(a) da \Delta t \right] + \mathbb{E} \left[\int z_t^i \boldsymbol{\pi}_t(a) da \right] \mathbb{E} [W_{t+\Delta t} - W_t] \\ &= \mathbb{E} \left[\int -f(t, \tilde{y}_t, \tilde{z}_t, a) \boldsymbol{\pi}_t(a) da \Delta t \right]. \end{aligned}$$

In the above deduction, we have assumed that both $\boldsymbol{\pi}$ and \tilde{z} are identically distributed over $[t, t + \Delta t]$ and independent of the increments of the sample paths of W . Then the entropy-regularized cost function appears as

$$(4.11) \quad J^\sigma(T, \xi; \boldsymbol{\pi}) = \mathbb{E} \left[\int_0^T \left(\int l(t, \tilde{y}_t, \tilde{z}_t, a) \boldsymbol{\pi}_t(a) da + \frac{\sigma^2}{2} \text{Ent}(\boldsymbol{\pi}_t | e^{-U}) \right) dt + \phi(\tilde{y}(0)) \right].$$

By Definition 2.3,

$$\frac{\tilde{f}(t, \tilde{y}_t, \tilde{z}_t, \boldsymbol{\pi}(a))}{\delta m} = f(t, \tilde{y}_t, \tilde{z}_t, a) \quad \text{and} \quad \frac{\tilde{l}(t, \tilde{y}_t, \tilde{z}_t, \boldsymbol{\pi}(a))}{\delta m} = l(t, \tilde{y}_t, \tilde{z}_t, a).$$

Hence

$$\frac{\delta H^0}{\delta m}(t, \tilde{y}_t, \tilde{z}_t, \tilde{p}_t, \boldsymbol{\pi}_t)(a) = -\tilde{p}_t f(t, \tilde{y}_t, \tilde{z}_t, a) + l(t, \tilde{y}_t, \tilde{z}_t, a),$$

where \tilde{p} is the solution to the corresponding adjoint equation, and a candidate optimal control for the cost functional (4.11) is

$$\mu_t(a) = \frac{e^{-U(a) - \frac{\sigma^2}{2} [-p_t f(t, \tilde{y}^\mu, \tilde{z}^\mu, a) + l(t, \tilde{y}^\mu, \tilde{z}^\mu, a)]}}{\int e^{-U(a) - \frac{\sigma^2}{2} [-p_t f(t, \tilde{y}^\mu, \tilde{z}^\mu, a) + l(t, \tilde{y}^\mu, \tilde{z}^\mu, a)]} da}.$$

5. Backward Stochastic Linear-Quadratic Control System with Entropy Regularization. Let \mathbb{S}^n be the set of all $n \times n$ symmetric matrices, \mathbb{S}_+^n be the set of all $n \times n$ positive semi-definite matrices, $\hat{\mathbb{S}}_+^n$ be the set of all $n \times n$ positive definite matrices, and \mathbb{I}_n be the $n \times n$ identity matrix. For $t \in [0, T]$, $\pi \in \mathcal{A}$, consider a linear controlled BSDE

$$(5.1) \quad \begin{cases} dY_t^\pi = (A_t Y_t^\pi + B_t \int a \pi_t(a) da + C_t Z_t^\pi) dt + Z_t^\pi dW_t, \\ Y_T = \xi \end{cases}$$

and its cost functional

$$(5.2) \quad J^\sigma(\pi) = \frac{1}{2} \mathbb{E} \left[\int_0^T \left(Y_t^\pi H_t Y_t^\pi + \int a R_t a \pi_t(a) da + Z_t^\pi N_t Z_t^\pi + \sigma^2 \text{Ent}(\pi_t | e^{-U}) \right) dt + Y_0^\pi G Y_0^\pi \right].$$

Then LQ problem is

(P2): to find an optimal $\mu \in \mathcal{A}$ such that

$$J^\sigma(\mu) = \inf_{\pi \in \mathcal{A}} J^\sigma(\pi).$$

We give the assumptions for the coefficients of LQ problem.

ASSUMPTION 5.1. (i) $\xi \in L^2_{\mathcal{F}_T}(\Omega; \mathbb{R}^n)$, $A, C \in L^\infty(0, T; \mathbb{R}^{n \times n})$ and $B \in L^\infty(0, T; \mathbb{R}^{n \times p})$.
(ii) $H, N \in L^\infty(0, T; \mathbb{S}_+^n)$, $R \in L^\infty(0, T; \mathbb{S}_+^p)$ and $G \in \mathbb{S}_+^n$.

Assumption 5.1 guarantees that linear BSDE (5.1) has a unique solution $(Y^\pi, Z^\pi) \in S^2_{\mathcal{F}}(0, T; \mathbb{R}^n) \times L^2_{\mathcal{F}}(0, T; \mathbb{R}^n)$.

From Theorem 3.7, (4.7) and (4.8), the necessary condition of optimality in LQ case follows.

THEOREM 5.2. Under Assumption 5.1, if $\mu \in \mathcal{A}$ is an optimal control of Problem (P2) with corresponding optimal state process (Y^μ, Z^μ) , then

$$\begin{cases} dP_t^\mu = -(A_t P_t^\mu + H_t Y_t^\mu) dt - (C_t P_t^\mu + N_t Z_t^\mu) dW_t, \\ P_0^\mu = -G Y_0^\mu, \end{cases}$$

has a unique solution $P^\mu \in L^2_{\mathcal{F}}(0, T; \mathbb{R}^n)$ such that for any $t \in [0, T]$,

$$\begin{cases} P_t^\mu B_t a + \frac{1}{2} a R_t a + \frac{\sigma^2}{2} (U(a) + \ln \mu_t(a) + 1) + \beta = 0, \quad \text{for any } a \in \mathbb{R}^p, \\ \int \mu_t(a) da = 1, \end{cases}$$

where β is a random variable coming from Lagrange multiplier method. Moreover, $(Y^\mu, Z^\mu, P^\mu, \mu)$ composes a stochastic Hamiltonian system

$$(5.3) \quad \begin{cases} dY_t^\mu = (A_t Y_t^\mu + B_t \int a \mu_t(a) da + C_t Z_t^\mu) dt + Z_t^\mu dW_t, \\ dP_t^\mu = -(A_t P_t^\mu + H_t Y_t^\mu) dt - (C_t P_t^\mu + N_t Z_t^\mu) dW_t, \\ Y_T^\mu = \xi, \quad P_0^\mu = -G Y_0^\mu, \\ P_t^\mu B_t a + \frac{1}{2} a R_t a + \frac{\sigma^2}{2} (U(a) + \ln \mu_t(a) + 1) + \beta = 0, \quad \text{for any } a \in \mathbb{R}^p, \\ \int \mu_t(a) da = 1, \end{cases}$$

and Hamiltonian system (5.3) gives an optimal control

$$(5.4) \quad \begin{cases} \mu_t(a) = \frac{e^{-U(a) - \frac{\sigma^2}{2}(P_t^\mu B_t a + \frac{1}{2}a R_t a)}}{\int e^{-U(a) - \frac{\sigma^2}{2}(P_t^\mu B_t a + \frac{1}{2}a R_t a)} da}, \\ \beta = \frac{\sigma^2}{2}(\ln(\int e^{-U(a) - \frac{\sigma^2}{2}(P_t^\mu B_t a + \frac{1}{2}a R_t a)} da) - 1). \end{cases}$$

We then consider a specific case by setting the reference measure to be a standard normal distribution, i.e. $e^{-U(a)} = \frac{e^{-\frac{|a|^2}{2}}}{\sqrt{(2\pi)^p}}$. Then since $(R + \frac{\sigma^2}{2}\mathbb{I}_p) \in L^\infty(0, T; \hat{\mathbb{S}}_+^p)$, (5.4) implies

$$(5.5) \quad \begin{aligned} \mu_t(a) &= \frac{e^{-\frac{1}{\sigma^2}(a + (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu)(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)(a + (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu)}}{\int e^{-\frac{1}{\sigma^2}(a + (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu)(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)(a + (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu)} da} \\ &= \frac{1}{\sqrt{(2\pi)^p |\det(\Sigma_t^\mu)|}} e^{-\frac{1}{2}(a + (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu)(\Sigma_t^\mu)^{-1}(a + (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu)}, \end{aligned}$$

where $\Sigma_t^\mu = \frac{\sigma^2}{2}(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}$. It appears that μ_t has a Gaussian distribution and Σ_t^μ is the covariance matrix of μ_t . Hence, from (5.5) we know $v_t^\mu := \int a \mu_t(a) da = -(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu$ and $\mu_t \in \mathcal{P}_2(\mathbb{R}^p)$.

Remark 5.3. If we only take $U(\cdot) \equiv 0$ and assume that $R \in L^\infty(0, T; \hat{\mathbb{S}}_+^p)$, then v_t coincides with the strict control, and the optimal control in this case satisfies

$$\mu_t(a) = \frac{1}{\sqrt{(2\pi)^p |\det(\Sigma_t^\mu)|}} e^{-\frac{1}{2}(a + R_t^{-1}B_t P_t^\mu)(\Sigma_t^\mu)^{-1}(a + R_t^{-1}B_t P_t^\mu)},$$

where the covariance matrix $\Sigma_t^\mu = \frac{\sigma^2}{2}R_t^{-1}\mathbb{I}_p$ and $\text{tr}(\Sigma_t^\mu R_t) = \frac{\sigma^2}{2}p$. Also we can define the cost of exploration (COE) as in [28],

$$\begin{aligned} COE &:= \frac{1}{2}\mathbb{E}[\int_0^T \int a R_t a \mu_t(a) da - v_t^* R_t v_t^\mu dt] \\ &= \frac{1}{2}\mathbb{E}[\int_0^T \int (a - v_t^\mu) R_t (a - v_t^\mu) \mu_t(a) da dt] \\ &= \frac{1}{2}\mathbb{E}[\int_0^T \text{tr}(\Sigma_t^\mu R_t) dt] \\ &= \frac{1}{2}\mathbb{E}[\int_0^T \frac{\sigma^2}{2} p dt] \\ &= \frac{\sigma^2}{4} p T. \end{aligned}$$

As $\sigma \rightarrow 0$, the cost of relaxed control degenerates to the cost of strict control in the following sense:

$$\begin{cases} \mu_t \rightarrow \delta_{v_t^\mu} \text{ weakly as } \sigma \rightarrow \text{a.e. a.s.}, \\ \lim_{\sigma \rightarrow 0} COE = 0, \end{cases}$$

where $\delta_{v_t^\mu}$ stands for the Dirac measure defined at v_t^μ (see also Exercise 14.4.2 in Klenke [19]).

Similar to Theorem 4.2, we give the sufficient condition for an optimal control in LQ case.

THEOREM 5.4. *Under Assumption 5.1, the control $\mu \in \mathcal{A}$ is an optimal control of Problem (P2) if for any $t \in [0, T]$, it satisfies the Hamiltonian system (5.3) with (Y^μ, Z^μ) and P^μ be the solutions to the corresponding BSDE (5.1) and adjoint equation (5.2), respectively.*

Then in the case $e^{-U(a)} = \frac{e^{-\frac{|a|^2}{2}}}{\sqrt{(2\pi)^p}}$, we study the existence and uniqueness of optimal control in LQ case following the decoupling technique for backward stochastic system proposed in Lim and Zhou [20].

To begin with, assume that Y^μ has a decoupling form like

$$(5.6) \quad Y_t^\mu = \Theta_t P_t^\mu + \phi_t,$$

where Θ is a deterministic process with $\Theta_T = 0$ and differentiable on t , and ϕ satisfies BSDE

$$(5.7) \quad \begin{cases} d\phi_t = \lambda_t dt + \eta_t dW_t, \\ \phi_T = \xi \end{cases}$$

for some adapted processes λ and η which will be determined later.

Applying Itô formula to Y_t^μ , by (5.3) and (5.6) we have

$$\begin{aligned} 0 &= dY_t^\mu - \dot{\Theta}_t P_t^\mu dt - \Theta_t dP_t^\mu - d\phi_t \\ &= \left(A_t Y_t^\mu + B_t \int a \mu_t(a) da + C_t Z_t^\mu \right) dt + Z_t^\mu dW_t \\ &\quad - \dot{\Theta}_t P_t^\mu dt + \Theta_t (A_t P_t^\mu + H_t Y_t^\mu) dt + \Theta_t (C_t P_t^\mu + N_t Z_t^\mu) dW_t - \lambda_t dt - \eta_t dW_t. \end{aligned}$$

Bearing in mind that in this case $\int a \mu_t(a) da = -(R_t + \frac{\sigma^2}{2} \mathbb{I}_p)^{-1} B_t P_t^\mu$, we further have

$$\begin{cases} \lambda_t = A_t Y_t^\mu - B_t (R_t + \frac{\sigma^2}{2} \mathbb{I}_p)^{-1} B_t P_t^\mu + C_t Z_t - \dot{\Theta}_t P_t^\mu + \Theta_t (A_t P_t^\mu + H_t Y_t^\mu), \\ Z_t^\mu + \Theta_t (C_t P_t^\mu + N_t Z_t^\mu) - \eta_t = 0, \end{cases}$$

which implies $Z_t^\mu = (\mathbb{I}_p + \Theta_t N_t)^{-1} (\eta_t - \Theta_t C_t P_t^\mu)$. Since the coefficient ahead of P_t^μ is 0, we get the Riccati equation

$$(5.8) \quad \begin{cases} \dot{\Theta}_t - A_t \Theta_t - \Theta_t A_t - \Theta_t H_t \Theta_t + (R_t + \frac{\sigma^2}{2} \mathbb{I}_p)^{-1} B_t + C_t (\mathbb{I}_p + \Theta_t N_t)^{-1} \Theta_t C_t = 0, \\ \Theta_T = 0. \end{cases}$$

It is well known that the above Riccati equation (5.8) has a unique solution $\Theta \in L^\infty(0, T; \mathbb{S}_+^n)$ (see e.g. [20]). Hence (5.7) can be rewritten as below:

$$(5.9) \quad \begin{cases} d\phi_t = ((A_t + \Theta_t H_t) \phi_t + C_t (\mathbb{I}_p + \Theta_t N_t)^{-1} \eta_t) dt + \eta_t dW_t, \\ \phi_T = \xi. \end{cases}$$

The solvability of BSDE (5.9) comes from the classical results in [22].

THEOREM 5.5. *If the reference measure is a standard normal distribution, i.e. $e^{-U(a)} = \frac{e^{-\frac{|a|^2}{2}}}{\sqrt{(2\pi)^p}}$, under Assumption 5.1, stochastic Hamiltonian system (5.3) has a unique solution $(Y^\mu, Z^\mu, P^\mu, \mu)$, where for any $t \in [0, T]$,*

$$(5.10) \quad \begin{cases} Y_t^\mu = \Theta_t P_t^\mu + \phi_t, \\ Z_t^\mu = (\mathbb{I}_n + \Theta_t N_t)^{-1}(\eta_t - \Theta_t C_t P_t^\mu), \\ Y_0^\mu = (\mathbb{I}_n + \Theta_0 G)^{-1}\phi_0, \end{cases}$$

μ_t is Gaussian with the covariance matrix $\Sigma_t^\mu = \frac{\sigma^2}{2}(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}$ and the mean v_t^μ satisfying $(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)v_t^\mu + B_t P_t^\mu = 0$, Θ is the solution to Riccati equation (5.8) and (ϕ, η) is the solution to BSDE (5.9).

Proof. We first verify that (5.10) and the Gaussian random variable μ_t give a solution to stochastic Hamiltonian systems (5.3). Consider SDE

$$(5.11) \quad \begin{cases} dP_t^\mu = -(A_t P_t^\mu + H_t(\Theta_t P_t^\mu + \phi_t)) dt \\ \quad - (C_t P_t^\mu + N_t((\mathbb{I}_n + \Theta_t N_t)^{-1}(\eta_t - \Theta_t C_t P_t^\mu))) dW_t, \\ P_0^\mu = -G(\mathbb{I}_n + \Theta_0 G)^{-1}\phi_0. \end{cases}$$

Obviously, the linear SDE (5.11) has a unique solution P^μ . Applying Itô formula to $Y_t^\mu = \Theta_t P_t^\mu + \phi_t$, we have

$$dY_t^\mu = \left(A_t Y_t^\mu - B_t(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}B_t P_t^\mu + C_t(\mathbb{I}_n + \Theta_t N_t)^{-1}(\eta_t - \Theta_t C_t P_t^\mu) \right) dt \\ + (\mathbb{I}_n + \Theta_t N_t)^{-1}(\eta_t - \Theta_t C_t P_t^\mu)dW_t,$$

and $Y_0^\mu = (\mathbb{I}_n + \Theta_0 G)^{-1}\phi_0$. Noticing $Z_t^\mu = (\mathbb{I}_n + \Theta_t N_t)^{-1}(\eta_t - \Theta_t C_t P_t^\mu)$, we know that (Y^μ, Z^μ, P^μ) satisfies (5.3). As for μ , its explicit form (5.5) deduced from (5.3) and the argument below it demonstrate that μ_t satisfies a Gaussian distribution with the covariance matrix $\Sigma_t^\mu = \frac{\sigma^2}{2}(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}$ and the mean v_t^μ satisfying $(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)v_t^\mu + B_t P_t^\mu = 0$. So we prove that $(Y^\mu, Z^\mu, P^\mu, \mu)$ is a solution to stochastic Hamiltonian system (5.3).

As for the uniqueness of optimal control, assume that stochastic Hamiltonian systems (5.3) has two solutions (Y, Z, P, μ) and (Y', Z', P', μ') . Let $\bar{\varphi} = \varphi - \varphi'$, $\varphi = Y, Z, P, \mu$ and $\bar{v} = \int a\bar{\mu}(a)da$. Then $(\bar{Y}, \bar{Z}, \bar{P}, \bar{v})$ satisfies

$$\begin{cases} d\bar{Y}_t = (A_t \bar{Y}_t + B_t \bar{v}_t + C_t \bar{Z}_t) dt + \bar{Z}_t dW_t, \\ d\bar{P} = -(A_t \bar{P}_t + H_t \bar{Y}_t) dt - (C_t \bar{P}_t + N_t \bar{Z}_t) dW_t, \\ \bar{Y}_T = 0, \quad \bar{P}_0 = -G\bar{Y}_0, \\ (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)\bar{v}_t + B_t \bar{P}_t = 0. \end{cases}$$

Applying Itô formula to $\bar{Y}_t \bar{P}_t$, we obtain

$$\mathbb{E}[\bar{Y}_0 G \bar{Y}_0] = -\mathbb{E}\left[\int_0^T \left(\bar{Y}_t H_t \bar{Y}_t + \bar{P}_t B_t (R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1} B_t \bar{P}_t + \bar{Z}_t N_t \bar{Z}_t \right) dt\right].$$

From Assumption 5.1, we know $G, H, N \in \mathbb{S}_+^n$ and $(R_t + \frac{\sigma^2}{2}\mathbb{I}_p) \in \hat{\mathbb{S}}_+^n$. Hence $B_t \bar{P}_t = 0$ a.s. Consequently, (\bar{Y}, \bar{Z}) satisfies

$$(5.12) \quad \begin{cases} d\bar{Y}_t = (A_t \bar{Y}_t + C_t \bar{Z}_t) dt + \bar{Z}_t dW_t, \\ \bar{Y}_T = 0. \end{cases}$$

Obviously, (5.12) has a unique solution $(\bar{Y}, \bar{Z}) \equiv 0$. So

$$\begin{cases} d\bar{P} = -A_t \bar{P}_t dt - C_t \bar{P}_t dW, \\ \bar{P}_0 = -G\bar{Y}_0, \end{cases}$$

also suggests $\bar{P} \equiv 0$, and then $\bar{\mu} \equiv 0$ follows immediately from the means and covariances of the optimal controls are identical. \square

Moreover, for a suitable reference measure, the solution to stochastic Hamiltonian system (5.3) is the optimal control of backward stochastic LQ control system with entropy regularization.

COROLLARY 5.6. *If the reference measure is a standard normal distribution, i.e. $e^{-U(a)} = \frac{e^{-\frac{|a|^2}{2}}}{\sqrt{(2\pi)^p}}$, under Assumption 5.1, the solution to stochastic Hamiltonian system (5.3) is the unique optimal control of Problem (P2).*

Proof. We only need to prove that the solution μ to stochastic Hamiltonian system (5.3) is an admissible control of Problem (P2).

Recall that μ_t is Gaussian with the covariance matrix $\Sigma_t^\mu = \frac{\sigma^2}{2}(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)^{-1}$ and the mean v_t^μ satisfying $(R_t + \frac{\sigma^2}{2}\mathbb{I}_p)v_t^\mu + B_t P_t^\mu = 0$. Noticing $(R + \frac{\sigma^2}{2}\mathbb{I}_p) \in L^\infty(0, T; \hat{\mathbb{S}}_+^p)$, $B \in L^\infty(0, T; \mathbb{R}^{n \times p})$ and $P^\mu \in S_{\mathcal{F}}^2(0, T; \mathbb{R}^n)$, we first have

$$\mathbb{E} \left[\int_0^T \int |a|^2 \mu_t(a) da dt \right] = \mathbb{E} \left[\int_0^T (|v_t^\mu|^2 + \text{tr}(\Sigma_t^\mu)) dt \right] \leq C \left(1 + \mathbb{E} \left[\int_0^T |P_t^\mu|^2 dt \right] \right) < \infty.$$

On the other hand,

$$\begin{aligned} \mathbb{E} \left[\int_0^T \text{Ent}(\mu_t | e^{-U}) da dt \right] &= \mathbb{E} \left[\int_0^T \frac{1}{2} (-\ln(\det(\Sigma_t^\mu)) + \text{tr}(\Sigma_t^\mu) + |v_t^\mu|^2 - p) dt \right] \\ &\leq C \left(1 + \mathbb{E} \left[\int_0^T |P_t^\mu|^2 dt \right] \right) < \infty. \end{aligned}$$

Therefore, $\mu \in \mathcal{A}$ follows. \square

REFERENCES

- [1] H. BECKER AND V. MANDREKAR, *On the existence of optimal random controls*, J. Math. Mech., 18 (1968/69), pp. 1151–1166.
- [2] R. BUCKDAHN, J. LI, S. PENG, AND C. RAINER, *Mean-field stochastic differential equations and associated pdes*, The Annals of Probability, 45 (2017), pp. 824–878.
- [3] R. CARMONA AND F. DELARUE, *Forward-backward stochastic differential equations and controlled McKean-Vlasov dynamics*, The Annals of Probability, 43 (2015), pp. 2647–2700.
- [4] ———, *Probabilistic Theory of Mean Field Games with Applications I*, Springer, 2018.
- [5] M. DAI, Y. DONG, Y. JIA, AND X. Y. ZHOU, *Learning merton strategies in an incomplete market: recursive entropy regularization and biased gaussian exploration*, arXiv:2312.11797, 43pp.
- [6] N. DOKUCHAEV AND X. Y. ZHOU, *Stochastic controls with terminal contingent conditions*, Journal of Mathematical Analysis and Applications, 238 (1999), pp. 143–165.
- [7] D. DUFFIE AND L. G. EPSTEIN, *Stochastic differential utility*, Econometrica, 60 (1992), pp. 353–394.
- [8] N. EL KAROUI, D. HÙU NGUYEN, AND M. JEANBLANC-PICQUÉ, *Compactification methods in the control of degenerate diffusions: Existence of an optimal control*, Stochastics, 20 (1987), pp. 169–219.
- [9] N. EL KAROUI, S. PENG, AND M. C. QUENEZ, *Backward stochastic differential equations in finance*, Mathematical Finance, 7 (1997), pp. 1–71.

- [10] D. FIROOZI AND S. JAIMUNGAL, *Exploratory LQG mean field games with entropy regularization*, Automatica, 139 (2022), Paper No. 110177, 12pp.
- [11] W. H. FLEMING, *Generalized solutions in optimal stochastic control*, in Differential games and control theory, II (Proc. 2nd Conf., Univ. Rhode Island, Kingston, R.I., 1976), vol. 30 of Lect. Notes Pure Appl. Math., Dekker, New York, 1977, pp. 147–165.
- [12] X. GAO, Z. Q. XU, AND X. Y. ZHOU, *State-Dependent temperature control for Langevin diffusions*, SIAM Journal on Control and Optimization, 60 (2022), pp. 1250–1268.
- [13] D. A. GOMES AND E. VALDINOCI, *Entropy penalization methods for Hamilton Jacobi equations*, Advances in Mathematics, 215 (2007), pp. 94–152.
- [14] K. HU, Z. REN, D. ŠIŠKA, AND L. SZPRUCH, *Mean-field Langevin dynamics and energy landscape of neural networks*, Annales de l’Institut Henri Poincaré, Probabilités et Statistiques, 57 (2021), pp. 2043 – 2065.
- [15] Y. JIA AND X. Y. ZHOU, *q-learning in continuous time*, Journal of Machine Learning Research, 24 (2023), pp. 1–61.
- [16] R. JORDAN, D. KINDERLEHRER, AND F. OTTO, *The variational formulation of the Fokker–Planck Equation*, SIAM Journal on Mathematical Analysis, 29 (1998), pp. 1–17.
- [17] C. KARNAM, J. MA, AND J. ZHANG, *Dynamic approaches for some time-inconsistent optimization problems*, The Annals of Applied Probability, 27 (2017), pp. 3435–3477.
- [18] B. KERIMKULOV, D. ŠIŠKA, L. SZPRUCH, AND Y. ZHANG, *Mirror descent for stochastic control problems with measure-valued controls*, arXiv:2401.01198, 37pp.
- [19] A. KLENKE, *Probability theory*, Springer, 2020.
- [20] A. LIM AND X. Y. ZHOU, *Linear-Quadratic control of backward stochastic differential equations*, SIAM Journal on Control and Optimization, 40 (2001), pp. 450–474.
- [21] P. LIONS, *Cours au coll ge de france: Th orie des jeu champs moyens.*, Available at [http://www.college-de-france.fr/default/EN/all/equ\[1\]der/audiovideo.jsp](http://www.college-de-france.fr/default/EN/all/equ[1]der/audiovideo.jsp), (2013).
- [22] E. PARDOUX AND S. PENG, *Adapted solution of a backward stochastic differential equation*, System Control Letters, 14 (1990), pp. 55–61.
- [23] S. PENG, *Probabilistic interpretation for systems of quasilinear parabolic partial differential equations*, Stochastics and Stochastics Reports, 37 (1991), pp. 61–74.
- [24] ———, *Backward stochastic differential equations and applications to optimal control*, Applied Mathematics and Optimization, 27 (1993), pp. 125–144.
- [25] C. REISINGER AND Y. ZHANG, *Regularity and stability of feedback relaxed controls*, SIAM Journal on Control and Optimization, 59 (2021), pp. 3118–3151.
- [26] D. ŠIŠKA AND L. SZPRUCH, *Gradient flows for regularized stochastic control problems*, SIAM Journal on Control and Optimization, 62 (2024), pp. 2036–2070.
- [27] A. TAKAHASHI, Y. TSUCHIDA, AND T. YAMADA, *A new efficient approximation scheme for solving high-dimensional semilinear pdes: control variate method for deep bsde solver*, Journal of Computational Physics, 454 (2022), Paper No. 110956, 39pp.
- [28] H. WANG, T. ZARIPHOPOULOU, AND X. Y. ZHOU, *Reinforcement learning in continuous time and space: A stochastic control approach*, Journal of Machine Learning Research, 21 (2020), pp. 1–34.
- [29] H. WANG AND X. Y. ZHOU, *Continuous-time mean-variance portfolio selection: A reinforcement learning framework*, Mathematical Finance, 30 (2020), pp. 1273–1308.
- [30] Q. ZHANG AND H. ZHAO, *Stationary solutions of SPDEs and infinite horizon BDSDEs*, Journal of Functional Analysis, 252 (2007), pp. 171–219.