

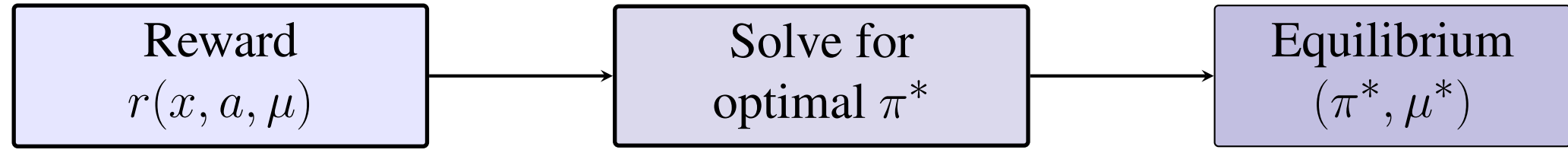
Average-Cost Mean-Field Games: Forward and Inverse Perspectives

Güneş Akbaş¹ and Şevket Kaan Alkır²

gunes.akbas@bilkent.edu.tr¹, kaan.alkir@bilkent.edu.tr²

Bilkent University - Department of Mathematics

1 Forward Reinforcement Learning



Mean-field Term

Let $(x_{1,t}, \dots, x_{N,t})$ be exchangeable. Then

$$\mu_t^N := \frac{1}{N} \sum_{j=1}^N \delta_{x_{j,t}} \rightarrow \mu_t \quad (\text{a.s.}) \text{ as } N \rightarrow \infty,$$

under Assumption (1). And conditioned on μ_t , agents are independent.

Representative-Agent Problem (Average Cost)

Fix a stationary population law μ (or a stationary flow $\{\mu_t\}$ that has converged). The representative agent chooses a stationary policy π to minimise the long-run average cost

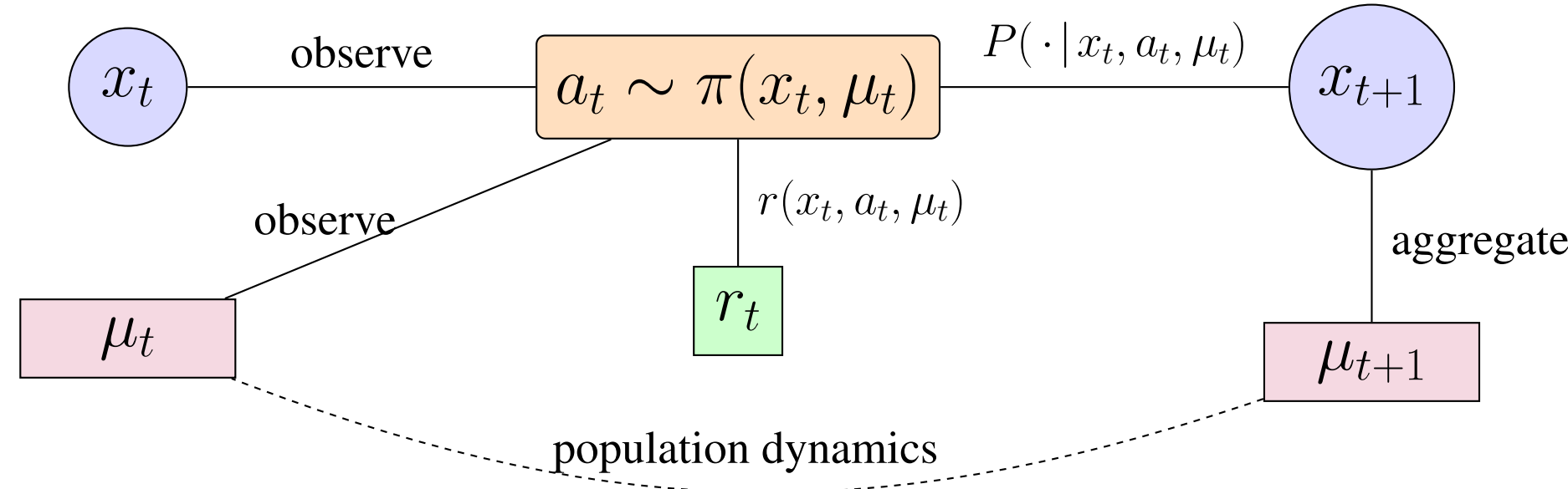
$$V_\mu(\pi, \mu_0) = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}^\pi \left[r(x_t, a_t, \mu) \right] \quad \text{where } x_0 \sim \mu_0.$$

Here, dynamics follows $x_{t+1} \sim P(\cdot | x_t, a_t, \mu)$, $a_t \sim \pi(\cdot | x_t)$. Now define

$$\rho(\mu) := \max_\pi V_\mu(\pi, \mu_0).$$

This criterion will serve as our reference for optimality.

Mean-Field Game — Schematic



State $X_t \in S$ (discrete)
Population $\mu_t \in \mathcal{P}(S)$ empirical distribution of states
Action $A_t \in \mathcal{A}(X_t)$ chosen by policy $\pi(\cdot | X_t, \mu_t)$
Reward $R_t = r(X_t, A_t, \mu_t)$
Dynamics $X_{t+1} \sim P(\cdot | X_t, A_t, \mu_t)$
Coupling: each agent uses (X_t, μ_t) , while μ_t aggregates everyone's states.

2 Mean-Field Average-Cost Soft Bellman Equation

Mean-Field Soft Average-Cost Bellman

For a fixed population law μ and parameter $\delta > 0$:

$$\rho(\mu) + h(x; \mu) = \max_{\pi(\cdot|x)} \left\{ \sum_a \pi(a | x) \left[r(x, a; \mu) + \sum_y P(y | x, a; \mu) h(y; \mu) \right] + \delta H(\pi(\cdot | x)) \right\},$$

where $H(\pi) := -\sum_a \pi(a) \log \pi(a)$.

Softmax Closed Form

$$\rho(\mu) + h(x; \mu) = \delta \log \sum_a \exp \left(\frac{1}{\delta} \left[r(x, a; \mu) + \sum_y P(y | x, a; \mu) h(y; \mu) \right] \right).$$

Optimal Stochastic Policy

$$\pi^*(a | x; \mu) = \frac{\exp \left(\frac{1}{\delta} \left[r(x, a; \mu) + \sum_y P(y | x, a; \mu) h(y; \mu) \right] \right)}{\sum_{a'} \exp \left(\frac{1}{\delta} \left[r(x, a'; \mu) + \sum_y P(y | x, a'; \mu) h(y; \mu) \right] \right)}.$$

3 Fixed-Point Iteration for Mean-Field Equilibrium

(1) Best-Response / Softmax Policy For a population measure μ and entropy weight $\delta > 0$, define

$$\begin{aligned} \pi_\mu(\cdot | x) &= \text{softmax}_{\frac{1}{\delta}}(Q_\mu(x, \cdot)) \\ &= \arg \max_{\pi(\cdot|x) \in \Delta(\mathcal{A})} \left[\sum_{a \in \mathcal{A}} \pi(a | x) (r(x, a; \mu) + \sum_{y \in \mathcal{X}} P(y | x, a; \mu) h_\mu(y)) + \delta H(\pi(\cdot | x)) \right] \end{aligned}$$

with h_μ solving the average-cost Bellman equation under μ .

(2) Population-Update Map Using π_μ , propagate the population forward:

$$\mu'(y) = \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} P(y | x, a, \mu) \pi_\mu(a | x) \mu(x), \quad y \in \mathcal{X}.$$

A mean-field equilibrium is a fixed point $\mu^* = \mu'$ of this two-step operator.

Corollary (Policy Lipschitz)

For the soft-max policy

$$\pi_\mu(a | x) \propto \exp(Q_\mu^*(x, a)/\tau),$$

there exists a constant $C_\pi > 0$ such that

$$\|\pi_\mu - \pi_{\mu'}\|_{\ell^1} \leq C_\pi \|\mu - \mu'\|_{\ell^1}.$$

Algorithm 1 Contraction Fixed-Point Iteration for Mean-Field Equilibrium

Require: Initial measure $\mu_0 \in \mathcal{P}(\mathcal{X})$, temperature $\delta > 0$, tolerance $\varepsilon > 0$

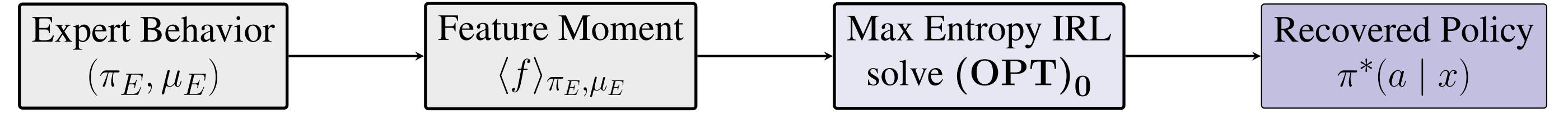
```

1:  $\mu \leftarrow \mu_0$ 
2: repeat
3:    $h \leftarrow \text{BELLMANSOLVE}(\mu)$  // span-contraction step
4:    $\pi \leftarrow \text{SOFTMAX}_{1/\delta}(h)$ 
5:    $\mu' \leftarrow \text{POPUPDATE}(\mu, \pi)$ 
6:    $\Delta \leftarrow \|\mu' - \mu\|_1$ ;  $\mu \leftarrow \mu'$ 
7: until  $\Delta \leq \varepsilon$ 
8: return Equilibrium measure  $\mu$  and policy  $\pi$ 
    
```

References

- [1] Óscar Hernández-Lerma and Jean B. Lasserre (1996), *Discrete-Time Markov Control Processes: Basic Optimality Criteria*, Springer, New York.
- [2] Berkay Anahtarci, Can Deha Kariksiz, and Naci Saldi, *Q-Learning in Regularized Mean-field Games, Dynamic Games and Applications*, vol. 13, no. 1, pp. 89–117, Mar. 2023.
- [3] Bolin Gao and Lacra Pavel, A game-theoretic approach to apprenticeship learning, in *On the Properties of the Softmax Function with Application in Game Theory and Reinforcement Learning*, arXiv preprint arXiv:1704.00805, 2017.

1 Inverse Reinforcement Learning



A widely used assumption in IRL for MFGs is that the reward function can be expressed as a linear combination of feature vectors, which depend on the state, action, and mean-field term:

$$\mathcal{R} := \left\{ r(x, a, \mu) = \langle \theta, f(x, a, \mu) \rangle \mid \theta \in \mathbb{R}^k, f : \mathcal{X} \times \mathcal{A} \times \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}^k \right\}$$

where $f(x, a, \mu) \in \mathbb{R}^k$ is the feature vector that includes information about the state-action pair (x, a) and the mean-field term μ .

Assumption 1 We assume that for any policy π and mean-field term μ , the state process is positive Harris recurrent and aperiodic, guaranteeing the convergence of empirical averages [1, Thm. 13.3.3]. In other words, the state process exhibits ergodic behavior under any policy and mean-field term, ensuring long-run statistical stability. Assuming the expert population distribution μ_E is known is standard, as it can be consistently estimated from long-run trajectories. [1, Thm. 13.3.3]

Definition 1 average-reward expected feature vector under the mean-field equilibrium (π_E, μ_E) for $x_0 \sim \mu_E$ is defined as,

$$\langle f \rangle_{\pi_E, \mu_E} := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}^{\pi_E, \mu_E} [f(x_t, a_t, \mu_E)]. \quad (1)$$

2 Entropy-Regularized Inverse Learning

Definition 2 Average-reward causal entropy $H(\pi)$ of the policy $\pi \in \Pi$ is defined as follows

$$H(\pi) := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}^{\pi, \mu_E} [-\log \pi(a_t | x_t)].$$

This principle selects the least biased policy among those consistent with expert demonstrations [2]. We define the average-cost maximum causal entropy IRL problem:

$$\begin{aligned} (\text{OPT})_0 \text{ maximize}_\pi \quad & H(\pi) \\ \text{subject to} \quad & \pi(a | x) \geq 0 \quad \forall (x, a) \in \mathcal{X} \times \mathcal{A} \\ & \sum_{a \in \mathcal{A}} \pi(a | x) = 1 \quad \forall x \in \mathcal{X} \\ & \mu_E(x) = \sum_{(a,y) \in \mathcal{A} \times \mathcal{X}} p(x | y, a, \mu_E) \pi(a | y) \mu_E(y) \quad \forall x \in \mathcal{X} \\ & \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}^{\pi, \mu_E} [f(x_t, a_t, \mu_E)] = \langle f \rangle_{\pi_E, \mu_E} \end{aligned}$$

3 Convex Reformulation via Occupation Measures

Definition 3 We define the state-action occupation measure ν_π for any policy $\pi \in \Pi$ and corresponding state occupation measure as

$$\nu_\pi(x, a) := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}^{\pi, \mu_E} \left[\mathbf{1}_{\{(x_t, a_t) = (x, a)\}} \right] \quad \nu_\pi^\mathcal{X}(x) := \sum_{a \in \mathcal{A}} \nu_\pi(x, a).$$

Under ergodicity, ν_π defines a valid probability measure over $\mathcal{X} \times \mathcal{A}$.

Lemma 1 Suppose $\pi \in \Pi$ is a feasible point for the optimization problem (OPT). Then:

1. **State Marginals:** $\nu_\pi^\mathcal{X}(x) = \mu_E(x) = \sum_{(y,a)} p(x | y, a, \mu_E) \nu_\pi(y, a)$ for all $x \in \mathcal{X}$.
2. **Entropy:** $H(\pi) = \sum_{(x,a)} -\log \left(\frac{\nu_\pi(x,a)}{\mu_E(x)} \right) \nu_\pi(x, a)$.
3. **Feature Expectation:** $\langle f \rangle_{\pi, \mu_E} = \sum_{(x,a)} f(x, a, \mu_E) \nu_\pi(x, a)$.

In view of above results, let us now reformulate (OPT)₀ in terms of state-action occupation measure [3]. Although (OPT)₀ is nonconvex, it can be reformulated as an equivalent convex program over occupation measures. This new formulation is denoted by (OPT), which is convex and yield the same solution as (OPT)₀:

$$\begin{aligned} (\text{OPT}) \text{ maximize}_\nu \quad & \sum_{(x,a) \in \mathcal{X} \times \mathcal{A}} -\log \left(\frac{\nu(x,a)}{\mu_E(x)} \right) \nu(x, a) \\ \text{subject to} \quad & \sum_{(x,a) \in \mathcal{X} \times \mathcal{A}} f(x, a, \mu_E) \nu(x, a) = \langle f \rangle_{\pi_E, \mu_E} \\ & \mu_E(x) = \sum_{(y,a) \in \mathcal{X} \times \mathcal{A}} p(x | y, a, \mu_E) \nu(y, a) \quad \forall x \in \mathcal{X} \\ & \nu^\mathcal{X}(x) = \mu_E(x) \quad \forall x \in \mathcal{X} \\ & \nu(x, a) \geq 0 \quad \forall (x, a) \in \mathcal{X} \times \mathcal{A}. \end{aligned}$$

The optimization problem (OPT) is convex: its objective combines a strongly concave entropy term with a linear component, and all constraints are linear in the occupation measure ν . Using convex duality [4] and Sion's minimax theorem [5], the problem (OPT) admits the dual formulation:

$$\min_{\alpha \in \mathbb{R}^k, \beta, \theta \in \mathbb{R}^\mathcal{X}} \left\{ \log \sum_{(x,a) \in \mathcal{X} \times \mathcal{A}} e^{k_{\alpha, \beta, \theta}(x, a)} - \langle \alpha, \langle f \rangle_{\pi_E, \mu_E} \rangle - \sum_{x \in \mathcal{X}} \theta_x \mu_E(x) \right\}$$

Here, $k_{\alpha, \beta, \theta}(x, a) := \langle \alpha, f(x, a, \mu_E) \rangle + \theta_x + \sum_{z \in \mathcal{X}} \beta_z (p(z | x, a, \mu_E) - \mu_E(z))$. For fixed dual variables, the inner maximization over ν is solved by the Boltzmann distribution, and the corresponding recovered policy is given by:

$$\nu^*(x, a) := \frac{e^{k_{\alpha, \beta, \theta}(x, a)}}{Z_{\alpha, \beta, \theta}}, \quad \pi^*(a | x) := \frac{\nu^*(x, a)}{\mu_E(x)}$$

which uniquely maximizes the entropy-regularized objective. This yields a fully tractable dual with no duality gap.

Theorem 1 The dual objective h is L -smooth and $\rho(D)$ -strongly convex over any compact subset $D \subset \mathbb{R}^m$, where $m = k + 2|\mathcal{X}|$. Specifically,

$$L := 2M^2 \sqrt{|\mathcal{X}| \cdot |\mathcal{A}|},$$

with M a uniform bound on the gradients of the function $k_{\alpha, \beta, \theta}(x, a)$. Moreover, if the set

$$\{(f(x, a, \mu_E), p(\cdot | x, a, \mu_E), e(\cdot | x, a)) : (x, a) \in \mathcal{X} \times \mathcal{A}\}$$

spans $\mathbb{R}^k \times \mathbb{R}^\mathcal{X} \times \mathbb{R}^\mathcal{X}$, then h is uniformly strongly convex on D .

Now we can introduce the gradient descent algorithm for finding the minimizer of h as follows.

Algorithm 1 Gradient Descent for Parameter Optimization

Require: $(\alpha_0, \beta_0, \theta_0) \in \mathbb{R}^m$ and $\delta \in \left(0, \frac{1}{L}\right]$ and $K \in \mathbb{N}$

- 1: Set $(\alpha, \beta, \theta) \leftarrow (\alpha_0, \beta_0, \theta_0)$
- 2: **for** $k = 0, 1, \dots, K-1$ **do**
- 3: $(\alpha, \beta, \theta) \leftarrow (\alpha, \beta, \theta) - \delta \nabla h(\alpha, \beta, \theta)$
- 4: **end for**
- 5: Compute $\nu_{\alpha, \beta, \theta}^*$ from (α, β, θ)
- 6: **return** (α, β, θ) and $\nu_{\alpha, \beta, \theta}^*$

References

- [1] Meyn, S. P. & Tweedie, R. L. (1993), *Markov Chains and Stochastic Stability*, Springer-Verlag London, Stochastic Modelling and Applied Probability.
- [2] Ziebart, B. D. (2008), *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*, Ph.D. thesis, Carnegie Mellon University.
- [3] Syed, U. & Schapire, R. E. (2007), A game-theoretic approach to apprenticeship learning, in *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 1449–1456.
- [4] Dupuis, P. & Ellis, R. S. (1997), Formulation of Large Deviation Theory in Terms of the Laplace Principle. In: *A Weak Convergence Approach to the Theory of Large Deviations*, Wiley-Interscience, ch. 1, pp. 1–47.
- [5] Sion, M. (1958), On general minimax theorems, *Pacific Journal of Mathematics*, 8, pp. 171–176.