

1. Preliminaries

Bellman Optim Operator: $T_0(s, a) = \max_{s', r} \mathbb{E}_{\pi^0}[r + \gamma V(s')]$

1.1 Stochastic Processes

OTMC Recm: 1) The transition matrix (stochastic Matrix) satisfies: $\sum_j P_{ij} = 1$. 2) Let $X_0 = s_0$, then $X_i = \pi_{s_0, \dots, s_{i-1}} P$. 3) π is called the stationary dist. of a OTMC with transition Mat P : $\pi = \pi P$.

Then: OTMC said to be irreducible if all states communicate, e.g. $\forall i, j \exists T \in \mathbb{N}: P_{ij}(T) > 0$

Then: OTMC said to be aperiodic if $\forall s \in S$ $\text{GCD}(\{T \mid \pi_{ss}(T) > 0\}) = 1$

Then: OTMC said to be ergodic if $\forall s \in S$ $\lim_{T \rightarrow \infty} \mathbb{P}_{ss}(T) = 1$

Then: $\forall k \in \mathbb{N}$ $\mathbb{P}_{ss}(k) > 0$ if transition graph has a closed loop from state s & length $k \leq k \leq N$

Then: Any finite state OTMC has a stationary dist. π s.t. $\pi = \pi P$. If finite state OTMC is irreducible, then π is unique. If OTMC is ergodic, π converges to it. $\lim_{T \rightarrow \infty} \mathbb{P}_{ss}(T) = 1$

1.2 Optimization:

Linear Programming: Primal: $\min \mathbf{c}^\top \mathbf{x} \leq \mathbf{b}$, $\mathbf{A}\mathbf{x} \leq \mathbf{b}$. Dual: $\max \mathbf{b}^\top \mathbf{y} \leq \mathbf{c}^\top \mathbf{y} = \mathbf{c}$, $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$

Then: If one of them (primary/dual) has an optimal sol., then so does the other & their optimal values are same.

Concavity: $f(x) = \frac{1}{2} \mathbf{x}^\top \mathbf{A}(\mathbf{x}) + b^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{P} \mathbf{x} \geq 0$

Descent Methods: Goal: $\min f(x)$. Scheme: $x_{t+1} = x_t - \eta_t \nabla f(x_t)$.

Alg: $d_t = \nabla f(x_t)$, $\alpha_t = -\eta_t$. 1) GO: $d_t = -\nabla f(x_t)$, 2) Newton GO: $d_t = -\nabla f(x_t) - Q^{-1} d_t$, 3) Newton's Method: $d_t = -[\nabla^2 f(x_t)]^{-1} \nabla f(x_t)$

Stoprule: $\|x_t - x_{t-1}\| < \epsilon$, $f(x_t) - f(x_{t-1}) < \epsilon$.

2) Diminishing Stepsize: $y_t = \frac{1}{t} \cdot \eta_t$. 3) Exact Line search: $y_t = \arg\min_y \|f(x_t) - f(y)\|$.

Convergence & Complexity: Given $\frac{\partial^2 f}{\partial x^2} \leq L$, Accuracy Measure: $E(x) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [f(\mathbf{x})] - g \rightarrow 0$ 1) Lin. Conv.: $\|\mathbf{x} - g\| \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\mathbf{x} - g\|] \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\mathbf{x} - g\|^2]^{1/2}$

2) Sublinear Conv.: $\|\mathbf{x} - g\| \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\mathbf{x} - g\|] + \frac{L}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\mathbf{x}\|^2]$

3) Superlinear conv.: $\|\mathbf{x} - g\| \leq O(\epsilon^{1/2})$

4) Conv of order p : $\|\mathbf{x} - g\| \leq O(\epsilon^{1/p})$, when $p=2$ quadratic conv.

e.g. $C(x, \epsilon) = O(\epsilon^{1/p})$

Then (Conv of GD): $x_{t+1} \leftarrow x_t - \eta_t \nabla f(x_t)$

Suppose: f is cont. diffible with L -Lipschitz gradients. Let $y_t = \frac{1}{t} \sum_{i=1}^t x_i$

1) general: $\|\mathbf{x}_t - \mathbf{y}_t\| \leq \frac{1}{t} \|\mathbf{x}_t - \mathbf{y}_t\| \leq \frac{1}{t} \|\mathbf{x}_t - \mathbf{y}_t\|$

2) If f is convex: $\|\mathbf{x}_t - \mathbf{y}_t\| \leq \frac{1}{t} \|\mathbf{x}_t - \mathbf{y}_t\|$

3) Stronger convex: If f is strongly convex ($f''(x) \geq \lambda \|\mathbf{x}\|^2$)

↳ $\|\mathbf{x}_t - \mathbf{y}_t\| \leq (\frac{1}{t} - \frac{1}{t^2}) \|\mathbf{x}_t - \mathbf{y}_t\|$

Stochastic Gradient Descent: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$

Unbiased estm: $\mathbf{x}_t = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}_t]$, $\mathbf{x}_t - \mathbf{y}_t = \mathbf{x}_t - \mathbf{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}_t]$

SGD: $\mathbf{x}_t \leftarrow \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$, $\mathbf{x}_t - \mathbf{y}_t = \text{unbiased estm}$

LR Rule: SGD is not a monotone descent method.

Then: Assume: $\eta_t = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\eta_t] \approx 1$ smooth ($\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\|\mathbf{x}_t - \mathbf{y}_t\|^2] \leq \epsilon$)

$\Rightarrow \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}_t - \mathbf{y}_t] = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}_t] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{y}_t] = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}_t] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}_t] = 0$

Then: $\mathbf{x}_t - \mathbf{y}_t = \eta_t \nabla f(\mathbf{x}_t)$

Relaxed FG-Alg: $\mathbf{x}_t \leftarrow (\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)) / (1 + \eta_t)$

Sublin. Conv: $\|\mathbf{x}_t - \mathbf{y}_t\| \leq \frac{\eta_t}{1 + \eta_t} \|\nabla f(\mathbf{x}_t)\|$

2. Dynamic & Linear Programming

2.1 MDPs:

Def: MDP = controlled MC + performance obj. - An MDP consists of S : State space, A : Action space, $P(s, a, s')$: $S \rightarrow A \rightarrow S'$ trans. model, $r(s, a)$: $S \times A \rightarrow R$ reward func., $\pi: S \times A \rightarrow \Delta A$: init. state dist., $\pi_0: S \rightarrow A$

Def: Controlled Markov Property: $P(s, a, s') = \mathbb{P}(s' | s, a)$

Performance Criteria: Finite Horizon: $\text{Cum. Rew. } E[\sum_{t=0}^{T-1} r_t]$

Average Reward: $E[\frac{1}{T} \sum_{t=0}^{T-1} r_t]$, $R \rightarrow 0$ \Rightarrow immediate reward

∞ -horizon: $y(t) = \text{Discounted Rew. } E[\sum_{t=0}^{\infty} \gamma^t r_t]$, $\gamma \in [0, 1]$

Q-Learning: $Q(s, a) = \text{Cum. Rew. } E[\sum_{t=0}^{\infty} \gamma^t r_t | s, a]$

State Value Fct: $V(s) = E[\sum_{t=0}^{\infty} \gamma^t r_t | s]$

State-Action Value Fct: $Q(s, a) = E[\sum_{t=0}^{\infty} \gamma^t r_t | s, a]$

Rem: $V^*(s) = \max_a Q(s, a)$, $Q^*(s, a) = \max_a Q(s, a)$

Def: $V^*(s) = \max_a V^*(s, a)$, $Q^*(s, a) = \max_a Q^*(s, a)$

Optimal policy: In dual formulation: 1) solve dual LP to obtain $\pi^* \Rightarrow 2) \Pi^*(s, a) = \frac{1}{Z(s)} \sum_{a \in A} \pi^*(s, a) \exp(Q^*(s, a))$

Def: $V^*(s) = \max_a V^*(s, a)$, $Q^*(s, a) = \max_a Q^*(s, a)$

Def: $V^*(s) = \max_a V^*(s, a)$, $Q^*(s, a) = \max_a Q^*(s, a)$

3. Value-Based Methods

3.1 From Planning to RL: learn π, Q, V

Value-Based: Learn $V^*(s)$, Policy-Based: Learn π , Model-Based: Function-based: $V(s), Q(s, a), \Pi(s, a) \Rightarrow \mathbb{P}(s'|s, a)$

Reinforcement Learning e.g. Linear: $V(s) = (V(s), \dots, V(s))$

Contrastive Properties: Bellman Optm, Opt. $\Rightarrow \mathbb{P}^*$ s.t. $\mathbb{P}^* \circ \mathbb{P} = \mathbb{P}^*$

Contraction Mapping: Bellman Optm, Opt. $\Rightarrow \mathbb{P}^*$ is a fixed point of \mathbb{P}^*

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \mathbb{E}_{\pi(s)}[\mathbb{P}^*(s', a)] = \mathbb{E}_{\pi(s)}[\mathbb{P}(s', a, \pi(s', a))]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \mathbb{E}_{\pi(s)}[\mathbb{P}^*(s', a)] = \mathbb{E}_{\pi(s)}[\mathbb{P}(s', a, \pi(s', a))]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

Gaussian Process: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.2 Bellman Eq's & Bellman Optimality

Bellman Consistency Eq: $V^*(s) = \mathbb{E}_{\pi(s)}[V^*(s', a)] = \mathbb{E}_{\pi(s)}[\mathbb{P}(s', a, \pi(s', a))]$

Bellman Optimality: $V^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Bellman Optimality Operator: $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[\mathbb{P}(s', a, \pi(s', a))]$

Bellman Expectation Operator: $\mathbb{P}^*: \mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \mathbb{E}_{\pi(s)}[\mathbb{P}^*(s', a)] = \mathbb{E}_{\pi(s)}[\mathbb{P}(s', a, \pi(s', a))]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.3 Model-free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.4 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.5 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.6 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.7 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.8 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.9 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.10 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.11 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.12 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.13 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.14 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.15 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.16 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.17 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$, $P = \mathbb{P}^{\text{trans}}$

Lemma: $(I - \mathbb{P})^{-1} = \mathbb{P}^*$, always invertible for $\gamma < 1$

then: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^* \circ \mathbb{P}^* \circ \dots \circ \mathbb{P}^*$

Double Learning (OCU): Goal: lower the maxim. bias.

General Case: $T_0(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a)) = \exp(-\frac{(s - \mathbb{P}(s))^2}{2\sigma^2})$

Then: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Double Learning: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.18 Value Fit Approx: If $|S| \ll |A|$ large \Rightarrow model approx.

Goal: $\mathbb{P}(s, a) = \mathbb{P}(s, a | \mathbb{P}(s, a))$

Learn π, Q, V : learn π, Q, V

3.19 Model-Free Prediction

Goal: $\mathbb{P}^* = \mathbb{P}^* \circ \mathbb{P}^*$, \mathbb{P}^* is a fixed point of \mathbb{P}^*

Given policy π : $\mathbb{P}^*(s) = \max_{\pi(s)} \mathbb{E}_{\pi(s)}[V^*(s', a)]$

Matrix Form: $\mathbb{P}^* = B + \gamma P$

Global Conv. of Projected Policy Grad Method:

Access occurs to exact gradient & def $\eta = \frac{(\gamma\pi)^2}{1-\gamma}$. Then:
 $J(\pi^*) - \max J(\pi) \geq \frac{\gamma^2(1-\gamma)}{1-\gamma} \frac{1}{\| \nabla \pi \|^2} \gg \text{large}$
 slow rate
 Proof sketch:
 Show that the objective $J(\pi)$ is smooth with L -smoothness and $J(\pi) \leq \frac{1}{L}$.
 Invoke convergence on gradient mapping: $\min_{\pi \in \Pi} \|G(\pi)\|_2 \leq 1 + L\|\pi\|_2$.
 Use the relationship between gradient mapping and approximation of stationary point [4]:
 $\max_{\pi \in \Pi} \langle \nabla J(\pi), \nabla J(\pi_1) \rangle = 1 + L\|\pi_1\|_2, \quad \|\pi_1 - \pi_0\|_2$.

Use the gradient dominance for global convergence.

Rem: PG update can be viewed as:

$$T_{t+1} = \text{Proj}_{\Pi} (T_t + \eta \nabla J(\pi_t)) - \text{argmax}_{\pi} \{ \langle \nabla J(\pi_t), \cdot \rangle \}_{\Pi} T_{t+1}$$

$$= \text{argmax}_{\pi} \left\{ \langle Q^*(\pi), \pi \rangle - \frac{1}{1-\gamma} \sum_i \pi_i^2 \right\}$$

where: $\pi_i(\pi) = \frac{1}{1-\gamma} \pi_i^*(\pi) C^*(\pi)$ & $\langle \cdot, \cdot \rangle$ is renormalized inner prod.

Policy Mirror Descent:
 $T_{t+1} = \text{argmax}_{\pi} \{ \langle Q^*(\pi), \pi \rangle - \frac{1}{1-\gamma} \sum_i \pi_i^2 \}_{\Pi} D(T_t, \pi_t, \pi_{t+1})$

Projected O-Accent: set $D(\pi) = K(\pi) \log \pi + \frac{1}{2} \pi \log \pi$

Problem: not guaranteed to converge

Good: Converges for zero-sum games.

Then for 2P ZSG the emp. dist. of FP conv. to NE:

$(\pi_1, \pi_2) \rightarrow (\pi_1^*, \pi_2^*)$ where (π_1^*, π_2^*) is NE. Conv. rate: $O(\frac{1}{\sqrt{T}})$

Grad. Accen: $T_{t+1}(\pi) = T_t(\pi) + \alpha \nabla_{\pi} D(\pi)$ & project to Π

Rem: $\alpha = \text{argmin}_{\alpha} \{ \langle \nabla_{\pi} D(\pi), \nabla_{\pi} D(\pi) \rangle \}_{\Pi} + \frac{1}{2}$

No Atm to NPG, so $\alpha \rightarrow 1$ reduces to prior iteration update.

4.8 Global Conv. of NPG Methods:

Let $T_t(\pi) = \exp(B(\pi)/2\alpha \log(\pi))$ & $\pi_t = T_t(\pi)$

$\Rightarrow \theta_t = \theta_t + \frac{1}{T} \nabla_{\theta}^T A^T$, $\pi_{t+1}(\pi) = \pi_t(\pi) + \frac{\exp(\theta_t^T \pi_t)}{Z(\theta_t)}$

Convergence of Uniform MFG: $V_t = V_{t+1}$ (log(π_t) & π_t)

$J(\pi^*) - J(\pi_t) \leq \frac{2}{T}$ Rem: Opt-free conv. No dep on θ_t !

No dep on state mismatch coeff

Lemma: $J(\pi) - J(\pi_t) \leq \frac{1}{T} \mathbb{E} \{ [KL(\pi^*(\cdot) || \pi_t(\cdot))] \}_{\Pi} + \frac{1}{T} \sum_{i=1}^T \mathbb{E}_{\pi_t(\cdot)} [\log Z_i(\cdot)]$

Proof: Performance difference lemma + $\log(2/\pi)$

Proof of Uniform MFG: Conv:

Setting $\pi = \pi_t$ in the previous lemma and telescoping from $t = 0, \dots, T-1$

$$\frac{1}{T} \sum_{i=0}^{T-1} (J(\pi^*) - J(\pi_i)) \leq \frac{1}{T} \mathbb{E}_{\pi_t(\cdot)} [KL(\pi^*(\cdot) || \pi_t(\cdot))] + \frac{1}{T} \sum_{i=0}^{T-1} \mathbb{E}_{\pi_t(\cdot)} [\log Z_i(\cdot)]$$

Setting $\pi = \pi_{T-1}$, in the previous lemma, we have

$$J(\pi_{T-1}) - J(\pi_t) \geq \frac{1}{T} \mathbb{E}_{\pi_t(\cdot)} [\log Z_i(\cdot)] - \mathbb{E}_{\pi_{T-1}(\cdot)} [\log Z_i(\cdot)] > 0, \forall i. (\text{why?})$$

Combining these two equations and the fact that $J(\pi) \geq \frac{1}{T}$ implies that

$$\frac{1}{T} \sum_{i=0}^{T-1} (J(\pi^*) - J(\pi_i)) \leq \log |\mathcal{A}| + \frac{1}{(1-\gamma)^2 T}.$$

Ref: $T_{t+1}(\pi) = \frac{\exp(\theta_t^T \pi_t)}{Z(\theta_t^T \pi_t)}$, $\pi_t = T_t(\pi)$,

$\theta_t = G + \frac{1}{T} \nabla_{\theta}^T A^T$, with $G = \text{argmin}_{\theta} \{ \langle \nabla_{\theta} D(\theta), \nabla_{\theta} D(\theta) \rangle \}_{\Pi} + \frac{1}{2} \mathbb{E}_{\pi_t(\cdot)} [\log Z(\cdot)]$

$\Rightarrow T_{t+1}(\pi) = T_t(\pi) + \frac{\exp(\theta_t^T \pi_t)}{Z(\theta_t^T \pi_t)} \exp(\frac{1}{2} \mathbb{E}_{\pi_t(\cdot)} [\log Z(\cdot)])$

Conv. of sample based MFG: $E[\min J(\pi_{T-1}) - J(\pi_0)] \leq$

$O(\frac{1}{T} \sqrt{2 \log(1/\epsilon)} + \text{Var}(V) + \sqrt{\epsilon \text{Var}(V)})$ where

Error: how close π_t to a w.r.t. statistical error

has good the best policy in the class

5. Multi-Agent RL & Markov Games:

5.1 From Single Agent to Multiple Agents:

Issues: exp. growing search space, Non-stationarity of env., Non-Markovian

Competing

Markov Games: Multi-agent & full state info

5.2 Normal Form Games & Repeated Games:

Normal Form Games: players/agents I, II, ..., jointed.

$I = (a_1, \dots, a_n)$, Reward $R(a, \dots)$ = NPG (Van der game)

If repeated too many times: repeated game

Strategies: $\pi_I, \pi_E, \Delta(\pi_I) = \pi_I(a)$, prob agent i selects a pure strategy (deterministic policy): only play one action mixed strategy (stochastic): a dist over set of actions

Strategy Profile: concord of players' strategies $\pi = (\pi_I, \pi_E)$

Expected Payoff of player i : $J_i(\pi) = E_{\pi_I}[\pi_I(a_i)]$

$\Rightarrow \pi_I(\cdot, \cdot) \pi_I(\cdot, \cdot) \pi_I(\cdot, \cdot)$ (Indep. policies)

Zero-Sum Games: In 2-player game: $\pi_I(a_1, a_2) = \pi_I(a_1) \pi_I(a_2)$

Payoff can be represented with matrix A : $A_{ij} = \pi_I(a_i) \pi_I(a_j)$

Expected Payoff of P : $\pi_I(\pi_I, \pi_I) = \pi_I(\pi_I)$

Response Models: Fix π_I , $\pi_I(\cdot, \cdot, \cdot, \cdot)$ Best response: $\pi_I(\cdot, \cdot, \cdot, \cdot) = \pi_I(\cdot, \cdot, \cdot, \cdot)$ vs. $\pi_I(\cdot, \cdot, \cdot, \cdot)$ vs. $\pi_I(\cdot, \cdot, \cdot, \cdot)$

Sol. linear response: $\pi_I(\cdot, \cdot) \propto \exp(J_I(\cdot, \cdot, \cdot, \cdot))$

Dominant Strategy Equilibrium: π^* : optimal π_I

It may not exist!

Nash Equilibrium: π^* : no player can improve exp payoff by changing policy if the others stick to their policy.

$\Rightarrow J_I(\pi^*) \geq J_I(\pi_I, \pi_I)$ $\forall i$

\Rightarrow π_I is mixed strategy, $\pi_I \in \text{argmax}_{\pi_I} J_I(\pi_I, \pi_I)$

Then in a normal form game with finite players/actions, $\exists \text{NE}$

Policy Gradient Methods / Gradient Ascent:

Max in 2P ZSG: $\max_{\pi \in \Pi} \langle \nabla J(\pi), \cdot \rangle_{\Pi} = \min_{\pi \in \Pi} \langle \nabla J(\pi), \pi \rangle_{\Pi}$
 \Leftrightarrow saddle point of $\text{Proj}_{\Pi} f(\pi)$: $f(\pi) = \langle \nabla J(\pi), \pi \rangle_{\Pi}$

Minimax Then $X \in \mathbb{R}^d$, $Y \in \mathbb{R}^d$, $f: X \times Y \rightarrow \mathbb{R}$ cont. s.t. $f(x, \cdot)$ is convex $\forall y$ & $f(\cdot, y)$ is concave $\forall x$. Then $\min_x \max_y f(x, y) = \max_y \min_x f(x, y)$

Iterated Best Response Alg: $\pi_I \leftarrow \cdot$ do each

player i updates $T_i(\cdot, \cdot)$ $\pi_I \leftarrow T_i(\cdot, \pi_I)$ $\forall i$

Non-Convergence: ICR^* (ICR) $\neq \text{NE}$ (ACR) $\neq \text{PDR}$ (PDR)

PG Alg: $K^* = K^* - \frac{\eta}{\eta - 1} \frac{\partial}{\partial K^*} \langle \nabla J(K^*), K^* \rangle_{\Pi}$

Then \exists LG Game that the set of int cord in a neighborhood of the Nash eq. from which gradient converges to the Nash eq. is of measure zero

Nash Q-Learning: V stage t do: 1. Agent i takes action a_t^* & observes every player's reward, every player's action & the next state s_{t+1} . 2. Compute the value of $\forall V$ players: $J(V) = \sum_i \mathbb{E}_{a_t^* \sim \pi_t(a_t^*)} [r_i + \gamma V(s_{t+1})]$

MLE guarantees: $E_{\pi_t} [V(s_{t+1})] = \log(J(V))$

Fairness Alg.: best bus option π^* : max $\{V, V^*\}$ $\forall V$

$\Rightarrow K^* = \max \{V, V^*\}$ $\forall V$

Can learn policy from demonstrations solving the following saddle point problem max π : $E_{\pi} [\log(J(\pi))]$

$\Rightarrow \min_{\pi} \max_{\pi'} \langle \nabla_{\pi} \log(\pi), \langle \nabla_{\pi'} \log(\pi'), \pi \rangle_{\Pi} \rangle$

Conv. of Gradient Descent: Jensen's inequality

Corollary: Let Π be a discrete & re-usable policy class, i.e. the Π with prob. $\pi \in \Pi$, behavioral cloning based on MLE returns a policy that obeys: $J(\Pi) - \langle \nabla J(\Pi), \pi \rangle_{\Pi} \geq \frac{1}{1-\gamma} \langle \nabla \log(J(\Pi)), \pi \rangle_{\Pi}$

MLE guarantees: $E_{\pi_t} [V(s_{t+1})] = \log(J(V))$

Fairness Alg.: best bus option π^* : max $\{V, V^*\}$ $\forall V$

$\Rightarrow K^* = \max \{V, V^*\}$ $\forall V$

Can learn policy from demonstrations solving the following saddle point problem max π : $E_{\pi} [\log(J(\pi))]$

$\Rightarrow \min_{\pi} \max_{\pi'} \langle \nabla_{\pi} \log(\pi), \langle \nabla_{\pi'} \log(\pi'), \pi \rangle_{\Pi} \rangle$

5.4 Two Player Zero Sum Markov Games:

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$

$\pi_1, \pi_2, \pi_3, \pi_4$ s.t. <