

Entropy-Regularized Stochastic Games

Yagiz Savas, Mohamadreza Ahmadi, Takashi Tanaka, Ufuk Topcu

Abstract—In two-player zero-sum stochastic games, where two competing players make decisions under uncertainty, a pair of optimal strategies is traditionally described by Nash equilibrium and computed under the assumption that the players have perfect information about the stochastic transition model of the environment. However, implementing such strategies may make the players vulnerable to unforeseen changes in the environment. In this paper, we introduce entropy-regularized stochastic games where each player aims to maximize the causal entropy of its strategy in addition to its expected payoff. The regularization term balances each player’s rationality with its belief about the level of misinformation about the transition model. We consider both entropy-regularized N -stage and entropy-regularized discounted stochastic games, and establish the existence of a value in both games. Moreover, we prove the sufficiency of Markovian and stationary mixed strategies to attain the value, respectively, in N -stage and discounted games. Finally, we present algorithms, which are based on convex optimization problems, to compute the optimal strategies. In a numerical example, we demonstrate the proposed method on a motion planning scenario and illustrate the effect of the regularization term on the expected payoff.

I. INTRODUCTION

A two-player zero-sum stochastic game (SG) [1] models sequential decision-making of two players with opposing objectives in a stochastic environment. An SG is played in stages. At each stage, the game is in a state, and the players choose one of their available actions simultaneously and receive payoffs. The game then transitions to a new random state according to a probability distribution which represents the stochasticity in the environment.

In an SG, each player aims to synthesize a strategy that maximizes the player’s expected payoff at the end of the game. Traditionally, a pair of optimal strategies is described by Nash equilibrium [2] according to which both players play their best-response strategies against the opponent’s strategy. The value of the game then corresponds to the expected payoff that each player receives at the end of the game, if they both play their respective equilibrium strategies.

The concept of Nash equilibrium is based on the assumptions that the players have perfect information about the environment and act rationally [3]. However, in certain scenarios, the information that a player has about the environment may not match the reality. For example, in a planning scenario, if the player obtains its information about the environment through surveillance missions performed in the

past, it may face with a significantly different environment during the execution of the play. In such scenarios, playing an equilibrium strategy may dramatically decrease the player’s actual expected payoff as the strategy is computed under the assumption of perfect information.

The principle of maximum entropy prescribes a probability distribution that is “maximally noncommittal with regard to missing information” [4]. The principle of maximum causal entropy extends the maximum entropy principle to settings where there is dynamically revealed side information that causally affects the evolution of a stochastic process [5], [6]. A distribution that maximizes the causal entropy of a stochastic process (in the absence of additional constraints) is the one that makes all admissible realizations equally probable regardless of the revealed information [7]. Therefore, the causal entropy of a player’s strategy provides a convenient way to quantify the dependence of its strategy to its level of information about the environment as well as the other player’s strategy.

In this paper, we propose a method to synthesize a pair of strategies that balances each player’s rationality with its belief about the level of missing information. Specifically, we regularize each player’s objective with the causal entropy of its strategy which is causally dependent on the history of play. Therefore, the proposed method allows the players to adjust their strategies according to different levels of misinformation by tuning a parameter that controls the importance of the regularization term. For example, in two extremes, it allows the player to be perfectly rational or to purely randomize its strategy.

We study both entropy-regularized N -stage and entropy-regularized discounted games, and show the existence of a value in both games. We first prove the sufficiency of Markovian and stationary strategies for both players, respectively, in N -stage and discounted games in order to maximize their entropy-regularized expected payoff. Then, we provide algorithms based on a sequence of convex optimization problems to compute a pair of equilibrium strategies. Finally, we demonstrate the proposed methods on a motion planning scenario, and illustrate that the introduced regularization term yields strategies that perform well in different environments.

Related work. In stochastic games literature, the idea of balancing the expected payoffs with an additional regularization term appeared recently in [8] and [9]. The work [8] proposes to bound the rationality of the players to obtain tunable behavior in video games. They study γ -discounted games and restrict their attention to stationary strategies to balance the expected payoffs with the Kullback-Leibner distance of the player’s strategies from reference strategies. In [9], authors

Y. Savas, T. Tanaka and U. Topcu are with the Department of Aerospace Engineering, University of Texas at Austin, TX, USA. E-mail: {yagiz.savas, ttanaka, utopcu}@utexas.edu

M. Ahmadi is with the Center for Autonomous Systems and Technologies, California Institute of Technology, CA, USA. E-mail: mrahmadi@caltech.edu

study N -stage games and consider only Markovian strategies to balance the expected payoffs with the player's sensing costs which are expressed as directed information from states to actions. Unlike this work, we introduce causal entropy of strategies as the regularization term. Additionally, we allow the player's to follow history-dependent strategies and prove the sufficiency of Markovian strategies to attain the value in N -stage games.

Regularization terms are also used in matrix and extensive form games generally to learn equilibrium strategies [10], [11], [12], [13]. When each player uses the same parameter to regularize its expected payoffs with the entropy of its strategy, an equilibrium strategy profile is called a quantal response equilibrium (QRE) [3], and an equilibrium strategy of a player is referred as quantal best response [11] or logit choice strategy [10]. From a theoretical perspective, the main difference between our approach and the well-studied QRE concept [14] is that we establish the existence of equilibrium strategies even if the players use different regularization parameters. Additionally, we provide an efficient algorithm based on a convex optimization problem to compute the equilibrium strategies.

Robust stochastic games [15], [16] concern the synthesis of equilibrium strategies when the uncertainty in transition probabilities and payoff functions can be represented by structured sets. Unlike robust SG models, the proposed method in this paper can still be used when it is not possible to form a structured uncertainty set.

In reinforcement learning literature, the use of regularization terms is extensively studied to obtain robust behaviors [17], improve the convergence rates [18], and compute optimal strategies efficiently [19]. As stochastic games model multi-player interactions, our approach leverages the ideas discussed in aforementioned work to environments where an adversary aims to prevent a player to achieve its objective.

II. BACKGROUND

We first review some concepts from game theory and information theory that will be used in the subsequent sections.

Notation: For a sequence x , we write x^t to denote (x_1, x_2, \dots, x_t) . Upper case symbols such as X denote random variables, and lower case symbols such as x denote a specific realization. The cardinality of a set \mathcal{X} is denoted by $|\mathcal{X}|$, and the probability simplex defined over the set \mathcal{X} is denoted by $\Delta(\mathcal{X})$. For $V_1, V_2 \in \mathbb{R}^n$, we write $V_1 \preceq V_2$ to denote the coordinate-wise inequalities. We use the index set $\mathbb{Z}_+ = \{1, 2, \dots\}$ and the natural logarithm $\log(\cdot) = \log_e(\cdot)$.

A. Two-Player Stochastic Games

A two-player stochastic game Γ [1] is played in stages. At each stage t , the game is in one of its finitely many states \mathcal{X} , and each player observes the current state x_t . At each state x_t , the players choose one of their finitely many actions, and the game transitions to a successor state x_{t+1} according to a probability distribution $\mathcal{P}: \mathcal{X} \times \mathcal{U} \times \mathcal{W} \rightarrow \Delta(\mathcal{X})$ where \mathcal{U} and \mathcal{W} are finite action spaces for player 1 and

player 2, respectively. The pair of actions, $u_t \in \mathcal{U}$ and $w_t \in \mathcal{W}$, together with the current state $x_t \in \mathcal{X}$ determine the payoff $\mathcal{R}(x_t, u_t, w_t) \leq \bar{\mathcal{R}} < \infty$ to be made by player 2 to player 1 at stage t .

A player's strategy is a specification of a probability distribution over available actions at each stage conditional on the history of the game up to that stage. Formally, let $\mathcal{H}_t = (\mathcal{X} \times \mathcal{U} \times \mathcal{W})^{t-1} \times \mathcal{X}$ be the set of all possible history of plays up to stage t . Then, the strategy of player 1 and player 2 are denoted by $\sigma = (\sigma_1, \sigma_2, \dots)$ and $\tau = (\tau_1, \tau_2, \dots)$, respectively, where $\sigma_t: \mathcal{H}_t \rightarrow \Delta(\mathcal{U})$ and $\tau_t: \mathcal{H}_t \rightarrow \Delta(\mathcal{W})$ for all t . If a player's strategy depends only on the current state for all stages, e.g., $\sigma_t: \mathcal{X} \rightarrow \Delta(\mathcal{U})$ for all t , the strategy is said to be *Markovian*. A *stationary* strategy depends only on the current state and is independent of the stage number, e.g., $\sigma = (\sigma, \sigma, \dots)$, where $\sigma: \mathcal{X} \rightarrow \Delta(\mathcal{U})$. We denote the set of all strategies, all Markovian strategies, and all stationary strategies for player $i \in \{1, 2\}$ by Γ_i , Γ_i^M , and Γ_i^S , respectively.

Let $\mu_{t+1}(x^{t+1}, u^t, w^t)$ be the joint probability distribution over the history \mathcal{H}_{t+1} of play which is uniquely determined by the initial state distribution $\mu_1(h_1)$ through the recursive formula

$$\begin{aligned} \mu_{t+1}(x^{t+1}, u^t, w^t) &= \mathcal{P}(x_{t+1}|x_t, u_t, w_t) \sigma_t(u_t|h_t) \\ &\quad \times \tau_t(w_t|h_t) \mu_t(h_t) \end{aligned} \quad (1)$$

where $h_t \in \mathcal{H}_t$ is the history of play up to stage t .

A stochastic game with the initial distribution $\mu_1(x_1)$ is called an N -stage game, if the game ends after N stages. The evaluation function for an N -stage game is

$$\begin{aligned} J(X^N, U^N, W^N) &:= \\ &\sum_{t=1}^N \mathbb{E}^{\bar{\mu}_t} \mathcal{R}(X_t, U_t, W_t) + \mathbb{E}^{\mu^{T+1}} \mathcal{R}(X_{N+1}) \end{aligned} \quad (2)$$

where $\bar{\mu}_t(\cdot) := \mu_t(\cdot) \sigma_t(\cdot) \tau_t(\cdot)$. Similarly, if the number of stages in the game is infinite, and the future payoffs are discounted by a factor $0 < \gamma < 1$, the game is called a γ -discounted game. The evaluation function for a γ -discounted game is

$$\sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}^{\bar{\mu}_t} \mathcal{R}(X_t, U_t, W_t). \quad (3)$$

The player 1's objective is to maximize the evaluation function, i.e., its expected payoff, whereas the player 2 aims to minimize it. A stochastic game is said to have the *value* \mathcal{V}^* , if for an evaluation function $f(\sigma, \tau)$, we have

$$\mathcal{V}^* = \max_{\sigma \in \Gamma_u} \min_{\tau \in \Gamma_w} f(\sigma, \tau) = \min_{\tau \in \Gamma_w} \max_{\sigma \in \Gamma_u} f(\sigma, \tau).$$

A pair of strategies (σ^*, τ^*) is said to be equilibrium strategies if it attains the value of the game.

It is well-known that both N -stage and γ -discounted games have a value for finite state and action sets [20]. Moreover, Markovian and stationary strategies are sufficient for players to attain the value in N -stage and γ -discounted games, respectively [1], [21].

B. Causal Entropy

For a sequential decision-making problem where decisions depend causally on the past information such as the history of play, the causal entropy of a strategy is a measure to quantify the randomness of the strategy. Let X^N , Y^N and Z^N be sequences of random variables with length N . The entropy of the sequence X^N causally conditioned on the sequences Y^N and Z^N is defined as [22]

$$H(X^N || Y^N, Z^N) := \sum_{t=1}^N H(X_t | X^{t-1}, Y^t, Z^t), \quad (4)$$

where

$$H(X_t | X^{t-1}, Y^t, Z^t) := - \sum_{x^N, y^N, z^N} \Pr(x^t, y^t, z^t) \log \Pr(x_t | x^{t-1}, y^t, z^t). \quad (5)$$

The concept of causal entropy has recently been used to infer correlated-equilibrium strategies in Markov games [23] and to recover cost functions in inverse optimal control problems [6]. In this study, we employ causal entropy to compute an equilibrium strategy profile that balances the players' expected payoff with the randomness of their strategies in stochastic games.

In the absence of additional constraints, a strategy $\sigma \in \Gamma_1$ that maximizes the causal entropy $H(U^N || X^N, W^{N-1})$ of the player 1, which is conditioned on the revealed history of play, is the stationary strategy $\sigma = (\sigma, \sigma, \dots)$ where $\sigma(x)(u) = 1/|\mathcal{U}|$. Therefore, a player that maximizes the entropy of its strategy acts purely randomly regardless of the history of play. On the other hand, a player that regularizes its expected payoff with the entropy of its strategy can be thought as a player that balances its rationality with its belief about the correctness of the underlying transition model of the environment.

III. PROBLEM STATEMENT

We first consider entropy-regularized N -stage games for which we define the evaluation function as

$$\Phi_N(\sigma, \tau) := J(X^N, U^N, W^N) + \frac{1}{\beta_1} H(U^N || X^N, W^{N-1}) - \frac{1}{\beta_2} H(W^N || X^N, U^{N-1}), \quad (6)$$

where $\beta_1, \beta_2 > 0$ are regularization parameters that adjust for players the importance of the randomness in their strategies. Note that, when $\beta_1 = \beta_2 = \infty$, both players act perfectly rational, and we recover the evaluation function (2). Additionally, since the play is simultaneous, the information of a player's strategy at a given stage is not revealed to the other player. Hence, at each stage, players are allowed to condition their strategies only to observed history of play.

Problem 1: Provide an algorithm to synthesize, if exists, equilibrium strategies in entropy-regularized N -stage games.

We next consider stochastic games that are played in infinite stages, and introduce entropy-regularized γ -discounted

games for which we define the evaluation function as

$$\Phi_\infty(\sigma, \tau) := \sum_{t=1}^{\infty} \gamma^{t-1} \left[\mathbb{E}^{\bar{\sigma}, \bar{\tau}} \mathcal{R}(X_t, U_t, W_t) + \frac{1}{\beta_1} H(U_t | H_t) - \frac{1}{\beta_2} H(W_t | H_t) \right], \quad (7)$$

where $H_t = (X^t, U^{t-1}, W^{t-1})$, i.e., the admissible histories of play at stage t . Note that in the evaluation function (7), we discount players' future entropy gains as well as the expected payoff in order to ensure the finiteness of the evaluation function.

Problem 2: Provide an algorithm to synthesize, if exists, equilibrium strategies in entropy-regularized γ -discounted games.

IV. EXISTENCE OF VALUES AND THE COMPUTATION OF OPTIMAL STRATEGIES

In this section, we analyze entropy regularized N -stage and γ -discounted games, and show that both games have values. Then, we provide algorithms to synthesize equilibrium strategies that attain the corresponding game values.

A. Entropy-Regularized N -Stage Games

Searching optimal strategies that solve a stochastic game with the evaluation function $\Phi_N(\sigma, \tau)$ in the space of all strategies can be intractable for large N . We begin with establishing the existence of optimal strategies for both players in the space of Markovian strategies.

Proposition 1: Markovian strategies are sufficient for both players to attain, if exists, the value in entropy-regularized N -stage games, i.e.,

$$\begin{aligned} \max_{\sigma \in \Gamma_u^M} \min_{\tau \in \Gamma_w^M} \Phi_N(\sigma, \tau) &= \max_{\sigma \in \Gamma_u} \min_{\tau \in \Gamma_w} \Phi_N(\sigma, \tau), \\ \min_{\tau \in \Gamma_w^M} \max_{\sigma \in \Gamma_u^M} \Phi_N(\sigma, \tau) &= \min_{\tau \in \Gamma_w} \max_{\sigma \in \Gamma_u} \Phi_N(\sigma, \tau). \end{aligned}$$

Proof: See Appendix A. \square

Next, we show that entropy-regularized N -stage games have a value. Let $\rho_t: \mathcal{X} \times \mathcal{U} \times \mathcal{W} \rightarrow \mathbb{R}$ be a function and x_t be a fixed state. Additionally, let

$$\mathcal{V}_t^{\sigma_t, \tau_t}(x_t) := \mathbb{E}^{\sigma_t, \tau_t} \left[\rho_t(x_t, u_t, w_t) - \frac{1}{\beta_1} \log \sigma_t(u_t | x_t) + \frac{1}{\beta_2} \log \tau_t(w_t | x_t) \right] \quad (8)$$

be the evaluation function for a "one-shot" game in which the game starts from the state x_t and ends after both players play their one-step strategy.

Proposition 2: A stochastic game with the evaluation function (8) has a value, i.e.,

$$\max_{\sigma_t \in \Delta(\mathcal{U})} \min_{\tau_t \in \Delta(\mathcal{W})} \mathcal{V}_t^{\sigma_t, \tau_t}(x_t) = \min_{\tau_t \in \Delta(\mathcal{W})} \max_{\sigma_t \in \Delta(\mathcal{U})} \mathcal{V}_t^{\sigma_t, \tau_t}(x_t).$$
 Proof: It is clear that $\mathcal{V}_t(x_t)$ is a continuous function that is concave in σ_t and convex in τ_t . Additionally, $\Delta(\mathcal{U})$ and $\Delta(\mathcal{W})$ are compact convex sets. The result follows from von Neumann's minimax theorem [24]. \square

The following proposition states that one can compute the value of the one shot game (8) and synthesize equilibrium strategies by solving a convex optimization problem.

Proposition 3: For a given one-shot game with the evaluation function (8), optimal strategies (σ_t^*, τ_t^*) satisfy

$$\begin{aligned} \sigma_t^*(u_t|x_t) \in \\ \arg \max_{\sigma_t \in \Delta(U)} \left[-\frac{1}{\beta_1} \sum_{u_t \in U} \sigma_t(u_t|x_t) \log \sigma_t(u_t|x_t) \right. \\ \left. - \frac{1}{\beta_2} \log \sum_{w_t \in \mathcal{W}} \exp\left(-\beta_2 \sum_{u_t \in U} \sigma_t(u_t|x_t) \rho_t(x_t, u_t, w_t)\right) \right], \end{aligned} \quad (9)$$

$$\tau_t^*(w_t|x_t) = \frac{\exp\left(-\beta_2 \sum_{u_t \in U} \sigma_t^*(u_t|x_t) \rho_t(x_t, u_t, w_t)\right)}{\sum_{w_t \in \mathcal{W}} \exp\left(-\beta_2 \sum_{u_t \in U} \sigma_t^*(u_t|x_t) \rho_t(x_t, u_t, w_t)\right)}. \quad (10)$$

Furthermore, the unique value $\mathcal{V}_t^*(x_t)$ of the game is given by

$$\begin{aligned} \mathcal{V}_t^*(x_t) = -\frac{1}{\beta_1} \sum_{u_t \in U} \sigma_t^*(u_t|x_t) \log \sigma_t^*(u_t|x_t) \\ - \frac{1}{\beta_2} \log \sum_{w_t \in \mathcal{W}} \exp\left(-\beta_2 \sum_{u_t \in U} \sigma_t^*(u_t|x_t) \rho_t(x_t, u_t, w_t)\right). \end{aligned} \quad (11)$$

Proof: See Appendix A. \square

It is worth noting that the objective function of the optimization problem given in (9) is strictly concave, and therefore, optimal strategies $\sigma_t^*(u_t|x_t)$ and $\tau_t^*(w_t|x_t)$ are unique. Additionally, an optimal strategy with the form (10) is known in the economics literature as quantal best response [14], and for $\beta_1 = \beta_2 < \infty$, the optimal strategies form the well-studied quantal response equilibrium strategies [3].

We remark that the optimization problem in (9) has a closed-form solution which is a function of the optimal strategy $\tau_t^*(w_t|x_t)$. However, since the closed-form expressions for equilibrium strategies constitute a system of coupled nonlinear equations, the convex optimization formulation provides a more convenient way to compute equilibrium strategies.

Utilizing the results of above propositions, we now reformulate an entropy-regularized N -stage game as a series of ‘‘one-shot’’ games through the use of Bellman recursions. Let

$$\begin{aligned} \rho_t(x_t, u_t, w_t) = \mathcal{R}(x_t, u_t, w_t) \\ + \sum_{x_{t+1} \in \mathcal{X}} \mathcal{P}(x_{t+1}|x_t, u_t, w_t) \mathcal{V}_{t+1}^{\sigma_t, \tau_t}(x_{t+1}), \end{aligned} \quad (12)$$

for $t=1, \dots, N$ where $\mathcal{V}_{N+1}^{\sigma_t, \tau_t}(x_{N+1}) = \mathcal{R}(x_{N+1})$. Then, it can be easily verified that

$$\Phi_N(\sigma, \tau) = \sum_{x_1 \in \mathcal{X}} \mu_1(x_1) \mathcal{V}_1(x_1) \quad (13)$$

for a given initial distribution $\mu_1(x_1)$. Consequently, we obtain the following result.

Theorem 1: Entropy-regularized N -stage games have a value.

Proof: Due to Proposition 1, we can focus on Markovian strategies to find an equilibrium point in N -stage games. We start from the stage $k=N$ and compute the value of the one shot game (8), which exists due to Proposition 2. Using (8)

with (12) for $k=N-1, N-2, \dots, 1$, we compute the value of $N-k+1$ stage games. As a result, the claim follows due to the equivalence given in (13). \square

Algorithm 1 summarizes the computation of the pair (σ^*, τ^*) of optimal strategies for entropy-regularized N -stage games.

Algorithm 1 Strategy computation for N -stage games

- 1: **Initialize:** $\mathcal{V}_{N+1}(x_{N+1}) = \mathcal{R}(x_{N+1})$ for all $x_{N+1} \in \mathcal{X}$.
 - 2: **for** $t = N, N-1, \dots, 1$ **do**
 - 3: Compute $\rho_t(x_t, u_t, w_t)$ for all $x_t \in \mathcal{X}$, $u_t \in U$, and $w_t \in \mathcal{W}$ as in (12).
 - 4: For all $x_t \in \mathcal{X}$,

$$\begin{aligned} \mathcal{V}_t^*(x_t) = \max_{\sigma_t \in \Delta(U)} & -\frac{1}{\beta_1} \sum_{u_t \in U} \sigma_t(u_t|x_t) \log \sigma_t(u_t|x_t) \\ & - \frac{1}{\beta_2} \log \sum_{w_t \in \mathcal{W}} \exp\left(-\beta_2 \sum_{u_t \in U} \sigma_t(u_t|x_t) \rho_t(x_t, u_t, w_t)\right) \end{aligned}$$
 - 5: For all $x_t \in \mathcal{X}$, compute σ_t^* and τ_t^* as in (9) and (10), respectively.
 - 6: **return** $\sigma^* = (\sigma_1^*, \dots, \sigma_T^*)$ and $\tau^* = (\tau_1^*, \dots, \tau_T^*)$.
-

Remark: In certain scenarios, one of the players may prefer to play perfectly rationally against a boundedly rational opponent, e.g., $\beta_2 = \infty$. In that case, it is still possible to compute equilibrium strategies by solving a convex optimization problem at each stage. The value of the one-shot game (8) still exists due to the arguments provided in the proof of Proposition 2. However, the form of optimal strategies slightly changes to

$$\begin{aligned} \tau_t^*(w_t|x_t) = \arg \min_{\tau_t \in \Delta(W)} \frac{1}{\beta_1} \log \sum_{u_t \in U} \exp\left(\right. \\ \left. \beta_1 \sum_{w_t \in \mathcal{W}} \tau_t(w_t|x_t) \rho_t(x_t, u_t, w_t) \right), \end{aligned} \quad (14)$$

$$\begin{aligned} \sigma_t^*(u_t|x_t) = \\ \frac{\exp\left(\beta_1 \sum_{w_t \in \mathcal{W}} \tau_t^*(w_t|x_t) \rho_t(x_t, u_t, w_t)\right)}{\sum_{u_t \in U} \exp\left(\beta_1 \sum_{w_t \in \mathcal{W}} \tau_t^*(w_t|x_t) \rho_t(x_t, u_t, w_t)\right)}. \end{aligned} \quad (15)$$

It is important to note that, if $\beta_i = \infty$ for some $i \in \{1, 2\}$, equilibrium strategies may be not unique since the function $\log \sum \exp(\cdot)$ is not strictly convex over its domain [25].

B. Entropy-Regularized γ -Discounted Games

In this section, we focus on Markovian strategies, whose optimality for N -stage games is shown in Proposition 1. Let $\mathcal{V} \in \mathbb{R}^{|\mathcal{X}|}$ be a real-valued function, and for a given $x \in \mathcal{X}$, $\mathcal{L}(\mathcal{V})(x, \cdot, \cdot) : U \times \mathcal{W} \rightarrow \mathbb{R}$ be a function such that

$$\begin{aligned} \mathcal{L}(\mathcal{V})(x, \sigma, \tau) := \mathbb{E}^{\sigma, \tau} \left[\mathcal{R}(x, u, w) - \frac{1}{\beta_1} \log \sigma(u|x) \right. \\ \left. + \frac{1}{\beta_2} \log \tau(w|x) + \gamma \sum_{x' \in \mathcal{X}} \mathcal{P}(x'|x, u, w) \mathcal{V}(x') \right]. \end{aligned} \quad (16)$$

As discussed in Proposition 2, a one-shot game with the evaluation function $\mathcal{L}(\mathcal{V})(x, \sigma, \tau)$ has a value. Therefore, we

can introduce the Shapley operator $\Psi: \mathcal{V} \rightarrow \Psi(\mathcal{V})$ from $\mathbb{R}^{|\mathcal{X}|}$ to itself specified, for all $x \in \mathcal{X}$, as

$$\Psi(\mathcal{V})[x] := \max_{\sigma \in \Delta(U)} \min_{\tau \in \Delta(W)} \mathcal{L}(\mathcal{V})(x, \sigma, \tau).$$

It is clear that the operator Ψ satisfies two key properties: *monotonicity*, i.e., $\mathcal{V} \preceq \bar{\mathcal{V}}$ implies $\Psi(\mathcal{V}) \preceq \Psi(\bar{\mathcal{V}})$, and *reduction of constants*, i.e., for any $k \geq 0$, $\Psi(\mathcal{V} + k\mathbf{1})[x] = \Psi(\mathcal{V})[x] + \gamma k$ for all $x \in \mathcal{X}$. Consequently, it is straightforward to show that the operator Ψ is a contraction mapping [26]. Specifically, we have

$$\|\Psi(\mathcal{V}) - \Psi(\bar{\mathcal{V}})\|_{\infty} \leq \gamma \|\mathcal{V} - \bar{\mathcal{V}}\|_{\infty},$$

where $\|\mathcal{V}\|_{\infty} = \max_{x \in \mathcal{X}} \mathcal{V}(x)$. We omit the details since similar results can be easily found in the literature [27]. Then, by Banach's fixed-point theorem [28], we conclude that the operator Ψ has a unique fixed point which satisfies $\mathcal{V}^* = \Psi\mathcal{V}^*$.

Next, we need to show that the fixed point \mathcal{V}^* is indeed the value of the entropy-regularized γ -discounted game. Let $\sigma^* = (\sigma^*, \sigma^*, \dots)$ be a stationary strategy such that σ^* is a one-step strategy for player 1 satisfying the fixed point equation, and τ be an arbitrary Markovian strategy for player 2. Denoting by h_t the history of play of length t , one has, by definition of Ψ and σ^* ,

$$\begin{aligned} & \mathbb{E}^{\bar{\mu}_t} \left[\mathcal{R}(x_t, u_t, w_t) - \frac{1}{\beta_1} \log \sigma^*(u_t | x_t) + \frac{1}{\beta_2} \log \tau(w_t | x_t) \right. \\ & \left. + \gamma \sum_{x' \in \mathcal{X}} \mathcal{P}(x' | x_t, u_t, w_t) \mathcal{V}^*(x') \middle| h_t \right] \geq \mathbb{E}^{\bar{\mu}_t} \left[\mathcal{V}^*(x_t) \middle| h_t \right]. \end{aligned}$$

This expression can further be written as

$$\begin{aligned} & \mathbb{E}^{\bar{\mu}_t, \bar{\mu}_{t+1}} \left[\mathcal{R}(x_t, u_t, w_t) - \frac{1}{\beta_1} \log \sigma^*(u_t | x_t) \right. \\ & \left. + \frac{1}{\beta_2} \log \tau(w_t | x_t) + \gamma \mathcal{V}^*(x_{t+1}) \middle| h_t \right] \geq \mathbb{E}^{\bar{\mu}_t} \left[\mathcal{V}^*(x_t) \middle| h_t \right]. \end{aligned}$$

Multiplying by γ^{t-1} , taking expectations and summing over $1 \leq t \leq k$, one obtains

$$\begin{aligned} & \sum_{t=1}^{k-1} \gamma^{t-1} \mathbb{E}^{\bar{\mu}_t} \left[\mathcal{R}(x_t, u_t, w_t) - \frac{1}{\beta_1} \log \sigma^*(u_t | x_t) + \right. \\ & \left. \frac{1}{\beta_2} \log \tau(w_t | x_t) \middle| x_1 \right] \geq \mathcal{V}^*(x_1) - \gamma^k \mathbb{E}^{\bar{\mu}_{k+1}} \left[\mathcal{V}^*(x_{k+1}) \middle| x_1 \right]. \end{aligned}$$

Taking the limit as $k \rightarrow \infty$ and using Proposition 1, we obtain

$$\Phi_{\infty}(\sigma^*, \tau) \geq \mathcal{V}^*(x_1). \quad (17)$$

Similarly, when player 2 plays the optimal stationary strategy $\tau^* = (\tau^*, \tau^*, \dots)$ against an arbitrary Markovian strategy σ of player 1, we have

$$\Phi_{\infty}(\sigma, \tau^*) \leq \mathcal{V}^*(x_1). \quad (18)$$

Then, the combination of (17) and (18) implies the following result.

Theorem 2: Entropy-regularized γ -discounted games have a value which satisfies $\Psi(\mathcal{V}) = \mathcal{V}$. Furthermore, stationary strategies are sufficient for both players to attain the game value, i.e.,

$$\max_{\sigma \in \Gamma} \min_{\tau \in \Gamma} \Phi_{\infty}(\sigma, \tau) = \max_{\sigma \in \Gamma^S} \min_{\tau \in \Gamma^S} \Phi_{\infty}(\sigma, \tau).$$

Computation of optimal strategies is just an extension of Algorithm 1. Note that for γ -discounted games, we use the same one-shot game introduced in (8). Therefore, optimal decision rules at each stage has the form (9) and (10). Consequently, to compute the optimal strategies, we initialize Algorithm 1 with an arbitrary value vector $\mathcal{V} \in \mathbb{R}^{|\mathcal{X}|}$ and iterate until convergence, which is guaranteed by the existence of a unique fixed point.

V. A NUMERICAL EXAMPLE

In this section, we demonstrate the proposed strategy synthesis method on a motion planning scenario that we model as an entropy-regularized γ -discounted game. To solve the convex optimization problems required for the computation of equilibrium strategies, we use ECOS solver [29] through the interface of CVXPY [30]. All computations are performed by setting $\gamma = 0.8$.

As the environment model, we consider a 5×5 grid world which is given in Figure 1 (top left). The brown grid denotes the initial position of the player 1 which aims to reach the goal (green) state. The red grid is the initial position of the player 2 whose aim is to catch the player 1 before reaching the goal state. Finally, black grids represent walls.

Let $x = (s_1, s_2)$ be the current state of the game such that $x[1] = s_1$ and $x[2] = s_2$ are the positions of the player 1 and the player 2, respectively. At each state, the action space for both players is given as $\mathcal{U} = \mathcal{W} = \{\text{right}, \text{left}, \text{up}, \text{down}, \text{stay}\}$. For simplicity, we assume deterministic transitions, i.e., $\mathcal{P}(x, u, w) \in \{0, 1\}$ for all $x \in \mathcal{X}$, $u \in \mathcal{U}$ and $w \in \mathcal{W}$. If a player takes an action for which the successor state is a wall, the player stays in the same state with probability 1.

For a given $(x, u, w) \in \mathcal{X} \times \mathcal{U} \times \mathcal{W}$, we encode the payoff function $\mathcal{R}(x, u, w)$ as the sum of two functions such that $\mathcal{R}(x, u, w) = \mathcal{R}_1(x, u, w) + \mathcal{R}_2(x, u, w)$ where

$$\begin{aligned} \mathcal{R}_1(x, u, w) &= \sum_{x'[1]=G} \mathcal{P}(x' | x, u, w), \\ \mathcal{R}_2(x, u, w) &= \begin{cases} -\mathcal{P}(x' | x, u, w) & \text{if } x'[1] = x'[2] \neq G \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Note that the payoff function defines a zero-sum game which is won by the player 1, if it reaches the goal state before getting caught, and by the player 2, if it catches the player 1 before reaching the goal state.

We first compute Nash equilibrium strategies in the absence of causal entropy terms, i.e., $\beta_1 = \beta_2 = \infty$, by employing standard linear programming formulation [31] for zero-sum games. Starting from the initial state, an equilibrium strategy for the player 1 is to move towards the goal state by taking the action *right* deterministically, and for the player 2 is to chase the player 1 by taking the action *up* in the first two stages, and then, take the action *right* until reaching the goal state. Therefore, a perfectly rational player 1 wins the game with probability 1 no matter what strategy is followed by the player 2.

To illustrate the drawback of playing with perfect rationality, we assume that there is another wall in the environment

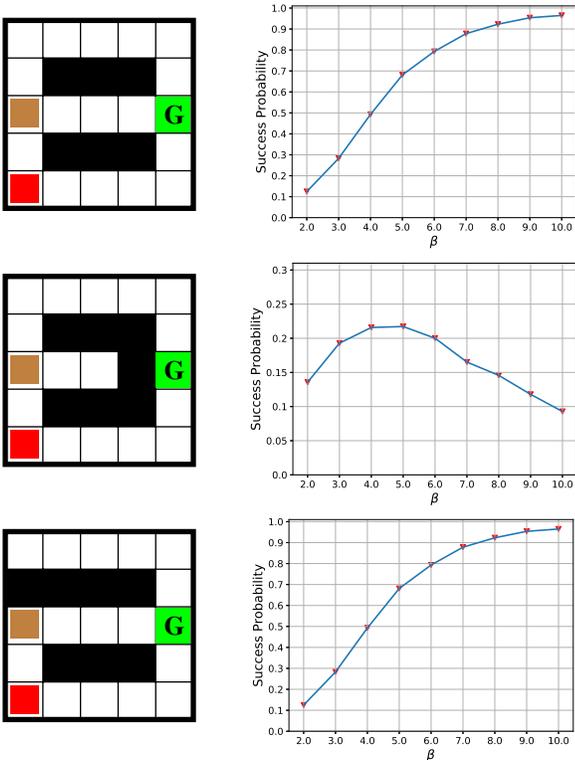


Fig. 1: (Top left) The nominal environment players use for computing their strategies. (Top right) The probability that the player 1 wins the game when it plays the strategy computed by using $\beta_1 = \beta_2 = \beta$ against the perfectly rational player 2. (Middle and bottom left) The actual environments where the game is played. (Middle and bottom right) The probability of winning for the player 1 when it employs strategies computed by using different β values against the perfectly rational player 2.

about which the players have no information while they compute the equilibrium strategies, i.e., the players use the nominal environment (top left) to compute the equilibrium strategies. First, we consider the case that the wall is between the goal state and the player 1, as shown in Figure 1 (middle left). In this case, if the player 1 follows the Nash equilibrium strategy, the probability that it reaches the goal state becomes zero. Therefore, following the Nash equilibrium strategy makes the player 1 significantly vulnerable to such changes in the environment.

To investigate the tradeoff between rationality and randomness, we compute 9 different strategies for player 1 by using $\beta_1 = \beta_2 = 2, 3, \dots, 10$, and let it play against the perfectly rational player 2 which follows its Nash equilibrium strategy computed on the nominal environment (top left). The winning probabilities of player 1 under different strategies are shown in Figure 1 (right) for the corresponding environments given in Figure 1 (left). This specific scenario demonstrates that, by choosing $\beta = 6$, the player 1 can obtain a robust behavior against unforeseen changes in the environment, i.e., the winning probability is around 20%, without sacrificing

too much from its optimal performance, i.e., around 15%, if the structure of the environment remains the same. It is worth noting that the asymptotical performance of the player 1 as $\beta \rightarrow \infty$ approaches to its performance under Nash equilibrium strategy as discussed in [14]. Additionally, the importance of the randomness for the player 1 increases as $\beta \rightarrow 0$, and using smaller β values negatively affects the performance after a critical point, i.e., $\beta = 4$.

Finally, one can argue that the tradeoff occurs in this specific scenario only if the unexpected wall is between the player 1 and the goal state. To justify the choice of β value, we also consider the scenario in which the unexpected wall occupies another state which is shown in Figure 1 (bottom left). In this case, as shown in Figure 1 (bottom right), the use of $\beta = 6$ result in a strategy that guarantees around 80% winning probability. Therefore, the entropy-regularized strategy of the player 1 still provides an advantage against unpredicted changes without sacrificing too much from the optimal performance.

VI. CONCLUSIONS AND FUTURE WORK

We consider the problem of two-player zero-sum stochastic games with entropy regularization, wherein the players aim to maximize the causal entropy of their strategies in addition to the conventional expected payoff. We show that equilibrium strategies exist for both entropy-regularized N -stage and entropy-regularized γ -discounted games, and can be computed by solving a convex optimization problem. In numerical examples, we applied the proposed approach to a motion planning scenario and observed that by tuning the regularization parameter, a player can synthesize robust strategies that perform well in different environments against a perfectly rational opponent.

Extending this work to multi-agent reinforcement learning settings, as discussed in [8], is an interesting future direction. Future work can also investigate the effect of entropy regularization term on the convergence rate of learning algorithms as discussed in [18].

REFERENCES

- [1] L. S. Shapley, "Stochastic games," *Proceedings of the National Academy of Sciences*, vol. 39, no. 10, pp. 1095–1100, 1953.
- [2] J. Nash, "Non-cooperative games," *Annals of Mathematics*, vol. 54, no. 2, pp. 286–295, 1951.
- [3] J. K. Goeree, C. A. Holt, and T. R. Palfrey, *Quantal Response Equilibrium: A Stochastic Theory of Games*. Princeton University Press, 2016.
- [4] E. T. Jaynes, "Information theory and statistical mechanics," *Physical Review*, vol. 106, no. 4, p. 620, 1957.
- [5] B. D. Ziebart, J. A. Bagnell, and A. K. Dey, "The principle of maximum causal entropy for estimating interacting processes," *IEEE Transactions on Information Theory*, vol. 59, no. 4, pp. 1966–1980, 2013.
- [6] —, "Modeling interaction via the principle of maximum causal entropy," in *Proceedings of the 27th International Conference on Machine Learning*, 2010, pp. 1255–1262.
- [7] Y. Savas, M. Ornik, M. Cubuktepe, and U. Topcu, "Entropy maximization for Markov decision processes under temporal logic constraints," *arXiv:1807.03223v2 [math.OA]*, 2018.
- [8] J. Grau-Moya, F. Leibfried, and H. Bou-Amman, "Balancing two-player stochastic games with soft Q-learning," in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, 2018.

- [9] M. Ahmadi, S. Bharadwaj, T. Tanaka, and U. Topcu, "Stochastic games with sensing costs," in *56th Annual Allerton Conference on Communication, Control, and Computing*, 2018.
- [10] P. Mertikopoulos and W. H. Sandholm, "Learning in games via reinforcement and regularization," *Mathematics of Operations Research*, vol. 41, no. 4, pp. 1297–1324, 2016.
- [11] D. S. Leslie and E. J. Collins, "Individual Q-learning in normal form games," *SIAM Journal on Control and Optimization*, vol. 44, no. 2, pp. 495–514, 2005.
- [12] C. K. Ling, J. Z. Kolter, and F. Fang, "What game are we playing? Differentiably learning games from incomplete observations," in *Proceedings of Conference on Neural Information Processing Systems*, 2017.
- [13] D. Fudenberg and D. Levine, *The Theory of Learning in Games*. The MIT Press, 1998.
- [14] R. D. McKelvey and T. R. Palfrey, "Quantal response equilibria for normal form games," *Games and Economic Behavior*, vol. 10, no. 1, pp. 6–38, 1995.
- [15] E. Kardeş, F. Ordóñez, and R. W. Hall, "Discounted robust stochastic games and an application to queueing control," *Operations Research*, vol. 59, no. 2, pp. 365–382, 2011.
- [16] M. Aghassi and D. Bertsimas, "Robust game theory," *Mathematical Programming*, vol. 107, no. 1-2, pp. 231–273, 2006.
- [17] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, "Reinforcement learning with deep energy-based policies," in *Proceedings of the Thirty-fourth International Conference on Machine Learning*, 2017.
- [18] R. Fox, A. Pakman, and N. Tishby, "Taming the noise in reinforcement learning via soft updates," in *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, 2016, pp. 202–211.
- [19] E. Todorov, "Efficient computation of optimal actions," *Proceedings of the National Academy of Sciences*, vol. 106, no. 28, pp. 11 478–11 483, 2009.
- [20] T. Bewley and E. Kohlberg, "The asymptotic theory of stochastic games," *Mathematics of Operations Research*, vol. 1, no. 3, pp. 197–208, 1976.
- [21] S. Sorin, "Discounted stochastic games: The finite case," in *Stochastic Games and Applications*. Springer, 2003, pp. 51–55.
- [22] G. Kramer, "Directed information for channels with feedback," Ph.D. dissertation, ETH Zurich, 1998.
- [23] B. D. Ziebart, J. A. Bagnell, and A. K. Dey, "Maximum causal entropy correlated equilibria for Markov games," in *International Conference on Autonomous Agents and Multiagent Systems*, 2011, pp. 207–214.
- [24] J. v. Neumann, "Zur theorie der gesellschaftsspiele," *Mathematische Annalen*, vol. 100, no. 1, pp. 295–320, 1928.
- [25] R. T. Rockafellar and R. J.-B. Wets, *Variational analysis*. Springer Science & Business Media, 2009, vol. 317.
- [26] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- [27] A. Neyman and S. Sorin, *Stochastic games and applications*. Springer Science & Business Media, 2003, vol. 570.
- [28] M. L. Puterman, *Markov decision processes: Discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [29] A. Domahidi, E. Chu, and S. Boyd, "ECOS: An SOCP solver for embedded systems," in *European Control Conference (ECC)*, 2013, pp. 3071–3076.
- [30] S. Diamond and S. Boyd, "CVXPY: A Python-embedded modeling language for convex optimization," *Journal of Machine Learning Research*, vol. 17, no. 83, pp. 1–5, 2016.
- [31] P. B. Miltersen and T. B. Sørensen, "Computing proper equilibria of zero-sum games," in *International Conference on Computers and Games*, 2006, pp. 200–211.
- [32] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2006.

APPENDIX A

Proof of Proposition 1: We show the sufficiency of Markovian strategies only for the maximin problem. The proof for the minimax formulation follows the same lines with the arguments provided below.

The proof is based on backward induction on the stage

number $1 \leq k \leq N$. Let

$$\mathcal{V}_k := \max_{\sigma \in \Gamma_u} \min_{\tau \in \Gamma_w} \sum_{l=k}^N \mathbb{E}^{\mu^l} \left[\mathcal{R}(X_l, U_l, W_l) - \frac{\log \sigma_l(U_l|H_l)}{\beta_1} + \frac{\log \tau_l(W_l|H_l)}{\beta_2} \right] + \mathbb{E}^{\mu^{T+1}} \mathcal{R}(X_{N+1})$$

be the value of the $N-k$ stage problem. Then, we can write the value of $N-k$ stage problem recursively as

$$\mathcal{V}_k = \max_{\sigma_k} \min_{\tau_k} \mathbb{E}^{\mu^k} \left[\mathcal{R}(X_k, U_k, W_k) - \frac{1}{\beta_1} \log \sigma_k(U_k|H_k) + \frac{1}{\beta_2} \log \tau_k(W_k|H_k) + \mathbb{E}^{\mathcal{P}}[\mathcal{V}_{k+1}] \right]. \quad (19)$$

Base step: $k=N$. Let σ_N and τ_N^* be an arbitrary strategy for player 1 and the optimal strategy for player 2 at stage N , respectively. Let

$$\lambda_N(h_N, u_N, w_N) := \mu_N(x^N, u^{N-1}, w^{N-1}) \times \sigma_N(u_N|h_N) \tau_N^*(w_N|h_N)$$

be the joint distribution induced by $\sigma_N(u_N|h_N)$ and $\tau_N^*(w_N|h_N)$. Additionally, let $\lambda_N(x_N, w_N)$ and $\lambda_N(x_N)$ be the marginal distributions of $\lambda_N(h_N, u_N, w_N)$. We construct a new strategy for player 2 as $\bar{\tau}_N(w_N|x_N) := \frac{\lambda_N(x_N, w_N)}{\lambda_N(x_N)}$. Let

$$\bar{\lambda}_N(h_N, u_N, w_N) := \mu_N(x^N, u^{N-1}, w^{N-1}) \times \sigma_N(u_N|h_N) \bar{\tau}_N(w_N|x_N)$$

be the joint distribution induced by $\bar{\tau}_N(w_N|x_N)$. Then, by construction, we have $\bar{\lambda}_N(x_N, u_N, w_N) = \lambda_N(x_N, u_N, w_N)$, which can be easily verified by calculating the corresponding marginal distributions. (For a similar strategy construction, see Theorem 5.5.1 in [28].)

The inner optimization problem in (19) for $k=N$ reads

$$\mathcal{V}_N = \min_{\tau_N} J_N^c(\lambda_N) - J_N^H(\lambda_N) \quad (20)$$

where

$$J_N^c(\lambda_N) = \mathbb{E}^{\lambda_N, \mathcal{P}} [\mathcal{R}(X_N, U_N, W_N) + \mathcal{R}(X_{N+1}) - \frac{1}{\beta_1} \log \sigma_N(U_N|H_N)]$$

$$J_N^H(\lambda_N) = \frac{1}{\beta_2} H_{\lambda_N}(W_N|X^N, U^{N-1}, W^{N-1}).$$

Since the strategy σ_N is arbitrarily chosen, it is sufficient to show that $J_N^c(\lambda_N) = J_N^c(\bar{\lambda}_N)$ and $J_N^H(\lambda_N) \leq J_N^H(\bar{\lambda}_N)$ in order to establish the sufficiency of Markovian strategies for player 2. The first equality holds by construction. (Note that the $\log(\cdot)$ term is indifferent to changes in the strategy of player 2.) The second inequality can be derived as

$$H_{\lambda_N}(W_N|X^N, U^{N-1}, W^{N-1}) \leq H_{\lambda_N}(W_N|X_N) \quad (21)$$

$$= H_{\bar{\lambda}_N}(W_N|X_N) \quad (22)$$

$$= H_{\bar{\lambda}_N}(W_N|X^N, U^{N-1}, W^{N-1}) \quad (23)$$

where (21) holds since conditioning reduces entropy [32], (22) is because $\lambda_N = \bar{\lambda}_N$ by construction, and (23) is due to the fact that $\bar{\tau}_N(w_N|x_N)$ is a Markovian strategy. Consequently, for any strategy chosen by player 1 in stage

N , player 2 has a best response strategy in the space of Markovian strategies.

Next, we can assume that player 2 uses a Markovian strategy and show through a similar strategy construction explained above that player 1 has an optimal strategy in the space of Markovian strategies. As a result, the value \mathcal{V}_N depends on the joint distribution $\mu_N(x^N, u^{N-1}, w^{N-1})$ only through its marginal $\mu_N(x_N)$ and becomes a function of the marginal $\mu_N(x_N)$ only.

Inductive step: $k=t$ Assume that Markovian strategies suffice for both players for $k=t+1, t+2, \dots, N$. Then, by induction hypothesis, \mathcal{V}_{t+1} is a function of $\mu_{t+1}(x_{t+1})$ only. Therefore, using the similar construction to the case $k=N$, we can construct Markovian strategies $\bar{\sigma}_t(u_t|x_t)$ and $\bar{\tau}_t(w_t|x_t)$ such that the objective function in the right hand side of (19) attained by $\bar{\sigma}_t$ and $\bar{\tau}_t$ is equal to the value of $N-t$ stage problem. As a result, we conclude that Markovian strategies are sufficient for both players to solve the maximin problem (6). \square

Proof of Proposition 3: Since the one-shot game has a value, without loss of generality, we focus on the problem $\max_{\sigma_t \in \Delta(\mathcal{U})} \min_{\tau_t \in \Delta(\mathcal{W})} \mathcal{V}_t(x_t)$. For notational convenience, we rewrite the problem as

$$\begin{aligned} \max_{Q^{ij}} \min_{Q^{ik}} \sum_{jk} Q^{ij} Q^{ik} \rho_{ijk} - \frac{1}{\beta_1} \sum_{jk} Q^{ik} Q^{ij} \log Q^{ij} \\ + \frac{1}{\beta_2} \sum_{jk} Q^{ij} Q^{ik} \log Q^{ik} \end{aligned}$$

$$\text{subject to } \sum_j Q^{ij} = 1, \sum_k Q^{ik} = 1, Q^{ij} \geq 0, Q^{ik} \geq 0,$$

where $Q^{ij} = \sigma_t(u_t|x_t)$, $Q^{ik} = \tau_t(w_t|x_t)$ and $\rho_{ijk} = \rho_t(x_t, u_t, w_t)$. Note that due to constraints $\sum_j Q^{ij} = 1$ and $\sum_k Q^{ik} = 1$, we can replace $\frac{1}{\beta_1} \sum_{jk} Q^{ik} Q^{ij} \log Q^{ij}$ and $\frac{1}{\beta_2} \sum_{jk} Q^{ij} Q^{ik} \log Q^{ik}$ by $\frac{1}{\beta_1} \sum_j Q^{ij} \log Q^{ij}$ and $\frac{1}{\beta_2} \sum_k Q^{ik} \log Q^{ik}$, respectively. For now, we neglect the non-negativity constraints and write the Lagrangian for the above optimization problem as

$$\begin{aligned} L = \sum_{jk} Q^{ij} Q^{ik} \rho_{ijk} - \frac{1}{\beta_1} \sum_j Q^{ij} \log Q^{ij} \\ + \frac{1}{\beta_2} \sum_k Q^{ik} \log Q^{ik} + \lambda^j (\sum_j Q^{ij} - 1) \\ + \lambda^k (\sum_k Q^{ik} - 1) \end{aligned}$$

where λ^j, λ^k are Lagrange multipliers. Then, taking derivative with respect to Q^{ik} and equating it to zero, we obtain

$$\frac{\partial L}{\partial Q^{ik}} = \sum_j Q^{ij} \rho_{ijk} + \frac{1}{\beta_2} \log Q^{ik} + \frac{1}{\beta_2} + \lambda^k = 0.$$

Rearranging terms and using the constraint $\sum_k Q^{ik} = 1$, we obtain

$$Q^{ik} = \frac{\exp(-\beta_2 \sum_j Q^{ij} \rho_{ijk})}{\sum_k \exp(-\beta_2 \sum_j Q^{ij} \rho_{ijk})}$$

which is the same as (10). Note that the resulting strategy also satisfies the non-negativity constraint. Plugging Q^{ik} into Lagrangian L , we obtain the optimization problem given in (9) for which the optimal variables correspond to the optimal strategy of player 1. Similarly, the optimal value (11) of the resulting optimization problem is the value of the game. Uniqueness of the value follows from the fact that the value of the game is the optimal value of a convex optimization problem given in (9). \square