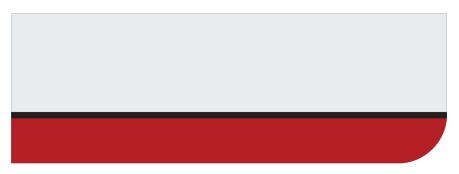


Multiple-timescales learning in games

February 2025 David S. Leslie



Introduction



- Prof. of Statistical Learning at Lancaster since 2014. Previously mathematics department at University of Bristol
- PI on EPSRC Data Science of the Natural Environment project (2018–2023)
- Researcher on EPSRC/BT Next Generation Converged Digital Infrastructure (2018–2023)
- Was consultant at Prowler.io (now Secondmind.ai), 2018–2020.





... to mathematical sciences at Lancaster

MARS (Maths for AI in Real-world Systems

- £15M investment to expand mathematical sciences at Lancaster (focus is AI especially with applications in health, environment, engineering, cybersecurity)
- 10 new permanent positions, 4 still to recruit (all levels)
- 8 post-doc positions, recruiting 2025/26

ProbAl research hub

- £10M to build collaborations across multiple universities and industry, focusing on probabilistic techniques for AI
- Recruiting post-docs imminently

Introduction



- Introduction
- Stochastic fictitious play and stochastic approximation
- Two-timescales stochastic approximation
- Applications:
 - Actor-critic learning
 - Player-dependent learning rates
 - Learning in stochastic games
 - Noise reduction in gradient estimation





Often we might want to run an inner loop between adaptations:

- Clinical trials Several treatments. Experiment enough with each treatment. Adapt the set of treatments and repeat.
 - Games Fix the (mixed) strategies. Play long enough to learn the strategies. Adapt the strategies and repeat.
- Deep learning Fix the weights. Gather enough observations with these weights. Adapt the weights and repeat.





Often we might want to run an inner loop between adaptations:

- Clinical trials Several treatments. Experiment enough with each treatment. Adapt the set of treatments and repeat.
 - Games Fix the (mixed) strategies. Play long enough to learn the strategies. Adapt the strategies and repeat.
- Deep learning Fix the weights. Gather enough observations with these weights. Adapt the weights and repeat.
 - Generally A system has parameters θ and a performance gradient $v(\theta)$. If v is not analytically available, fix θ for long enough to reliably estimate $v(\theta)$ on the basis of observations, update θ and repeat.

Two timescales helps to avoid "fix" and "enough"



Ფ◍ᅄ

$$\begin{array}{cccc} R & S & P \\ R & \begin{pmatrix} 0,0 & (1,-1) & (-1,1) \\ (-1,1) & (0,0) & (1,-1) \\ (1,-1) & (-1,1) & (0,0) \end{pmatrix} \end{array}$$

- Finite set of players, labeled i
- Each player has an action space A^i ; joint action space $A = A^1 \times \cdots \times A^N$
- Usually we consider mixed strategies πⁱ ∈ Δ(Aⁱ); joint mixed strategies π ∈ Δ
- Reward functions extend to $r^i : \Delta \to \mathbb{R}$

Normal form games Equilibrium



Response to (beliefs about) other players becomes key. Define the best response correspondence

$$b^i(\pi^{-i}) = \operatorname*{argmax}_{\pi^i \in \Delta(A^i)} r^i(\pi^i, \pi^{-i})$$

{Nash equilibria} := {fixed points $\pi^i \in b^i(\pi^{-i})$ }

Normal form games Equilibrium



Response to (beliefs about) other players becomes key. Define the best response correspondence

$$b^i(\pi^{-i}) = \operatorname*{argmax}_{\pi^i \in \Delta(\mathcal{A}^i)} r^i(\pi^i, \pi^{-i})$$

{Nash equilibria} := {fixed points $\pi^i \in b^i(\pi^{-i})$ }

Often we want a continuous response. The canonical example is the smooth best response function $\beta_{\tau}^{i}(\pi^{-i})$ satisfying

$$eta^i_{ au}(\pi^{-i})(a^i)\propto \exp(r^i(a^i,\pi^{-i})/ au)$$

{fixed points $\pi^{i} = \beta^{i}_{\tau}(\pi^{-i})$ } =: {smoothed Nash equilibrium}

Normal form games Fictitious play



Even if the game is fully known, things are non-trivial!

Fictitious play

- Repeatedly play the game
- On iteration *n*, estimate πⁱ by σⁱ_n, the empirically observed distribution of opponent actions so far
- Play a best response $a_{n+1}^i \in b^i(\sigma_n^{-i})$

$$\sigma_{n+1}^{i}(a^{i}) = \frac{1}{n+1} \sum_{m=1}^{n+1} \mathbb{I}_{\{a_{m}^{i}=a^{i}\}} = \sigma_{n}^{i}(a^{i}) + \frac{1}{n+1} \left[\mathbb{I}_{\{a_{n+1}^{i}=a^{i}\}} - \sigma_{n}^{i}(a^{i}) \right]$$

Normal form games Fictitious play



Even if the game is fully known, things are non-trivial!

Fictitious play

- Repeatedly play the game
- On iteration *n*, estimate πⁱ by σⁱ_n, the empirically observed distribution of opponent actions so far
- Play a best response $a_{n+1}^i \in b^i(\sigma_n^{-i})$

$$\sigma_{n+1}^{i}(a^{i}) = \frac{1}{n+1} \sum_{m=1}^{n+1} \mathbb{I}_{\{a_{m}^{i}=a^{i}\}} = \sigma_{n}^{i}(a^{i}) + \frac{1}{n+1} \left[\mathbb{I}_{\{a_{n+1}^{i}=a^{i}\}} - \sigma_{n}^{i}(a^{i}) \right]$$

$$\sigma_{n+1}^{i} = \sigma_{n}^{i} + \frac{1}{n+1} \left[b^{i}(\sigma_{n}^{-i}) - \sigma_{n}^{i} \right]$$

Normal form games Fictitious play



Even if the game is fully known, things are non-trivial!

Fictitious play Stochastic fictitious play

- Repeatedly play the game
- On iteration *n*, estimate πⁱ by σⁱ_n, the empirically observed distribution of opponent actions so far
- Play a best response $a_{n+1}^i \in b^i(\sigma_n^{-i})$ Play $a_{n+1}^i \sim \beta^i(\sigma_n^{-i})$

$$\sigma_{n+1}^{i}(a^{i}) = \frac{1}{n+1} \sum_{m=1}^{n+1} \mathbb{I}_{\{a_{m}^{i}=a^{i}\}} = \sigma_{n}^{i}(a^{i}) + \frac{1}{n+1} \left[\mathbb{I}_{\{a_{n+1}^{i}=a^{i}\}} - \sigma_{n}^{i}(a^{i}) \right]$$

$$\sigma_{n+1}^{i} = \sigma_{n}^{i} + \frac{1}{n+1} \left[\beta^{i}(\sigma_{n}^{-i}) - \sigma_{n}^{i} + M_{n+1}^{i} \right]$$



$$\theta_{t+1} = \theta_t + \alpha_{t+1} \{ F(\theta_t) + \boldsymbol{e}_t + \boldsymbol{M}_{t+1} \}$$

- Robbins-Monro
- Kiefer-Wolfowicz
- Ljung
- Kushner
- Benveniste, Metivier and Priouret
- Duflo
- Borkar
- Benaïm



$$\theta_{t+1} = \theta_t + \alpha_{t+1} \left\{ F(\theta_t) + \boldsymbol{e}_t + \boldsymbol{M}_{t+1} \right\}$$

Rearrange:

$$\frac{\theta_{t+1} - \theta_t}{\alpha_t} = F(\theta_t) + \boldsymbol{e}_t + \boldsymbol{M}_{t+1}$$



$$\theta_{t+1} = \theta_t + \alpha_{t+1} \{ F(\theta_t) + \boldsymbol{e}_t + \boldsymbol{M}_{t+1} \}$$

Rearrange:

$$\frac{\theta_{t+1} - \theta_t}{\alpha_t} = F(\theta_t) + \underline{e_t} + \underline{M_{t+1}}$$

Looks like a discretisation of

 $\dot{\theta} = F(\theta).$



$$\theta_{t+1} = \theta_t + \alpha_{t+1} \left\{ F(\theta_t) + \boldsymbol{e}_t + \boldsymbol{M}_{t+1} \right\}$$

Rearrange:

$$\frac{\theta_{t+1} - \theta_t}{\alpha_t} = F(\theta_t) + \underline{e_t} + \underline{M_{t+1}}$$

Looks like a discretisation of

$$\dot{\theta} = F(\theta).$$

Theorem (ish)

If the ODE has a unique globally attracting fixed point θ^* then the stochastic approximation iterates converge almost surely to θ^*

Normal form games Smooth best response dynamics



Recall stochastic fictitious play (SFP):

$$\sigma_{t+1} = \sigma_t + \frac{1}{t+1} \left\{ \beta(\sigma_t) - \sigma_t + M_{t+1} \right\}$$

This is a stochastic approximation with $F(\sigma_t) = \beta(\sigma_t) - \sigma_t$

Hence SFP converges if the smooth best response dynamics

$$\dot{\sigma} = \beta(\sigma) - \sigma$$

are globally convergent

Convergence in zero-sum-games, potential games, some other less obvious classes (Benaïm and Hirsch, Hofbauer, others)

Normal form games Radically uncoupled



Suppose can't observe opponent actions and don't know the payoff matrix. Now what?!

Normal form games Radically uncoupled



Suppose can't observe opponent actions and don't know the payoff matrix. Now what?!

Each player now faces a bandit problem \Rightarrow Use RL in bandits approach \Rightarrow Individual *Q*-learning (Leslie and Collins 2006)

Normal form games Radically uncoupled



Suppose can't observe opponent actions and don't know the payoff matrix. Now what?!

Each player now faces a bandit problem \Rightarrow Use RL in bandits approach \Rightarrow Individual *Q*-learning (Leslie and Collins 2006)

Can mixed strategies behave like fictitious play beliefs

$$\pi_{t+1} = \pi_t + \alpha_{t+1} \{\beta(\pi_t) - \pi_t\}?$$

Yes, if each player can calculate $\beta^i(\pi_t^{-i}) \propto \exp(r^i(\cdot, \pi_t^{-i})/\tau)$

Normal form games Estimating $r^i(\cdot, \pi^{-i})$



- "Wait everybody, don't move your π^i , we're all going to observe for a while"
- Play repeatedly and estimate rⁱ(aⁱ, π⁻ⁱ) to be the average reward obtained with *i* play action aⁱ
- When these have converged, everybody adjust π^i a little bit

Normal form games Estimating $r^i(\cdot, \pi^{-i})$



- "Wait everybody, don't move your π^i , we're all going to observe for a while"
- Play repeatedly and estimate rⁱ(aⁱ, π⁻ⁱ) to be the average reward obtained with *i* play action aⁱ
- When these have converged, everybody adjust π^i a little bit

Like in fictitious play, the averages can be calculated "online":

$$Q_{n+1}^{i}(a^{i}) = Q_{n}^{i}(a^{i}) + rac{\mathbb{I}_{\{a_{n}^{i}=a^{i}\}}}{\kappa_{n}^{i}(a^{i})} \left\{ R_{n}^{i} - Q_{n}^{i}(a^{i})
ight\}$$

Stochastic approximation ^{Mathematics} Lancaster & Statistics University Two timescales (Borkar 1997)

Two SA processes, with $\alpha_n/\gamma_n \rightarrow 0$

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \{ F(\theta_n, \phi_n) + e_n + M_{n+1} \}$$

$$\phi_{n+1} = \phi_n + \gamma_{n+1} \{ G(\theta_n, \phi_n) + h_n + N_{n+1} \}$$

Stochastic approximation ^{Mathematics} Lancaster & Statistics University Two timescales (Borkar 1997)

Two SA processes, with $\alpha_n/\gamma_n \rightarrow 0$

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \{ F(\theta_n, \phi_n) + e_n + M_{n+1} \}$$

$$\phi_{n+1} = \phi_n + \gamma_{n+1} \{ G(\theta_n, \phi_n) + h_n + N_{n+1} \}$$

Rewrite as a single SA, with learning parameters γ_n

$$\begin{pmatrix} \theta_{n+1} \\ \phi_{n+1} \end{pmatrix} = \begin{pmatrix} \theta_n \\ \phi_n \end{pmatrix} + \gamma_{n+1} \begin{pmatrix} \frac{\alpha_{n+1}}{\gamma_{n+1}} \left\{ F(\theta_n, \phi_n) + e_n + M_{n+1} \right\} \\ G(\theta_n, \phi_n) + h_n + N_{n+1} \end{pmatrix}$$

Stochastic approximation ^{Mathematics} Lancaster & Statistics University Two timescales (Borkar 1997)

Two SA processes, with $\alpha_n/\gamma_n \to 0$

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \{ F(\theta_n, \phi_n) + e_n + M_{n+1} \}$$

$$\phi_{n+1} = \phi_n + \gamma_{n+1} \{ G(\theta_n, \phi_n) + h_n + N_{n+1} \}$$

Rewrite as a single SA, with learning parameters γ_n

$$\begin{pmatrix} \theta_{n+1} \\ \phi_{n+1} \end{pmatrix} = \begin{pmatrix} \theta_n \\ \phi_n \end{pmatrix} + \gamma_{n+1} \begin{pmatrix} 0 + \tilde{e}_n \\ G(\theta_n, \phi_n) + h_n + N_{n+1} \end{pmatrix}$$

The approximated differential equation is

$$\begin{pmatrix} \theta \\ \phi \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ G(\theta, \phi) \end{pmatrix}$$

Stochastic approximation Mathematics & Statistics University

$$\begin{pmatrix} \dot{\theta} \\ \phi \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ G(\theta, \phi) \end{pmatrix}$$

Assumption:

For each θ there is a unique, globally attracting, fixed point of the "fast ODE" $\dot{\phi} = G(\theta, \phi)$. Call this $\phi^*(\theta)$.

Under this assumption, the set

$$\left\{ \begin{pmatrix} \theta \\ \phi^{\star}(\theta) \end{pmatrix} \, : \, \theta \in \Theta \right\}$$

is globally attracting; $\binom{\theta_n}{\phi_n}$ converges to this set

Stochastic approximation Authematics Statistics University

We have shown that $\phi_n = \phi^*(\theta_n) + \epsilon_n$. So

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \{ F(\theta_n, \phi_n) + e_n + M_{n+1} \}$$

= $\theta_n + \alpha_{n+1} \{ F(\theta_n, \phi^*(\theta_n) + \epsilon_n) + e_n + M_{n+1} \}$
= $\theta_n + \alpha_{n+1} \{ F(\theta_n, \phi^*(\theta_n)) + \eta_n + e_n + M_{n+1} \}$

The "slow ODE" is

$$\dot{ heta} = F(heta, \phi^{\star}(heta))$$

If the fast and slow ODEs both converge, then we're in business!



Put the inner loop estimation of $r^i(\cdot, \pi^{-i})$ on the fast timescale:

$$\pi_{n+1} = \pi_n + \alpha_{n+1} \{\beta(Q_n) - \pi_n + M_{n+1}\}$$
$$Q_{n+1}^i(a^i) = Q_n^i(a^i) + \gamma_{n+1} \mathbb{I}_{\{a_n^i = a^i\}} \{R_{n+1}^i - Q_n^i(a^i)\}$$

with learning parameters such that $\frac{\alpha_n}{\gamma_n} \rightarrow 0$



Put the inner loop estimation of $r^i(\cdot, \pi^{-i})$ on the fast timescale:

$$\pi_{n+1} = \pi_n + \alpha_{n+1} \{\beta(Q_n) - \pi_n + M_{n+1}\}$$
$$Q_{n+1}^i(a^i) = Q_n^i(a^i) + \gamma_{n+1} \mathbb{I}_{\{a_n^i = a^i\}} \{R_{n+1}^i - Q_n^i(a^i)\}$$

with learning parameters such that $\frac{\alpha_n}{\gamma_n} \rightarrow 0$

Fast timescale: fix π

$$\dot{Q}^{i}(a^{i}) = \pi^{i}(a^{i})\left\{r^{i}(a^{i},\pi^{-i}) - Q^{i}(a^{i})\right\}$$

This ODE converges: $Q^i(a^i) \rightarrow r^i(a^i, \pi^{-i}) =: Q^{\star,i}(\pi)(a^i)$

Therefore Q_n^i will be close to $Q^{\star,i}(\pi_n)$ for large n



Put the inner loop estimation of $r^i(\cdot, \pi^{-i})$ on the fast timescale:

$$\pi_{n+1} = \pi_n + \alpha_{n+1} \{\beta(Q_n) - \pi_n + M_{n+1}\}$$
$$Q_{n+1}^i(a^i) = Q_n^i(a^i) + \gamma_{n+1} \mathbb{I}_{\{a_n^i = a^i\}} \{R_{n+1}^i - Q_n^i(a^i)\}$$

with learning parameters such that $\frac{\alpha_n}{\gamma_n} \rightarrow 0$

Slow timescale: analyse as if $Q^{i}(a^{i}) = r^{i}(a^{i}, \pi^{-i})$

$$\dot{\pi}^i = \beta(\boldsymbol{Q}^{\star,i}(\pi)) - \pi^i) = \beta(\boldsymbol{r}^i(\cdot,\pi^{-i})) - \pi^i = \beta(\pi^{-i}) - \pi^i$$

which is the smooth best response dynamics

The actor–critic algorithm converges in the same games as stochastic fictitious play



Revert to stochastic fictitious play, two player games:

$$\sigma_{n+1}^{1} = \sigma_{n}^{1} + \frac{1}{n+1} \left\{ \beta^{1}(\sigma_{n}^{2}) - \sigma_{n}^{1} + M_{n+1}^{1} \right\}$$
$$\sigma_{n+1}^{2} = \sigma_{n}^{2} + \frac{1}{n+1} \left\{ \beta^{2}(\sigma_{n}^{1}) - \sigma_{n}^{2} + M_{n+1}^{2} \right\}$$



Revert to stochastic fictitious play, two player games:

$$\sigma_{n+1}^{1} = \sigma_{n}^{1} + \alpha_{n+1} \left\{ \beta^{1}(\sigma_{n}^{2}) - \sigma_{n}^{1} + M_{n+1}^{1} \right\}$$

$$\sigma_{n+1}^{2} = \sigma_{n}^{2} + \alpha_{n+1} \left\{ \beta^{2}(\sigma_{n}^{1}) - \sigma_{n}^{2} + M_{n+1}^{2} \right\}$$



Revert to stochastic fictitious play, two player games:

$$\sigma_{n+1}^{1} = \sigma_{n}^{1} + \alpha_{n+1} \left\{ \beta^{1}(\sigma_{n}^{2}) - \sigma_{n}^{1} + M_{n+1}^{1} \right\}$$
$$\sigma_{n+1}^{2} = \sigma_{n}^{2} + \gamma_{n+1} \left\{ \beta^{2}(\sigma_{n}^{1}) - \sigma_{n}^{2} + M_{n+1}^{2} \right\}$$

with $\alpha_n/\gamma_n \rightarrow 0$

$$\sigma_{n+1}^{1} = \sigma_{n}^{1} + \alpha_{n+1} \left\{ \beta^{1}(\sigma_{n}^{2}) - \sigma_{n}^{1} + M_{n+1}^{1} \right\}$$

$$\sigma_{n+1}^{2} = \sigma_{n}^{2} + \gamma_{n+1} \left\{ \beta^{2}(\sigma_{n}^{1}) - \sigma_{n}^{2} + M_{n+1}^{2} \right\}$$

Fast timescale: fix σ^1

$$\dot{\sigma}^2 = \beta^2(\sigma^1) - \sigma^2$$

 $\sigma^2 \to \beta^2(\sigma^1)$



$$\sigma_{n+1}^{1} = \sigma_{n}^{1} + \alpha_{n+1} \left\{ \beta^{1}(\sigma_{n}^{2}) - \sigma_{n}^{1} + M_{n+1}^{1} \right\}$$

$$\sigma_{n+1}^{2} = \sigma_{n}^{2} + \gamma_{n+1} \left\{ \beta^{2}(\sigma_{n}^{1}) - \sigma_{n}^{2} + M_{n+1}^{2} \right\}$$

Mathematics & Statistics University

Slow timescale: analyse as if $\sigma_n^2 = \beta^2(\sigma_n^1)$

$$\dot{\sigma}^1 = \beta^1(\beta^2(\sigma^1)) - \sigma^1$$

- This ODE has a globally attracting fixed point for zero-sum games, potential games and Shapley's game
- The ODE falls outside Hart and Mas-Colell's impossibility framework
- I have yet to find a game in which it does not converge

$$\sigma_{n+1}^{1} = \sigma_{n}^{1} + \alpha_{n+1}^{1} \left\{ \beta^{1}(\sigma_{n}^{-1}) - \sigma_{n}^{1} + M_{n+1}^{1} \right\}$$

$$\sigma_{n+1}^{2} = \sigma_{n}^{2} + \alpha_{n+1}^{2} \left\{ \beta^{2}(\sigma_{n}^{-2}) - \sigma_{n}^{2} + M_{n+1}^{2} \right\}$$

$$\sigma_{n+1}^{3} = \sigma_{n}^{3} + \alpha_{n+1}^{3} \left\{ \beta^{3}(\sigma_{n}^{-3}) - \sigma_{n}^{3} + M_{n+1}^{3} \right\}$$

$$\sigma_{n+1}^{4} = \sigma_{n}^{4} + \alpha_{n+1}^{4} \left\{ \beta^{4}(\sigma_{n}^{-4}) - \sigma_{n}^{4} + M_{n+1}^{4} \right\}$$

$$\vdots \quad \text{with } \alpha_{n}^{i} / \alpha_{n}^{i+1} \to 0$$

Mathematics & Statistics University

Theorem-ish

If the fast strategies $\sigma^{>i}$ converge to a unique $\beta^{>i}(\sigma^{\leq i})$ for fixed $\sigma^{\leq i}$, for each *i*, then the system converges iff $\dot{\sigma}^1 = \beta^1(\beta^{>1}(\sigma^1)) - \sigma^1$ converges

Stochastic games Setup



Stochastic game framework (Shapley 1953):

- Finite set of players $i \in \{1, \ldots, N\}$
- Finite set of states $s \in S$
- Finite set of actions $A^{i}(s)$ for each player *i* in each state s
- Transitions $P_{s,s'}(a)$ and rewards $r^i(s,a)$ for $a = (a^1, \ldots, a^N)$
- · Players attempt to maximise cumulative discounted reward

Key concept: auxiliary games

At each state s, all players choose actions, receive reward and move to next state. Next state has 'continuation payoffs' $V^i(s')$. Auxiliary game at s, with continuation payoffs V has payoff matrix

$$q_{s,\mathbf{V}}^{i}(a) = r^{i}(s,a) + \delta \sum_{s'} P_{ss'}(a) V^{i}(s')$$

Learning in stochastic gamestics Lancaster Introduction

"Normal-forming" (Stochastic game strategy ↔ Normal form action): Not v interesting! Finding a best response is solving an MDP. Mixed strategies are weird, except perhaps in evolutionary interpretation.

Per-state fictitious play: This can work (Sayin, Parise, Ozdaglar, SICON 2022 building on Leslie, Perkins, Xu, JET 2020). But can we do radically-uncoupled learning?

Simple *Q*-learning: Many hint this is solved. It is not!

Learning in stochastic games tistics Lancaster With University Key idea

Challenge: There are many moving parts. State values we are yet to receive are affected by current strategies

Solution (ish): Fixing the "continuation payoffs" and learning in just the "auxiliary games" makes things much easier

Finishing off: If the auxiliary games are all played 'at' equilibrium, then the state values will converge

Learning in stochastic games Lancaster Reinforcement learning

$$V_n^i(s) = \max_{a^i} Q_n^i(s, a^i)$$

$$Q_{n+1}^{i}(s, a^{i}) = Q_{n}^{i}(s, a^{i}) + \gamma_{n} \mathbb{I}_{\{(s_{n}, a_{n}^{i}) = (s, a)\}} \left\{ r_{n}^{i} + \delta V_{n}^{i}(s_{n+1}) - Q_{n}^{i}(s, a) \right\}$$

Learning in stochastic games Lancaster Reinforcement learning

$$V_n^i(s) = \sum_b \pi_n^i(s,b) Q_n^i(s,b)$$

$$\begin{aligned} Q_{n+1}^{i}(s,a^{i}) &= Q_{n}^{i}(s,a^{i}) + \\ & \frac{\gamma_{n}}{\pi_{n}^{i}(s,a)} \mathbb{I}_{\{(s_{n},a_{n}^{i})=(s,a)\}} \left\{ r_{n}^{i} + \delta V_{n}^{i}(s_{n+1}) - Q_{n}^{i}(s,a) \right\} \end{aligned}$$

where $\pi_n^i(s, a) \propto \exp(Q_n^i(s, a)/\tau_n)$

Learning in stochastic games Lancaster Reinforcement learning

$$V_{n+1}^i(s) = V_n^i(s) + \alpha_n \mathbb{I}_{\{s_n=s\}} \left\{ \sum_b \pi_n^i(s,b) Q_n^i(s,b) - V_n^i(s) \right\}$$

$$Q_{n+1}^{i}(s, a^{i}) = Q_{n}^{i}(s, a^{i}) + \frac{\gamma_{n}}{\pi_{n}^{i}(s, a)} \mathbb{I}_{\{(s_{n}, a_{n}^{i}) = (s, a)\}} \left\{ r_{n}^{i} + \delta V_{n}^{i}(s_{n+1}) - Q_{n}^{i}(s, a) \right\}$$

where $\pi_n^i(s, a) \propto \exp(Q_n^i(s, a)/\tau_n)$ and $\alpha_n/\gamma_n \to 0$

Learning in stochastic gamestistics Lancaster Decoupling step

Two-timescale approach decouples the states:

	$Q_{n+1}^i(s,a) - Q_n^i(s,a)$		$\int q_n^i(s,(a,\pi_n^{-i}(s))) - Q_n^i(s,a)$	
E	$V_{n+1}^i(s) - V_n^i(s)$	$= \alpha$	0	+ <i>e</i> _{n+1}
	$\tau_{n+1} - \tau_n$		0	

where $q_n^i(s, \boldsymbol{a}) = r^i(s, \boldsymbol{a}) + \delta \sum_{s'} P_{ss'}(\boldsymbol{a}) V_n^i(s')$.

This fast timescale corresponds to considering "individual Q-learning" (Leslie and Collins 2005) in an arbitrary fixed auxiliary game with payoffs $q^i(s, \cdot)$.



New Lyapunov function (fast timescale)

- $\dot{Q}^i(a) = q^i(a, \pi^{-i}) Q^i(a)$ with $\pi^i(a) \propto \exp(Q^i(a)/\tau)$
- Introduce auxiliary vars σ^i defined by $\dot{\sigma}^i = \pi^i \sigma^i$.
- New Lyapunov function:

$$\begin{split} L(\boldsymbol{Q}^1, \boldsymbol{Q}^2, \sigma^1, \sigma^2) &= \left[\sum_{i=1,2} \left\{ \pi^i \cdot \boldsymbol{Q}^i + \tau \boldsymbol{v}^i(\pi^i) \right\} - \lambda \zeta \right]_+ \\ &+ \sum_{i=1,2} \| \boldsymbol{Q}^i - \boldsymbol{q}^i(\cdot, \sigma^{-i}) \|^2 \end{split}$$

where $\lambda \in (1, \gamma^{-1})$ and $\zeta = \|q^1 + (q^2)^T\|_{\max} + \tau \log(|A^1||A^2|)$.



Learning in stochastic games statistics U

$$\begin{split} \mathcal{L}(\mathcal{Q}^1, \mathcal{Q}^2, \sigma^1, \sigma^2) &= \left[\sum_{i=1,2} \left\{ \pi^i \cdot \mathcal{Q}^i + \tau \mathbf{v}^i(\pi^i) \right\} - \lambda \zeta \right]_+ \\ &+ \sum_{i=1,2} \| \mathcal{Q}^i - \mathcal{q}^i(\cdot, \sigma^{-i}) \|^2 \end{split}$$

- Start with standard Lyapunov function for smooth BR learning
- $\lambda \zeta$ term means we only make this small, not 0
- Second summation shows Q are asymptotically belief based



Learning in stochastic gamestistics New Lyapunov function (fast timescale)

$$\begin{split} L(\boldsymbol{Q}^1, \boldsymbol{Q}^2, \sigma^1, \sigma^2) &= \left[\sum_{i=1,2} \left\{ \pi^i \cdot \boldsymbol{Q}^i + \tau \boldsymbol{v}^i(\pi^i) \right\} - \lambda \zeta \right]_+ \\ &+ \sum_{i=1,2} \| \boldsymbol{Q}^i - \boldsymbol{q}^i(\cdot, \sigma^{-i}) \|^2 \end{split}$$

So, there exists $\epsilon_n \rightarrow 0$, such that

$$\sum_{i=1,2} \left\{ \pi_n^i \cdot \mathcal{Q}_n^i + \tau_n \mathbf{v}^i(\pi_n^i) \right\} \le \lambda \left\{ \|\mathcal{Q}_n^1 + (\mathcal{Q}_n^2)^T\|_{\max} + \tau_n \log(|\mathcal{A}^1||\mathcal{A}^2|) \right\} + \epsilon_n$$

Stochastic games



- For fixed continuation payoffs *V*, we have shown convergence (admittedly to a set)
- The two-timescales theory allows us to analyse *V* as if the *Q* values are always in this set
- Convergence follows in two-player zero-sum games

Refs:

- Leslie, Perkins, Xu, JET 2020
- Sayin, Parise, Ozdaglar, SICON 2022
- Sayin, Zhang, Leslie, Basar, Ozduglar, NeurIPS 2020

In continuous games, we use a very different notation:

- Actions are $x = (x^1, \dots, x^N)$
- Payoffs are $u^i : \mathcal{X} \to \mathbb{R}$
- Individual payoff gradients are $v^i(x^i, x^{-i}) = \nabla_{x^i} u^i(x)$
- Pseudogradient is $v(x) = (v^1(x), \dots, v^N(x))$

Often players need to estimate $v^{i}(x)$. Estimates may have very high variance.

Averaging several observations $v^i(x) + \epsilon_n^i$ would reduce the variance So....

$$\begin{aligned} x_{n+1}^{i} &= x_{n}^{i} + \alpha_{n+1} V_{n}^{i} \\ V_{n+1}^{i} &= V_{n}^{i} + \gamma_{n+1} \left\{ v_{n}^{i}(x_{n}^{i}) + \epsilon_{n}^{i} - V_{n}^{i} \right\} \end{aligned}$$

For fixed x, the V_n^i converge to v_n^i . Then the slow equation follows the gradient nicely.



Rate of convergence is the elephant in the two-timescales room!

My student Miles Elvidge is working on some really cool ideas along these lines

Essentially, work out what the fast timescale analysis tells you, then plug that into a finite time analysis on the slow timescale





- Whenever an inner loop would be useful, think about using two timescales
- Has been deployed in:
 - actor–critic learning
 - player-dependent learning rates
 - stochastic games
 - gradient smoothing
- Convergence rates are hard, but very recent work is getting there



February 2025 David S. Leslie

