



Decentralized Stochastic Approximation, Optimization, and Multi-Agent Reinforcement Learning

Justin Romberg, Georgia Tech ECE
CAMDA/TAMIDS Seminar, Texas A& M
College Station, Texas Streaming live from Atlanta, Georgia
March 16, 2020 October 30, 2020

Collaborators



Think Doan
Virginia Tech, ECE

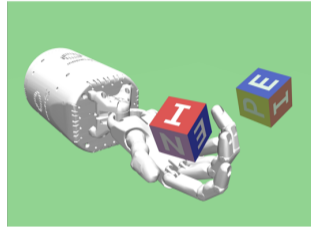


Siva Theja Maguluri
Georgia Tech, ISyE



Sihan Zeng
Georgia Tech, ECE

Reinforcement Learning



Distributed RL is a combination of:

- stochastic approximation
- Markov decision processes
- function representation
- network consensus

Distributed RL is a combination of:

- stochastic approximation
- Markov decision processes
- function representation
- network consensus
- (complicated probabilistic analysis)

Distributed RL is a combination of:

- **stochastic approximation**
- Markov decision processes
- function representation
- network consensus
- (complicated probabilistic analysis)

Classical result (Banach fixed point theorem): when $H(\cdot) : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is a contraction

$$\|H(\mathbf{u}) - H(\mathbf{v})\| \leq \delta \|\mathbf{u} - \mathbf{v}\|, \quad \delta < 1,$$

then there is a unique **fixed point** \mathbf{x}^* such that

$$\mathbf{x}^* = H(\mathbf{x}^*),$$

and the iteration

$$\mathbf{x}_{k+1} = H(\mathbf{x}_k),$$

finds it

$$\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}^*.$$

Choose any point \mathbf{x}_0 , then take

$$\mathbf{x}_{k+1} = H(\mathbf{x}_k)$$

so

$$\begin{aligned}\mathbf{x}_{k+1} - \mathbf{x}^* &= H(\mathbf{x}_k) - \mathbf{x}^* \\ &= H(\mathbf{x}_k) - H(\mathbf{x}^*)\end{aligned}$$

and

$$\begin{aligned}\|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \|H(\mathbf{x}_k) - H(\mathbf{x}^*)\| \\ &\leq \delta \|\mathbf{x}_k - \mathbf{x}^*\| \\ &\leq \delta^{k+1} \|\mathbf{x}_0 - \mathbf{x}^*\|,\end{aligned}$$

so the convergence is **geometric**

Choose any point \mathbf{x}_0 , then take

$$\mathbf{x}_{k+1} = H(\mathbf{x}_k),$$

then

$$\begin{aligned}\|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \|H(\mathbf{x}_k) - H(\mathbf{x}^*)\| \\ &\leq \delta^{k+1} \|\mathbf{x}_0 - \mathbf{x}^*\|,\end{aligned}$$

Gradient descent takes

$$H(\mathbf{x}) = \mathbf{x} - \alpha \nabla f(\mathbf{x})$$

for some differentiable f .

Take

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha(H(\mathbf{x}_k) - \mathbf{x}_k), \quad 0 < \alpha \leq 1.$$

(More conservative, convex combination of new iterate and old.)

Then again

$$\mathbf{x}_{k+1} = (1 - \alpha)\mathbf{x}_k + \alpha H(\mathbf{x}_k)$$

and

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\| &\leq (1 - \alpha)\|\mathbf{x}_k - \mathbf{x}^*\| + \alpha\|H(\mathbf{x}_k) - H(\mathbf{x}^*)\| \\ &\leq (1 - \alpha - \delta\alpha)\|\mathbf{x}_k - \mathbf{x}^*\|. \end{aligned}$$

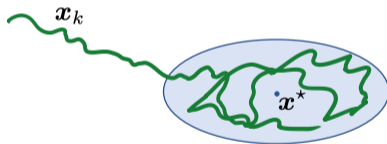
Still converge, albeit a little more slowly for $\alpha < 1$.

What if there is noise?

If our observations of $H(\cdot)$ are *noisy*,

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha (H(\mathbf{x}_k) - \mathbf{x}_k + \boldsymbol{\eta}_k), \quad \mathbb{E}[\boldsymbol{\eta}_k] = \mathbf{0},$$

then we don't get convergence for fixed α ,



but we do converge to a “ball” around at a geometric rate

Stochastic approximation

If our observations of $H(\cdot)$ are *noisy*,

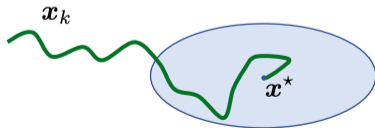
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (H(\mathbf{x}_k) - \mathbf{x}_k + \boldsymbol{\eta}_k), \quad \mathbb{E}[\boldsymbol{\eta}_k] = \mathbf{0},$$

then we need to take $\alpha_k \rightarrow 0$ as we approach the solution.

If we take $\{\alpha_k\}$ such that

$$\sum_{k=0}^{\infty} \alpha_k^2 < \infty, \quad \sum_{k=0}^{\infty} \alpha_k = \infty$$

then we so get (much slower) convergence



Example: $\alpha_k = C/(k+1)$

Distributed RL is a combination of:

- stochastic approximation
- **Markov decision processes**
- function representation
- network consensus
- (complicated probabilistic analysis)

Markov decision process

At time t ,

- 1 An agent finds itself in a **state** s_t
- 2 It takes action $\mathbf{a}_t = \mu(s_t)$
- 3 It moves to state s_{t+1} according to $P(s_{t+1}|s_t, \mathbf{a}_t) \dots$
- 4 ... and receives reward $R(s_t, \mathbf{a}_t, s_{t+1})$.



Markov decision process

At time t ,

- 1 An agent finds itself in a **state** s_t
- 2 It takes action $\mathbf{a}_t = \mu(s_t)$
- 3 It moves to state s_{t+1} according to $P(s_{t+1}|s_t, \mathbf{a}_t) \dots$
- 4 ... and receives reward $R(s_t, \mathbf{a}_t, s_{t+1})$.



Long-term reward of **policy** μ :

$$V_{\mu}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \mu(s_t), s_{t+1}) \mid s_0 = s \right]$$

Markov decision process

At time t ,

- 1 An agent finds itself in a **state** s_t
- 2 It takes action $\mathbf{a}_t = \mu(s_t)$
- 3 It moves to state s_{t+1} according to $P(s_{t+1}|s_t, \mathbf{a}_t) \dots$
- 4 ... and receives reward $R(s_t, \mathbf{a}_t, s_{t+1})$.



Bellman equation: V_μ obeys

$$V_\mu(s) = \underbrace{\sum_{z \in \mathcal{S}} P(z|s, \mu(s)) [R(s, \mu(s), z) + \gamma V_\mu(z)]}_{\mathbf{b}_\mu + \gamma \mathbf{P}_\mu \mathbf{V}_\mu}$$

This is a *fixed point equation* for V_μ

Markov decision process

At time t ,

- 1 An agent finds itself in a **state** s_t
- 2 It takes action $\mathbf{a}_t = \mu(s_t)$
- 3 It moves to state s_{t+1} according to $P(s_{t+1}|s_t, \mathbf{a}_t) \dots$
- 4 ... and receives reward $R(s_t, \mathbf{a}_t, s_{t+1})$.



State-action value function (Q function):

$$Q_{\mu}(\mathbf{s}, \mathbf{a}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}_t, \mu(\mathbf{s}_t) \mathbf{s}_{t+1}) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a} \right]$$

Markov decision process

At time t ,

- 1 An agent finds itself in a **state** s_t
- 2 It takes action $\mathbf{a}_t = \mu(s_t)$
- 3 It moves to state s_{t+1} according to $P(s_{t+1}|s_t, \mathbf{a}_t) \dots$
- 4 ... and receives reward $R(s_t, \mathbf{a}_t, s_{t+1})$.



State-action value for the optimal policy obeys

$$Q^*(s, \mathbf{a}) = \mathbb{E} \left[R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}') \mid s_0 = s, \mathbf{a}_0 = \mathbf{a} \right]$$

and we take $\mu^*(s) = \arg \max_{\mathbf{a}} Q^*(s, \mathbf{a}) \dots$

... this is another *fixed point equation*

Fixed point iteration for finding $V_\mu(\mathbf{s})$:

$$V_{t+1}(\mathbf{s}) = V_t(\mathbf{s}) + \alpha \underbrace{\left(\sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{s}) [R(\mathbf{s}, \mathbf{z}) + \gamma V_t(\mathbf{z})] - V_t(\mathbf{s}) \right)}_{H(\mathbf{V}_t) - \mathbf{V}_t}$$

Fixed point iteration for finding $V_\mu(\mathbf{s})$:

$$V_{t+1}(\mathbf{s}) = V_t(\mathbf{s}) + \alpha \underbrace{\left(\sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{s}) [R(\mathbf{s}, \mathbf{z}) + \gamma V_t(\mathbf{z})] - V_t(\mathbf{s}) \right)}_{H(\mathbf{V}_t) - \mathbf{V}_t}$$

In practice, we don't have the model $P(\mathbf{z}|\mathbf{s})$, only observed data $\{(\mathbf{s}_t, \mathbf{s}_{t+1})\}$

Fixed point iteration for finding $V_\mu(\mathbf{s})$:

$$V_{t+1}(\mathbf{s}) = V_t(\mathbf{s}) + \alpha \underbrace{\left(\sum_{\mathbf{z}} P(\mathbf{z}|\mathbf{s}) [R(\mathbf{s}, \mathbf{z}) + \gamma V_t(\mathbf{z})] - V_t(\mathbf{s}) \right)}_{H(\mathbf{V}_t) - \mathbf{V}_t}$$

Stochastic approximation iteration

$$V_{t+1}(\mathbf{s}_t) = V_t(\mathbf{s}_t) + \alpha_t (R(\mathbf{s}_t, \mathbf{s}_{t+1}) + \gamma V_t(\mathbf{s}_{t+1}) - V_t(\mathbf{s}_t))$$

The “noise” is that \mathbf{s}_{t+1} is sampled, rather than averaged over

Stochastic approximation for policy evaluation

Fixed point iteration for finding $V_\mu(\mathbf{s})$:

$$V_{t+1}(\mathbf{s}) = V_t(\mathbf{s}) + \alpha \underbrace{\left(\sum_z P(z|\mathbf{s}) [R(\mathbf{s}, z) + \gamma V_t(z)] - V_t(\mathbf{s}) \right)}_{H(\mathbf{V}_t) - \mathbf{V}_t}$$

Stochastic approximation iteration

$$V_{t+1}(\mathbf{s}_t) = V_t(\mathbf{s}_t) + \alpha_t \underbrace{(R(\mathbf{s}_t, \mathbf{s}_{t+1}) + \gamma V_t(\mathbf{s}_{t+1}) - V_t(\mathbf{s}_t))}_{H(\mathbf{V}_t) - \mathbf{V}_t + \boldsymbol{\eta}_t}$$

The “noise” is that \mathbf{s}_{t+1} is sampled, rather than averaged over

Stochastic approximation for policy evaluation

Fixed point iteration for finding $V_\mu(\mathbf{s})$:

$$V_{t+1}(\mathbf{s}) = V_t(\mathbf{s}) + \alpha \underbrace{\left(\sum_z P(z|\mathbf{s}) [R(\mathbf{s}, z) + \gamma V_t(z)] - V_t(\mathbf{s}) \right)}_{H(\mathbf{V}_t) - \mathbf{V}_t}$$

Stochastic approximation iteration

$$V_{t+1}(\mathbf{s}_t) = V_t(\mathbf{s}_t) + \alpha_t \underbrace{(R(\mathbf{s}_t, \mathbf{s}_{t+1}) + \gamma V_t(\mathbf{s}_{t+1}) - V_t(\mathbf{s}_t))}_{H(\mathbf{V}_t) - \mathbf{V}_t + \boldsymbol{\eta}_t}$$

The “noise” is that \mathbf{s}_{t+1} is sampled, rather than averaged over

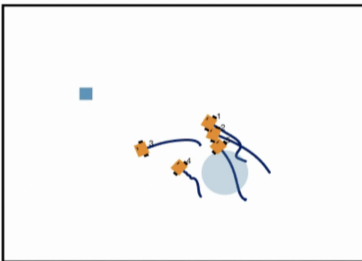
This is different from stochastic gradient descent, since $H(\cdot)$ is in general not a gradient map

Distributed RL is a combination of:

- stochastic approximation
- Markov decision processes
- **function representation**
- network consensus
- (complicated probabilistic analysis)

Function approximation

State space can be large (or even infinite) ...



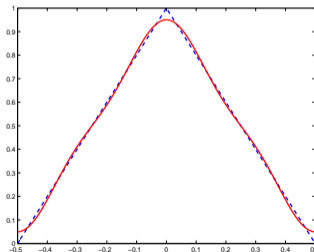
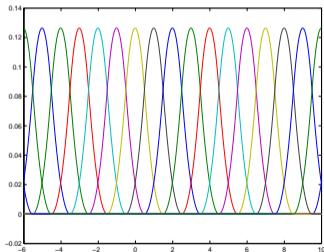
... we need a natural way to parameterize/simplify

Simple (but powerful) model: linear representation

$$V(\mathbf{s}; \boldsymbol{\theta}) = \sum_{k=1}^K \theta_k \phi_k(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\theta}, \quad \boldsymbol{\phi}(\mathbf{s}) = \begin{bmatrix} \phi_1(\mathbf{s}) \\ \vdots \\ \phi_K(\mathbf{s}) \end{bmatrix}$$

Simple (but powerful) model: linear representation

$$V(\mathbf{s}; \boldsymbol{\theta}) = \sum_{k=1}^K \theta_k \phi_k(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\theta}, \quad \boldsymbol{\phi}(\mathbf{s}) = \begin{bmatrix} \phi_1(\mathbf{s}) \\ \vdots \\ \phi_K(\mathbf{s}) \end{bmatrix}$$



Bellman equation:

$$V(\mathbf{s}) = \sum_{\mathbf{z} \in \mathcal{S}} P(\mathbf{z}|\mathbf{s}) [R(\mathbf{s}, \mu(\mathbf{s}), \mathbf{z}) + \gamma V(\mathbf{z})]$$

Linear approximation:

$$V(\mathbf{s}; \boldsymbol{\theta}) = \sum_{k=1}^K \theta_k \phi_k(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\theta}$$

These can conflict

Bellman equation:

$$V(\mathbf{s}) = \sum_{\mathbf{z} \in \mathcal{S}} P(\mathbf{z}|\mathbf{s}) [R(\mathbf{s}, \mu(\mathbf{s}), \mathbf{z}) + \gamma V(\mathbf{z})]$$

Linear approximation:

$$V(\mathbf{s}; \boldsymbol{\theta}) = \sum_{k=1}^K \theta_k \phi_k(\mathbf{s}) = \boldsymbol{\phi}(\mathbf{s})^T \boldsymbol{\theta}$$

These can conflict

... but the following iterations

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t + \alpha_t (R(\mathbf{s}_t, \mathbf{s}_{t+1}) + \gamma V(\mathbf{s}_{t+1}; \boldsymbol{\theta}_t) - V(\mathbf{s}_t; \boldsymbol{\theta}_t)) \nabla_{\boldsymbol{\theta}} V(\mathbf{s}_t, \boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t + \alpha_t (R(\mathbf{s}_t, \mathbf{s}_{t+1}) + \gamma \boldsymbol{\phi}(\mathbf{s}_{t+1})^T \boldsymbol{\theta}_t - \boldsymbol{\phi}(\mathbf{s}_t)^T \boldsymbol{\theta}_t) \boldsymbol{\phi}(\mathbf{s}_t) \end{aligned}$$

converge to a “near optimal” $\boldsymbol{\theta}^*$

Tsitsiklis and Roy, '97

Distributed RL is a combination of:

- stochastic approximation
- Markov decision processes
- function representation
- **network consensus**
- (complicated probabilistic analysis)

- Each node in a network has a number $x(i)$
- We want each node to agree on the average

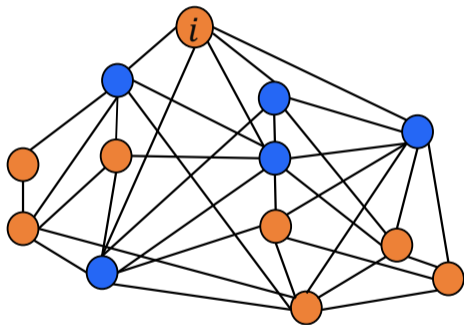
$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x(i) = \mathbf{1}^T \mathbf{x}$$

Node i communicates with its neighbors \mathcal{N}_i

- Iterate, take $\mathbf{v}_0 = \mathbf{x}$, then

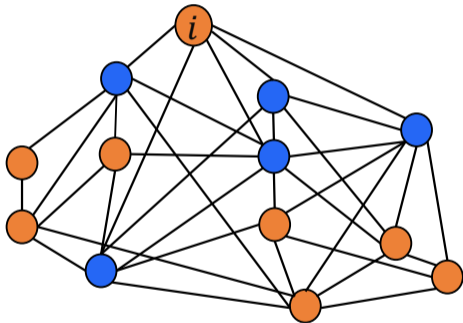
$$v_{k+1}(i) = \sum_{j \in \mathcal{N}_i} W_{ij} v_k(j)$$

$$\mathbf{v}_{k+1} = \mathbf{W} \mathbf{v}_k, \quad \mathbf{W} \text{ doubly stochastic}$$



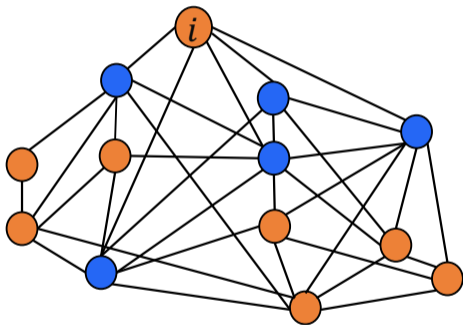
Nodes reach “consensus” quickly:

$$\begin{aligned}\mathbf{v}_{k+1} &= \mathbf{W} \mathbf{v}_k \\ \mathbf{v}_{k+1} - \bar{x} \mathbf{1} &= \mathbf{W} \mathbf{v}_k - \bar{x} \mathbf{1} \\ &= \mathbf{W} (\mathbf{v}_k - \bar{x} \mathbf{1}) \\ \|\mathbf{v}_{k+1} - \bar{x} \mathbf{1}\| &= \|\mathbf{W} (\mathbf{v}_k - \bar{x} \mathbf{1})\|\end{aligned}$$



Nodes reach “consensus” quickly:

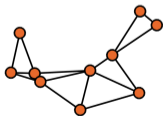
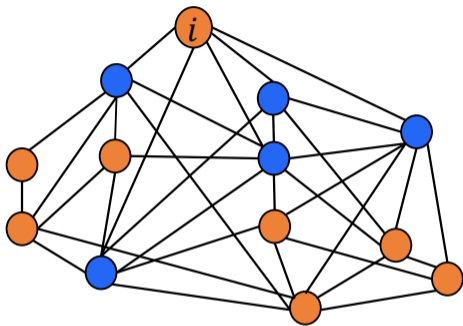
$$\begin{aligned}\mathbf{v}_{k+1} &= \mathbf{W} \mathbf{v}_k \\ \mathbf{v}_{k+1} - \bar{x} \mathbf{1} &= \mathbf{W} \mathbf{v}_k - \bar{x} \mathbf{1} \\ &= \mathbf{W} (\mathbf{v}_k - \bar{x} \mathbf{1}) \\ \|\mathbf{v}_{k+1} - \bar{x} \mathbf{1}\| &= \|\mathbf{W} (\mathbf{v}_k - \bar{x} \mathbf{1})\| \\ &\leq \sigma_2 \|\mathbf{v}_k - \bar{x} \mathbf{1}\| \\ &\leq \sigma_2^{k+1} \|\mathbf{v}_0 - \bar{x} \mathbf{1}\|\end{aligned}$$



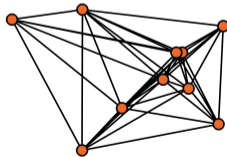
Network consensus convergence

Nodes reach “consensus” quickly:

$$\begin{aligned}\mathbf{v}_{k+1} &= \mathbf{W}\mathbf{v}_k \\ \mathbf{v}_{k+1} - \bar{x}\mathbf{1} &= \mathbf{W}\mathbf{v}_k - \bar{x}\mathbf{1} \\ &= \mathbf{W}(\mathbf{v}_k - \bar{x}\mathbf{1}) \\ \|\mathbf{v}_{k+1} - \bar{x}\mathbf{1}\| &= \|\mathbf{W}(\mathbf{v}_k - \bar{x}\mathbf{1})\| \\ &\leq \sigma_2 \|\mathbf{v}_k - \bar{x}\mathbf{1}\| \\ &\leq \sigma_2^{k+1} \|\mathbf{v}_0 - \bar{x}\mathbf{1}\|\end{aligned}$$



σ larger



σ smaller

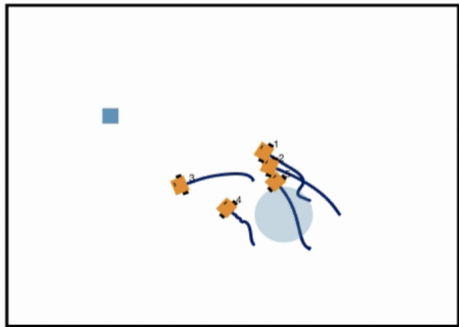
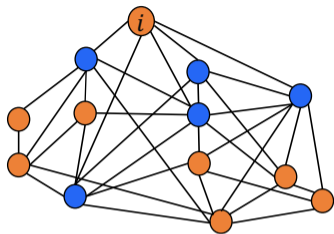
Multi-agent Reinforcement Learning, Scenario 1:

Multiple agents in a single environment, common state, different rewards

What is the value of a particular policy?

Multi-agent reinforcement learning

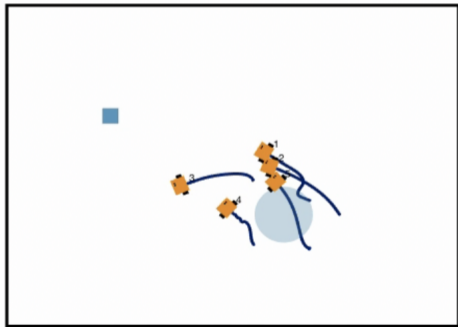
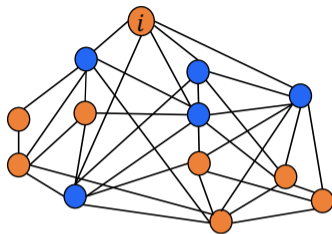
- N agents, communicating on a network
- One environment, common state $\mathbf{s}_t \in \mathcal{S}$
transition probabilities $P(\mathbf{s}_{t+1}|\mathbf{s}_t)$
- Individual actions $\mathbf{a}_t^i \in \mathcal{A}^i$
- Individual rewards $R^i(\mathbf{s}_t, \mathbf{s}_{t+1})$
- evaluate policies $\mu^i : \mathcal{S} \rightarrow \mathcal{A}^i$



Multi-agent reinforcement learning

- N agents, communicating on a network
- One environment, common state $\mathbf{s}_t \in \mathcal{S}$
transition probabilities $P(\mathbf{s}_{t+1}|\mathbf{s}_t)$
- Individual actions $\mathbf{a}_t^i \in \mathcal{A}^i$
- Individual rewards $R^i(\mathbf{s}_t, \mathbf{s}_{t+1})$
- compute average cumulative reward

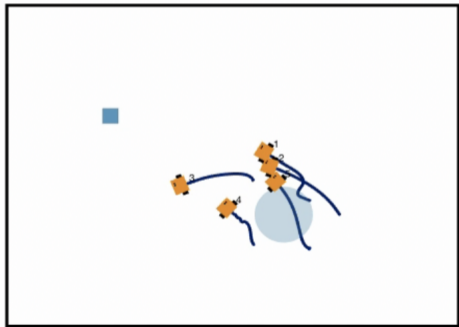
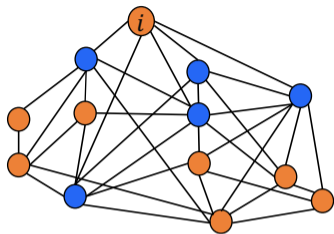
$$V(\mathbf{s}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \frac{1}{N} \sum_{i=1}^N R^i(\mathbf{s}_t, \mathbf{s}_{t+1}) \mid \mathbf{s}_0 = \mathbf{s} \right]$$



Multi-agent reinforcement learning

- N agents, communicating on a network
- One environment, common state $\mathbf{s}_t \in \mathcal{S}$
transition probabilities $P(\mathbf{s}_{t+1}|\mathbf{s}_t)$
- Individual actions $\mathbf{a}_t^i \in \mathcal{A}^i$
- Individual rewards $R^i(\mathbf{s}_t, \mathbf{s}_{t+1})$
- find \mathbf{V} that satisfies

$$V(\mathbf{s}) = \sum_{\mathbf{z} \in \mathcal{S}} P(\mathbf{z}|\mathbf{s}) \left[\frac{1}{N} \sum_{n=1}^N R^i(\mathbf{s}, \mathbf{z}) + \gamma V(\mathbf{z}) \right]$$



Distributed temporal difference learning

Initialize: Each agent starts at θ_0^i

Iterations: Observe: s_t , take action to go to s_{t+1} ,
get reward $R(s_t, s_{t+1})$

Communicate: average estimates from neighbors

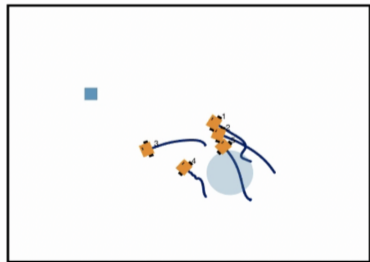
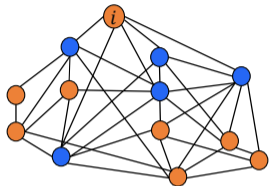
$$\mathbf{y}_t^i = \sum_{j \in \mathcal{N}_i} W_{ij} \theta_t^j$$

Local updates:

$$\theta_{t+1}^i = \mathbf{y}_t^i + \alpha_t d_t^i \phi(s_t),$$

where

$$d_t^i = R^i(s_t, s_{t+1}) + \gamma \phi(s_{t+1})^\top \theta_t^i - \phi(s_t)^\top \theta_t^i$$



Distributed RL is a combination of:

- stochastic approximation
- Markov decision processes
- function representation
- network consensus
- **(complicated probabilistic analysis)**

Subset of existing results:

- Unified convergence theory: Borkar and Meyn '00
- Convergence rates with “independent noise” (centralized):
Thoppe and Borkar '19, Dalal et al '18, Lakshminarayanan and Szepesvari '18
- Convergence rates under Markovian noise (centralized):
Bhandari et al COLT '18. Srikant and Ying COLT '19
- Multi-agent RL: Mathkar and Bokar '17, Zhang et al '18, Kar et al '13, Stankovic and Stankovic '16, Macua et al '15

- Fixed **step size** $\alpha_t = \alpha$, for small enough α

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O(\eta^{t-\tau}) + O(\alpha)$$

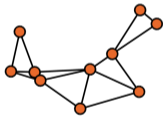
where $\sigma < 1$ is **network connectivity**, $\eta < 1$ are **problem parameters**, and τ is the **mixing time** for the underlying Markov chain

Rate of convergence for distributed TD

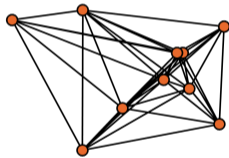
- Fixed **step size** $\alpha_t = \alpha$, for small enough α

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O(\eta^{t-\tau}) + O(\alpha)$$

where $\sigma < 1$ is **network connectivity**, $\eta < 1$ are **problem parameters**, and τ is the **mixing time** for the underlying Markov chain



σ larger



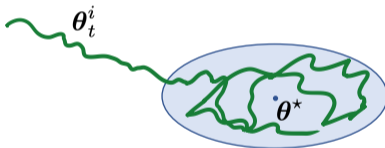
σ smaller

Rate of convergence for distributed TD

- Fixed **step size** $\alpha_t = \alpha$, for small enough α

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O(\eta^{t-\tau}) + O(\alpha)$$

where $\sigma < 1$ is **network connectivity**, $\eta < 1$ are **problem parameters**, and τ is the **mixing time** for the underlying Markov chain



- Fixed **step size** $\alpha_t = \alpha$, for small enough α

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O(\eta^{t-\tau}) + O(\alpha)$$

where $\sigma < 1$ is **network connectivity**, $\eta < 1$ are **problem parameters**, and τ is the **mixing time** for the underlying Markov chain

- Time-varying step size $\alpha_t \sim 1/(t+1)$

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O\left(\frac{T}{(1-\sigma_2)^2} \frac{\log(t+1)}{t+1}\right)$$

Rate of convergence for distributed TD

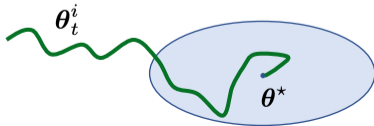
- Fixed **step size** $\alpha_t = \alpha$, for small enough α

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O(\eta^{t-\tau}) + O(\alpha)$$

where $\sigma < 1$ is **network connectivity**, $\eta < 1$ are **problem parameters**, and τ is the **mixing time** for the underlying Markov chain

- Time-varying step size $\alpha_t \sim 1/(t+1)$

$$\mathbb{E} [\|\boldsymbol{\theta}_t^i - \boldsymbol{\theta}^*\|] \leq O(\sigma^{t-\tau}) + O\left(\frac{T}{(1-\sigma_2)^2} \frac{\log(t+1)}{t+1}\right)$$



Goal: Find $\boldsymbol{\theta}^*$ such that $\bar{F}(\boldsymbol{\theta}^*) = \mathbf{0}$, where

$$\bar{F}(\boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}[F_i(X_i; \boldsymbol{\theta})],$$

using decentralized communications between agents with access to $F_i(X_i; \boldsymbol{\theta})$.

Using the iteration

$$\boldsymbol{\theta}_i^{k+1} = \sum_{j \in \mathcal{N}(i)} W_{i,j} \boldsymbol{\theta}_j^k + \epsilon F_i(X_i^k, \boldsymbol{\theta}_i^k)$$

gives us

$$\max_j \mathbb{E} \left[\|\boldsymbol{\theta}_i^k - \boldsymbol{\theta}^*\|_2^2 \right] \rightarrow O \left(\frac{\epsilon \log(1/\epsilon)}{1 - \sigma_2^2} \right) \quad \text{at a linear rate}$$

when the F_i are Lipschitz, \bar{F}_i are strongly monotone, and the $\{X_i^k\}$ are Markov

Multi-agent Reinforcement Learning, Scenario 2:

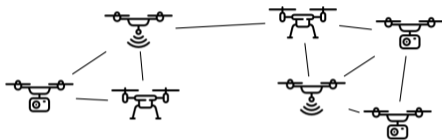
Multiple agents in different environments (dynamics, rewards)

Can we find a jointly optimal policy?

Policy Optimization, Framework

We will set this up as a distributed optimization program with decentralized communications

- One agent explores each environment
- Agent collaborate by sharing their models
- Performance guarantees:
number of gradient iterations
sample complexity (future)

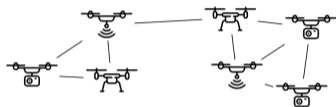


Policy Optimization, Framework

Environments $i = 1, \dots, N$, each with similar state/action spaces

Key quantities:

- $\pi(\cdot|s)$: policy that maps states into actions
- $r_i(s, a)$: reward function in environment i
- $\rho_i(s)$: initial state distribute in environment i
- $L_i(\pi)$: long-term reward of π in environment i



$$L_i(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_i(s_i^k, a_i^k) \right], \quad a_i^k \sim \pi(\cdot|s_i^{k-1}), \quad s_i^0 \sim \rho_i$$

We want to solve

$$\underset{\pi}{\text{maximize}} \quad \sum_{i=1}^N L_i(\pi)$$

$$\text{maximize}_{\pi} \sum_{i=1}^N L_i(\pi) \rightarrow \text{maximize}_{\theta} \sum_{i=1}^N L_i(\theta), \quad \pi_{\theta}(a|s) = \frac{e^{\theta_{s,a}}}{\sum_{a'} e^{\theta_{s,a'}}$$

- Natural parameterization (softmax) is **ill-conditioned at solution**

$$\underset{\pi}{\text{maximize}} \sum_{i=1}^N L_i(\pi) \rightarrow \underset{\theta}{\text{maximize}} \sum_{i=1}^N L_i(\theta) - \lambda \text{RE}(\theta), \quad \pi_{\theta}(a|s) = \frac{e^{\theta_{s,a}}}{\sum_{a'} e^{\theta_{s,a'}}}$$

- Natural parameterization (softmax) is **ill-conditioned at solution**

$$\underset{\pi}{\text{maximize}} \sum_{i=1}^N L_i(\pi) \rightarrow \underset{\theta}{\text{maximize}} \sum_{i=1}^N L_i(\theta) - \lambda \text{RE}(\theta), \quad \pi_{\theta}(a|s) = \frac{e^{\theta_{s,a}}}{\sum_{a'} e^{\theta_{s,a'}}}$$

- Natural parameterization (softmax) is **ill-conditioned at solution**
- Even for a single agent, this problem is nonconvex ...
... ability to find global optimum tied to “exploration conditions” (Agarwal et al '19)

$$\underset{\pi}{\text{maximize}} \sum_{i=1}^N L_i(\pi) \rightarrow \underset{\theta}{\text{maximize}} \sum_{i=1}^N L_i(\theta) - \lambda \text{RE}(\theta), \quad \pi_{\theta}(a|s) = \frac{e^{\theta_{s,a}}}{\sum_{a'} e^{\theta_{s,a'}}}$$

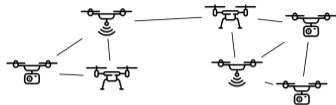
- Natural parameterization (softmax) is **ill-conditioned at solution**
- Even for a single agent, this problem is nonconvex ...
... ability to find global optimum tied to “exploration conditions” (Agarwal et al '19)
- Agents have competing interests (global solution suboptimal for every agent)

$$\underset{\pi}{\text{maximize}} \sum_{i=1}^N L_i(\pi) \rightarrow \underset{\theta}{\text{maximize}} \sum_{i=1}^N L_i(\theta) - \lambda \text{RE}(\theta), \quad \pi_{\theta}(a|s) = \frac{e^{\theta_{s,a}}}{\sum_{a'} e^{\theta_{s,a'}}}$$

- Natural parameterization (softmax) is **ill-conditioned at solution**
- Even for a single agent, this problem is nonconvex ...
... ability to find global optimum tied to “exploration conditions” (Agarwal et al '19)
- Agents have competing interests (global solution suboptimal for every agent)
- Gradients can only be computed imperfectly for large or partially specified problems

Algorithm: Decentralized Policy Optimization

$$\text{maximize}_{\{\theta_i\}} \sum_{i=1}^N L_i(\theta_i), \quad \text{subject to } \theta_i = \theta_j, (i, j) \in \mathcal{E}$$



- Each agent stores a local version of policy θ_i , initialized to θ_i^0
- At each node, iterate from policy $\pi_{\theta_i^k}$
 - ▶ Compute “advantage function” $A(s, a) = Q(s, a) - V(s)$
 - ▶ Compute gradient

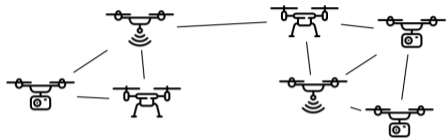
$$\nabla L_i(\theta_i^k) = (\text{complicated function of } \pi_{\theta_i^k} \text{ and } A(s, a))$$

- ▶ Meanwhile, exchange θ_i^k with neighbors
- ▶ Update policy

$$\theta_i^{k+1} = \sum_{j \in \mathcal{N}(i)} W_{i,j} \theta_j^k + \alpha_k \nabla L_i(\theta_i^k)$$

Algorithm: Mathematical Guarantees

$$\theta_i^{k+1} = \sum_{j \in \mathcal{N}(i)} W_{i,j} \theta_j^k + \alpha_k \nabla L_i(\theta_i^k)$$



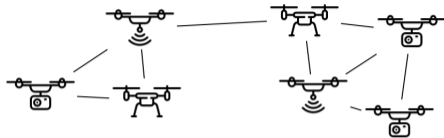
For small enough step sizes α_k , after k iterations we have

$$\left\| \frac{1}{N} \sum_{i=1}^N \nabla L_i(\theta_i^k) \right\|^2 \leq O\left(\frac{1}{\sqrt{k}} + \frac{C_g}{k}\right)$$

- Convergence to stationary point (not global max)
- Graph properties expressed in C_g
- Other constants come from λ , N , and MDP properties

Algorithm: Mathematical Guarantees

$$\theta_i^{k+1} = \sum_{j \in \mathcal{N}(i)} W_{i,j} \theta_j^k + \alpha_k \nabla L_i(\theta_i^k)$$

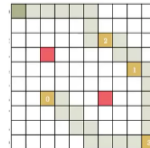
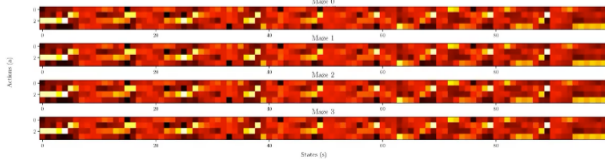
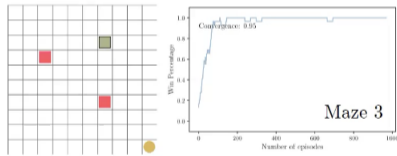
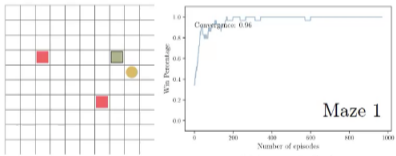
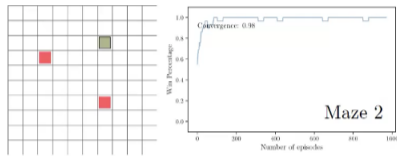
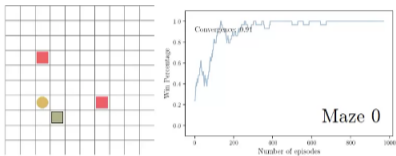


If common states are “equally explored” across environments, then after k iterations

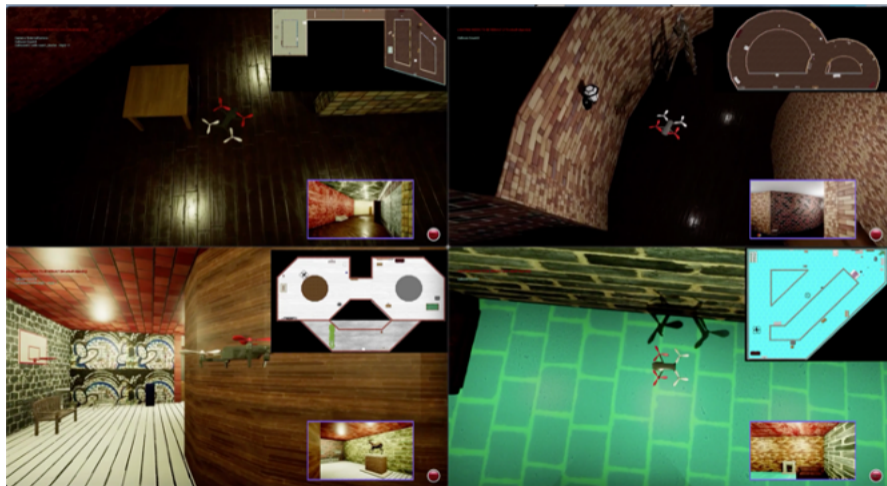
$$\max_j \left\{ \sum_{i=1}^N L_i(\theta^*) - L_i(\theta_j^k) \right\} \leq \epsilon \quad \text{when} \quad k \geq \frac{C}{\epsilon^2}$$

- Convergence to **global optimum**
- Requires careful choice of regularization parameter λ
- “Equal exploration” hard to verify
- Can make this stochastic, but not with finite-sample guarantee

MultiTask RL4 – Unconflicted Goals



Simulation: Drones in D-PEDRA



Simulation: Drones in D-PEDRA

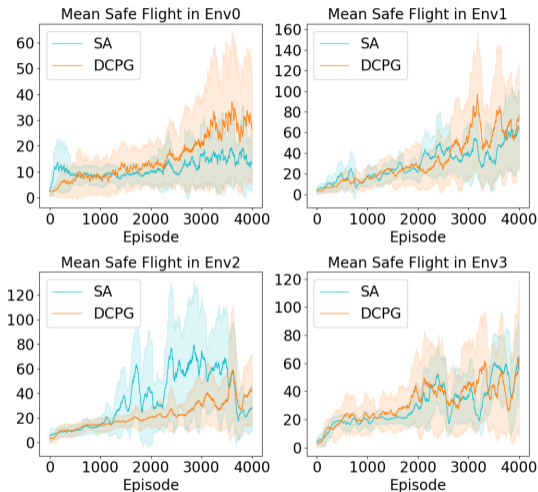


Table 1: MSF of the learned policy

Policy	Env0	Env1	Env2	Env3	Sum
SA-0	15.9	4.5	4.1	3.6	28.1
SA-1	3.0	55.4	9.7	8.1	76.2
SA-2	1.5	0.8	21.1	2.0	25.4
SA-3	2.3	0.8	8.6	40.1	51.8
DCPG	25.2	67.9	40.5	61.8	195.4
Random	2.5	3.9	4.7	3.7	14.8

Interesting, unexplained result:

Learning a joint policy is easier than learning individual policies

Thank you!

References:

T. T. Doan, S. M. Maguluri, and J. Romberg, “Finite-time performance of distributed temporal difference learning with linear function approximation,” submitted January 2020, arXiv:1907.12530.

T. T. Doan, S. M. Maguluri, and J. Romberg, “Finite-time analysis of distributed TD(0) with linear function approximation for multi-agent reinforcement learning,” ICML 2019.

S. Zeng, T. T. Doan, and J. Romberg, “Finite-time analysis of decentralized stochastic approximation with applications in multi-agent and multi-task learning,” submitted October 2020, arxiv:2010.15088 .

S. Zeng, et al, “ A decentralized policy gradient approach to multi-task reinforcement learning,” submitted October 2020, arxiv:2006.04338.