# Lecture 20

CS 59000_RL

MDP

- Infinite horizon MDP - Undiscounted
- Optimism

we defined

$$\rho^\Pi = \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} P_\Pi^{t-1} r_\Pi$$

$$\to \text{Let } \overline{P}_\Pi = \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} P_\Pi^{t-1}$$

$$\Rightarrow \rho^\Pi = \overline{P}_\Pi r_\Pi$$

$$V_\Pi^T = \sum_{t=1}^{T} P_\Pi^{t-1} \left( r_\Pi \xrightarrow{\downarrow} \rho^\Pi \right) \xrightarrow{\overline{\rho}^\Pi} \mathbb{1}$$

$$\Rightarrow V_\Pi = \lim_{T \to \infty} \frac{1}{T} \sum_{T > 0} V_\Pi^T$$

Hardy 1945

Lemma: $V_\Pi = \left( (I - P_\Pi - \overline{P}_\Pi)^{-1} - \overline{P}_\Pi \right) r_\Pi$

A value function $V: \mathcal{X} \to \mathbb{R}$

$$\Rightarrow span(V) = \max_{\eta \in \mathcal{X}} V(\eta) - \min_{\eta \in \mathcal{X}} V(\eta)$$

Lemma: For any memoryless policy $\pi$

$$\bar{P}_\pi V_\pi = 0$$

(why ?)

Proof: note that $\boxed{\bar{P}_\pi P_\pi = \bar{P}_\pi}$

$$\Rightarrow \bar{P}_\pi V_\pi^T = \sum_{t=1}^{T} \bar{P}_\pi P_\pi^{t-1} \left( r_\pi - \rho^\pi \right) = T \left( \bar{P}_\pi \left( r_\pi - \rho^\pi \right) \right)$$

$\underbrace{\qquad}_{\bar{P}_\pi}$

$$= T \left( \underbrace{\bar{P}_\pi r_\pi}_{\rho^\pi} - \rho^\pi \right)$$

$\boxed{\bar{P}^\pi 1}$

$$= T \left( 0 \right) = 0$$

$$\Rightarrow V_\pi^T = 0 \Rightarrow V_\pi = 0$$

---

Lemma: For any memoryless policy, we have

$$\rho^\pi + V_\pi = r_\pi + P_\pi V_\Pi$$

Historically this equation is known as Poisson equation
Nowdays we call it Bellman equation.

---

proof: $\quad V_\pi = (I - P_\pi + \bar{P}_\pi)^{-1} r_\pi - \bar{P}_\pi r_\pi$

$\rho^\pi$

$$V_\pi + \rho^\pi = (I - P_\pi + \bar{P}_\pi)^{-1} r_\pi$$

$\underbrace{\qquad\qquad}$

$$\left( I - P_\Pi + \overline{P}_\Pi \right) \left( V_\Pi + \rho^\Pi \right) = \gamma_\Pi$$

$$V_\Pi - P_\Pi V_\Pi + \overline{P}_\Pi V_\Pi + \underbrace{\rho^\Pi - P_\Pi \rho^\Pi}_{\smash{\overbrace{\text{}}}} + \underbrace{\overline{P}_\Pi \rho^\Pi}_{\rho^\Pi} = \gamma_\Pi$$

$$\underset{0}{\underbrace{}}$$

$$\Rightarrow \quad V_\Pi - P_\Pi V_\Pi + \rho^\Pi = \gamma_\Pi$$

---

## Bellman Optimality equation          inner product.

$$\longrightarrow \quad \rho + V(x) = \max_{a \in \mathcal{A}} \left( \overline{r}(a,x) + \left\langle P(\cdot | x, a), V \right\rangle \right)$$

Bellman operator: $T(V) = \max_a \left( \overline{r}(a,x) + \left\langle P(\cdot | x, a), V \right\rangle \right)$

For $V$, define greedy policy $\Pi_V(x) \in \arg\max_{a \in \mathcal{A}} \overline{r}(a,x) + \left\langle P(\cdot | x, a), V \right\rangle$

## Theorem:
i) Bellman Optimality equation has a solution
ii) A solution $(\rho, V)$ satisfies $\rho = \rho^*$, $\Pi^* = \Pi_V$

---

Why not unique? if $V^\Pi$ is solution $\Rightarrow V^\Pi + \alpha 1$ is also a solution.

---

Now, how to solve it?
   - Policy iteration
   - Value iteration.
   - Linear program.

minimize $\rho$
$\rho \in R, V \in R^{|x|}$

$$\text{s.t.: } \rho + V(n) \geq \bar{r}(n, a) + \langle P(\cdot | n, a), V \rangle \quad \forall (n, a) \in X \times A$$

$\Rightarrow$ The solution is $\rho^*$. But, the $V$ and $\Pi_V$ might be useless.

having $\rho^*$, now so for $V^*$:

minimize $\langle V, 1 \rangle$
$V \in R^{|x|}$

s.t., $\begin{cases} \rho^* + V(n) \geq \bar{r}(n, a) + \langle P(\cdot | n, a), V \rangle \quad \forall (n, a) \in X \times A \\ \\ V(n_0) = 0 \qquad\qquad n_0 \Rightarrow \text{a state that all } \Pi^* \text{ visib.} \end{cases}$

Now consider, we neither have $P$, nor $R$.

We need to intract with the environment, enplore it, learn it, and exploit what we have.
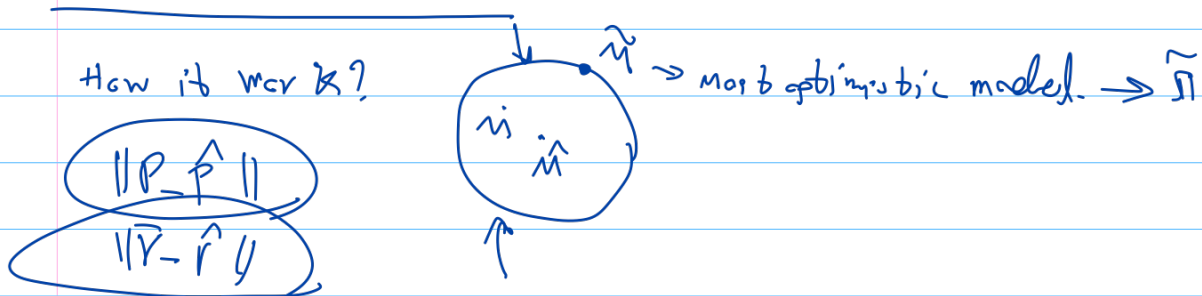
$\hookrightarrow$ Tradeoff exploration and exploitation.

Define regret as: $\hat{R}_T = T\rho^* - \sum_{t=1}^{I} r(X_t, A_t)$

Find an algorithm with reasonable upper bound on regret.

We use optimism

### Alg: Upper Confidence bound for reinforcement learning $\underline{2}$
### $(UCRL2)$

How it work?

$\boxed{\|P - \hat{P}\|}$

$\boxed{\|\bar{r} - \hat{r}\|}$

$\tilde{M} \Rightarrow$ Most optimistic model. $\Rightarrow \tilde{\pi}$

True model $M$, $P, \bar{r}$
Estimated model $\hat{M}: \hat{P}, \hat{\bar{r}}$

For now: $0 \leq r \leq 1$, we assume $\bar{r} = r$ and we know $\bar{r}$

At time $t$, define $T_t(x, a) = \sum_{i=1}^{t} \mathbb{I}\{X_i = x, A_i = a\}$

Then $P_t(x' \mid x, a) = \dfrac{\sum_{i=1}^{t} \mathbb{I}\{X_i = x, A_i = a, X_{i+1} = x'\}}{1 \vee \underset{\smile}{T_t(x, a)}}$

$\Rightarrow$ i.e empirical estimate.

Define Confidence set for $x, a$

$\Rightarrow C_t^{\delta}(x, a) = \left\{ P' \in \Delta^{|X|} : \| P' - \underset{t-1}{P}(\cdot \mid x, a) \|_1 \leq \sqrt{\dfrac{|X| L_{t-1}(x, a)}{1 \vee T_{t-1}(x, a)}} \right\}$

For $L_{t-1}(x, a) = 2 \log \left( \dfrac{4 |X| |A| T_t(x, a)(1 + T_t(x, a))}{\delta} \right)$

$$P_{t-1} \in R^{|X|} \qquad \sum_{i=1}^{|X|} P_{t-1}(i) = 1 \Rightarrow P_{t-1} \in \Delta^{|X|}$$

$\Delta^3$

## UCRL2:

For $k = 1, \cdots$

- Set $\tau_k = t+1$
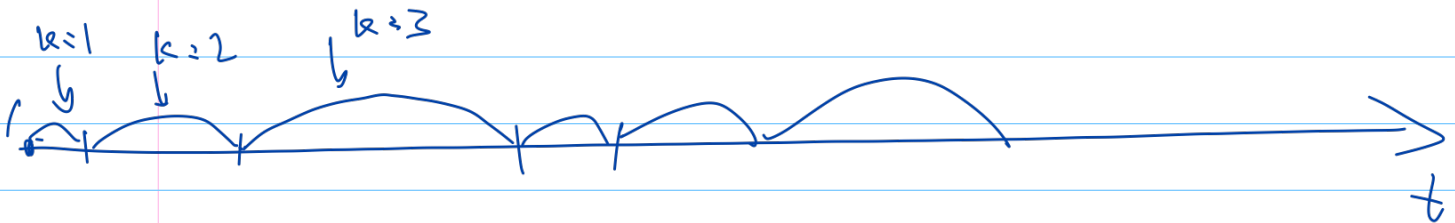- Find an optisbic model, return, value its optimal memory less, determisbic policy

$$\Rightarrow \rho_k = \max_{\pi \in X} \left( \max_{\pi} \max_{P' \in C_{P_k}} \rho_n^{\Pi}(P') \right)$$

- Optimistic model $P_k$, $\Pi_k \rightarrow$ optimistic policy

- els

 $t \leftarrow t+1$, observe $X_t$, take $A_t = \Pi_k(X_t)$

while $T_t(X_t, A_t) < 2 T_{\tau_{k-1}}(X_t, A_t)$

$k=1 \quad k=2 \qquad k=3$

$t$

Theorem: UCRL2 acheives regret of

$$\hat{R}_T < c \, D \, |\mathcal{X}| \sqrt{|\mathcal{A}| \, T} \, \log\left(\frac{T \, |\mathcal{X}| \, |\mathcal{A}|}{\delta}\right)$$

universal constant

with probability at least $1-\delta$

Lower bound on MDP

For $|\mathcal{X}| \geq 3$, $|\mathcal{A}| \geq 2$, $D \geq 6 + 2 \log_{|\mathcal{A}|} |\mathcal{X}|$,

and $T \geq D \, |\mathcal{X}| \, |\mathcal{A}|$,

then for any algorithm, there exists an MDP, such that

$$E[\hat{R}_T] \geq c' \sqrt{D \, |\mathcal{X}| \, |\mathcal{A}| \, T}$$

universal constant.

UCRL2 is order optimal; But off $\sqrt{|\mathcal{X}|} D$

up to log factors.