

Final presentation schedule in the third sheet of sign-up google sheet.

Piazza poll @ 13.

Final report:

* Introduction

- what is the research problem?
- why is this fun/interesting
- prior work. w/ citations
- what ~~are~~ the goals your project tries to answer?
- what you've accomplished. (at a high level).

* Preliminaries: (notation & setting)

- define models / notations
- basic math facts. Lemma / Fact.

* Body

- details.

§ Reinforcement Learning

account for long term effects of actions.

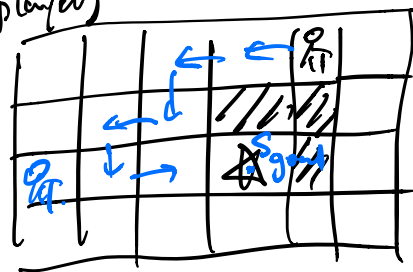
Mathematical model: episodic Markov Decision Processes (MDPs).

Defns: state space: \mathcal{S} (e.g. location of player)

action space: A .

(left / right / up / down / stay)

episode length H : (# steps in episode) $\{ \dots H \}$



reward fn $r_1, \dots, r_H : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

(e.g. $r_h(s_{goal}, a) = 1$)

$r_h(s, a) = 0 \quad s \neq s_{goal}$

transition fn: $P_1, \dots, P_H : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$

Learning protocol:

Assume: r_1, \dots, r_H known.

P_1, \dots, P_H unknown

(e.g. know the goal's location)

For episodes $t=1, 2, \dots, T$:

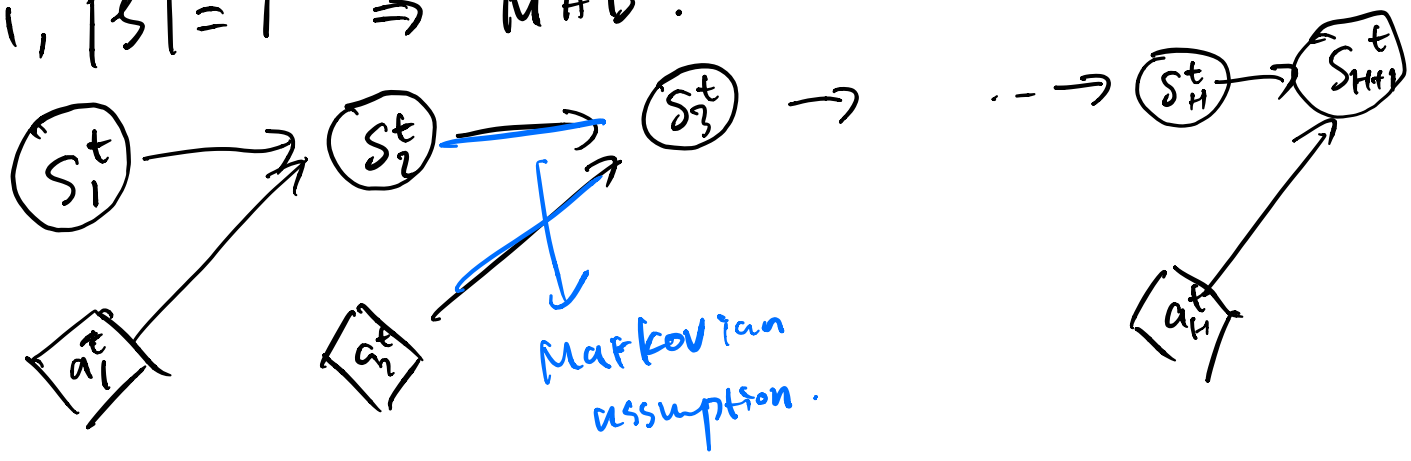
- learner sees some initial state for this episode s_1^t

For steps $h=1, \dots, H$:

- learner takes action $a_h^t \in \mathcal{A}$, receive reward $r_h(s_h^t, a_h^t)$, observe next state $s_{h+1}^t \sim P_h(\cdot | s_h^t, a_h^t)$.

$H=1 \Rightarrow$ contextual bandits

$H=1, |\mathcal{S}|=1 \Rightarrow$ MAB.

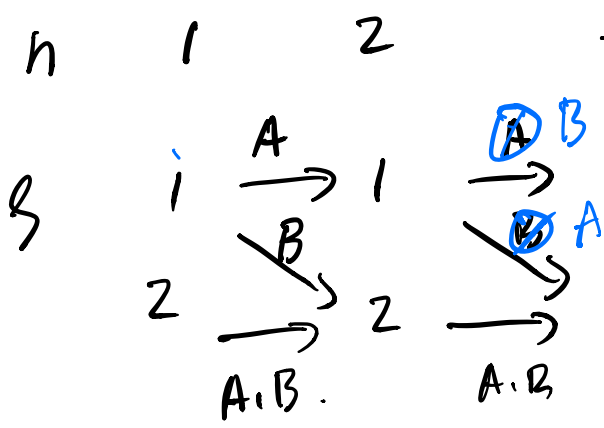


What makes RL uniquely challenging?

Ex: combination lock.

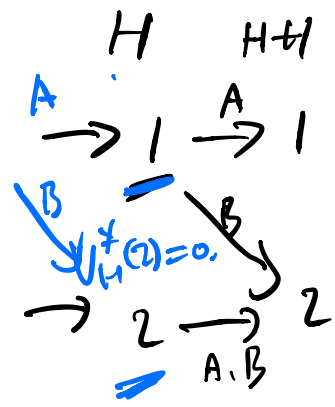
$\mathcal{L} = \{1, 2\}$ $\mathcal{A} = \{A, B\}$ H B

$(P_h)_{h=1}^H$



$V_H^*(1) = 1 \iff Q_H^*(1, A) = 1$

$Q_H^*(1, A) = 1$
 $Q_H^*(1, B) = 0$



$Q_H^*(2, A) = 0$
 $Q_H^*(2, B) = 0$

$\forall h, s, a,$

$r_h(s, a) = 0$. except that $r_H(1, A) = 1$
 $r_H(1, B) = 1$

① If random actions are taken, chance of receiving nonzero reward is $\frac{1}{2^{H-1}}$ (exponentially small)

② "A...A" as password.

If we view this problem as a MAB problem (w/ 2^{H-1} actions) then need to try 2^{H-1} "actions" to get nonzero reward.

Optimal value & policy in MDPs.

$\sum_{h=1}^H r_h$

$V_h^*(s)$ = highest expected reward we can get, given that we are at state s at step h .

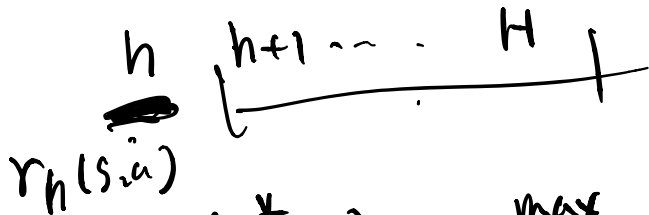
$Q_h^*(s, a)$ = highest expected reward we can get, given that we are at state s and take action a at step h .

Solve V_n^* Q_n^* 's using dynamic programming.

$h = H+1$. Q_{H+1}^* undefined $V_{H+1}^*(s) = 0$.

$h \leq H$. $Q_h^*(s, a) = r_h(s, a) + \sum_{s' \in \mathcal{S}_h} P_h(s'|s, a) \cdot V_{h+1}^*(s')$

Bellman optimality eqn.



$V_h^*(s) = \max_{a \in \mathcal{A}} Q_h^*(s, a)$

$V_1^* \leq \dots \leq V_H^* \leq Q_H^* \leq V_{H+1}^*$

Optimal policy: given we are at state s at h .

$a = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q_h^*(s, a) =: \pi_h^*(s)$

$\pi^* = (\pi_1^* \dots \pi_H^*)$ In general, π_h^* 's are

not equal.

$\pi = (\pi_1, \dots, \pi_h, \dots, \pi_H)$

$V_h^\pi(s)$ = expected reward given we are at state

s at step h , and we execute policy π .

Solve V_h^π using DP:

$V_{h+1}^\pi(s')$

$$\underline{V}_h^{\pi}(s) = r_h(s, \pi_h(s)) + \sum_{s' \in \mathcal{S}} p_h(s' | s, \pi_h(s)) \underline{V}_{h+1}^{\pi}(s')$$

$$V_{H+1}^{\pi}(s) = 0.$$

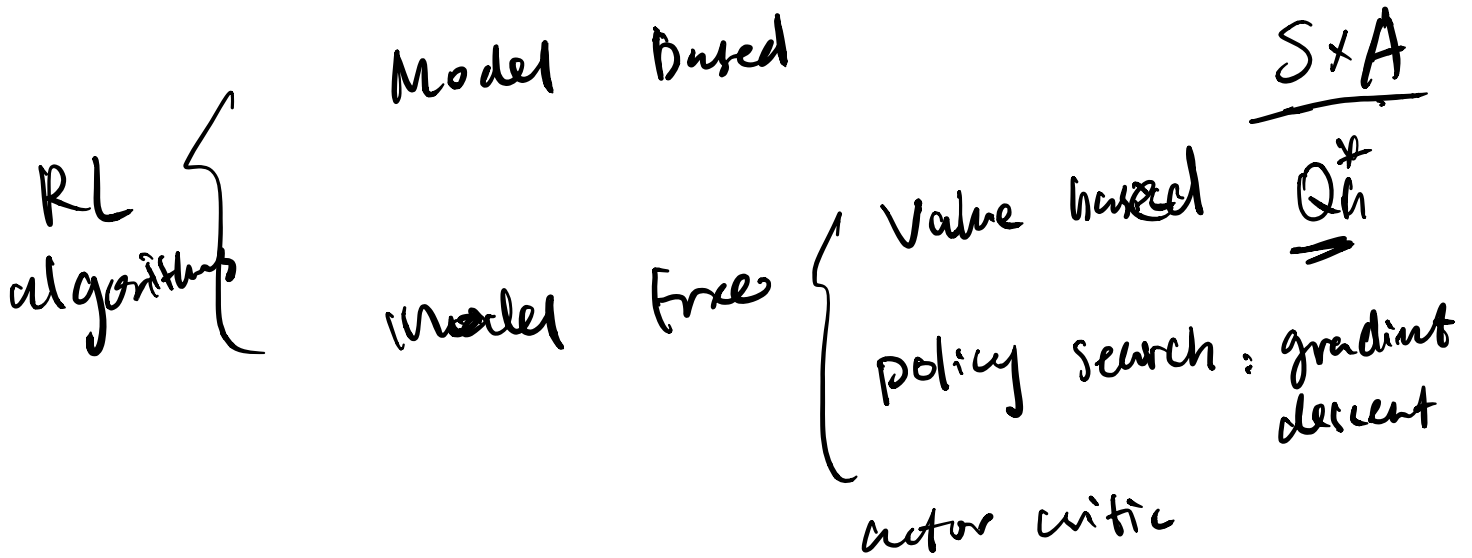
performance measure of online RL: $V_1^{\pi^t}(s_1^t)$

$$R_T = \mathbb{E} \left[\sum_{t=1}^T V_1^{\pi^t}(s_1^t) - \sum_{t=1}^T \sum_{h=1}^H r_h(s_h^t, a_h^t) \right]$$

learner takes policy π^t for episode t

$$= \mathbb{E} \left[\sum_{t=1}^T \left(V_1^{\pi^t}(s_1^t) - V_1^{\pi^t}(s_1^t) \right) \right].$$

would like $R_T = o(T)$.



This lecture: value-based method.

optimistic Q-learning (Jin, Allen-Zhu, Bubeck, Jordan, 2018).

maintain optimistic estimate \underline{Q}_n for Q_n^* online.

Alg: Initialization: $V_h(s) = H$. $\forall h, s, a$
 $Q_n(s, a) = H$. $\forall h, s, a$

$$V_{H+1}(s) \equiv 0.$$

$$m_n(s, a) = 0 \quad \forall s, a, h.$$

For episodes $t = 1, 2, \dots, T$:

execute $\pi_h^t(s) = \arg \max_a Q_n(s, a)$.
 initial state s_1^t
 For step $h = 1 \dots H$:

take action $a_h^t = \arg \max_{a \in A} Q_n(s_h^t, a)$.

transition to s_{h+1}^t . $r_h(s_h^t, a_h^t)$

$$m = m_n(s_h^t, a_h^t) \leftarrow m(s_h^t, a_h^t) + 1$$

update:

$$Q_n(s_h^t, a_h^t) \leftarrow (1 - \alpha_m) Q_n(s_h^t, a_h^t) + \alpha_m (r_h(s_h^t, a_h^t) + V_{H+1}(s_{h+1}^t) + b_m)$$

$$\alpha_m = \frac{1}{m}$$

$$\alpha_m = \frac{H+1}{H+m} > \frac{1}{m}$$

$$V_h(s_h^t) \leftarrow \min_{a \in A} \max_{a \in A} Q_n(s_h^t, a)$$

$$Q_n^*(s, a) = r_h(s, a) + \langle P_h(\cdot | s, a), V_{h+1}^* \rangle$$

~~$$Q_n^*(s, a) = r_h(s, a) + \langle P_h(\cdot | s, a), V_{h+1}^* \rangle$$~~

\mathbb{E}

$$V_h^*(s_h^t)$$

$$C \leq \sqrt{H^3 \ln \frac{SAHT}{\delta}}$$

$$r_n(s_{n+1}^a) + V_{n+1}(s_{n+1}^a)$$

$$r_n(s_n^a) + \overset{\uparrow}{V_{n+1}(s_{n+1}^t)}$$

Thm: optimistic Q-learning has:

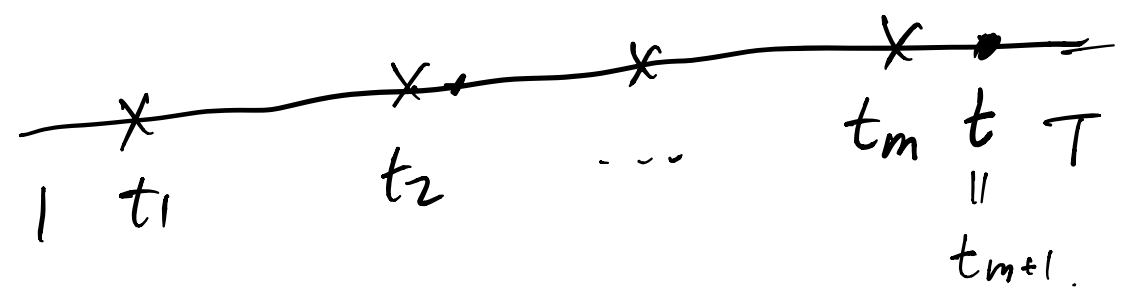
$$R_T \leq O(\sqrt{T \cdot S \cdot A \cdot H^5})$$

(Lower bound: \forall alg, \exists episodic RL environment, $R_T \geq \Omega(\sqrt{T \cdot S \cdot A \cdot H^3})$)

defs: $Q_n^t(s, a) = Q_n(s, a)$, at the beginning of episode t

$$V_n^t(s) = V_n(s) \dots \dots \dots$$

$$m_n^t(s, a) = m_n(s, a) \dots \dots \dots$$

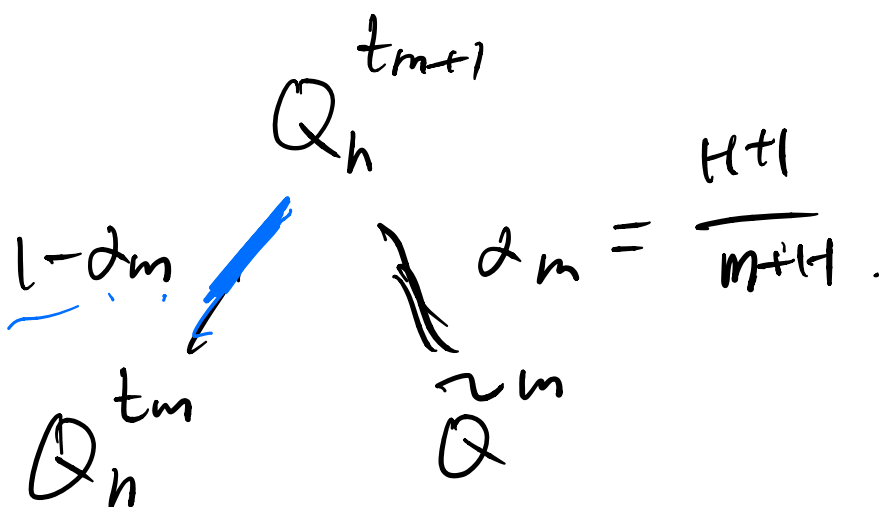


$$Q_n^{t_1}(s, a) = H$$

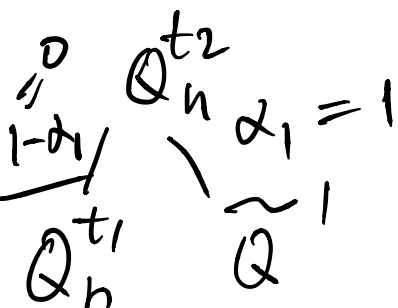
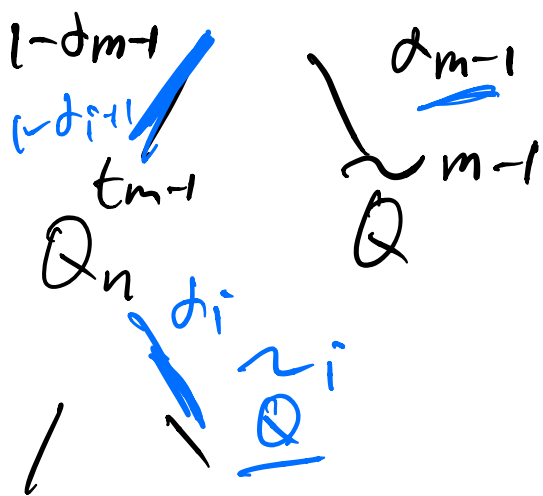
$$\forall i = 1, 2, \dots, m.$$

$$Q_n^{t_{i+1}}(s, a) = (1 - \alpha_i) Q_n^{t_i}(s, a) + \alpha_i \cdot \left(r_n(s, a) + V_{n+1}^{t_i}(s_{n+1}) + b_i \right)$$

\tilde{Q}^i



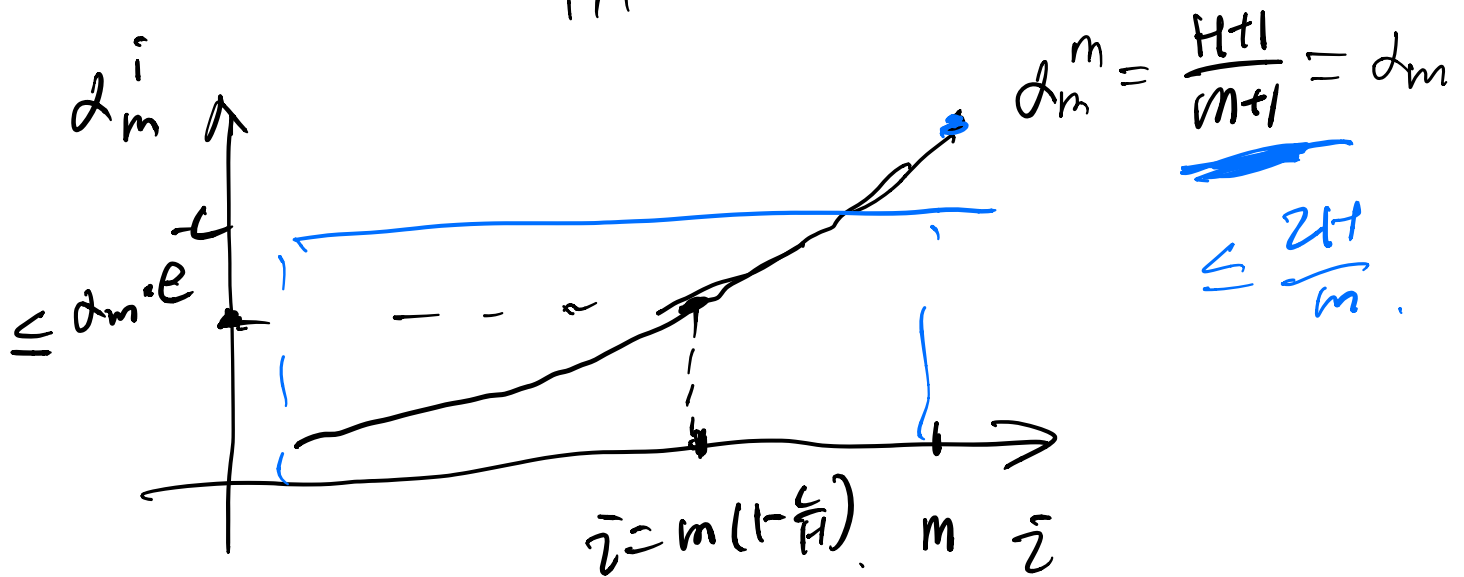
$$\frac{1}{m-1} \times \left(1 - \frac{1}{m}\right) = \frac{1}{m}$$



$$Q_n^t(s, u) = \sum_{i=1}^m d_m^i z_i^t \quad m \geq 1.$$

$$d_m^i = d_i \cdot \prod_{j=i+1}^m (1 - d_j)$$

For $d_i = \frac{H+1}{H+i}$ d_m^i ?



properties of d_m^i :

① $\sum_{i=1}^m d_m^i = 1 \quad d_m^i \geq 0 \quad \forall i$

② $\sum_{i=1}^m d_m^i \cdot \frac{1}{\sqrt{i}} \in \left[\frac{1}{\sqrt{m}}, \frac{2}{\sqrt{m}} \right]$

③ $\sum_{i=1}^m (d_m^i)^2 \leq \frac{2H}{m}$

$\leq \left(\sum_{i=1}^m d_m^i \right) \cdot \left(\max d_m^i \right)$

$$\textcircled{4} \quad \sum_{m=i}^{\infty} \alpha_m^i = 1 + \frac{1}{H}$$

$$m \leq i \left(1 + \frac{1}{H}\right) \cdot \sum_{m=i}^{i(1+\frac{1}{H})} \frac{H}{m}$$

$$Q_n^t(s, u) = \sum_{i=1}^n \alpha_m^i (r_n(s, u) + V_{n+1}^{t_i}(s_{n+1}^{t_i}) + b_i) + I(m=0) \cdot H$$

Lemma: \exists event E , $P(E) \geq 1 - \delta$, $\forall t, \forall h$

$\forall s, a$

$$0 \leq \textcircled{1} \quad Q_n^t(s, a) - Q_n^*(s, u)$$

$$\leq \textcircled{2} \quad I(m=0) \cdot H + 3 \cdot \frac{C}{\sqrt{m}} I(m > 0)$$

$$= m_n^t(s, a)$$

$$+ \sum_{i=1}^m \alpha_m^i \phi_{n+1}^{t_i}$$

unique to RL.

Notation

$$\left[\phi_n^t = V_n^t(s_n^t) - V_n^*(s_n^t) \right]$$

Pf: Recall:

$$Q_h^t(s, a) = \sum_{i=1}^m \alpha_m^i (r_h(s, a) + V_{h+1}^{ti}(s_{h+1}^{ti}) + b_i) + I(m=0) \cdot H$$

$$Q_h^x(s, a) = r_h(s, a) + \langle P(\cdot | s, a), V_{h+1}^x \rangle$$

$$= \sum_{i=1}^m \alpha_m^i (\downarrow) + I(m=0) Q_h^x(s, a)$$

$$Q_h^t(s, a) - Q_h^x(s, a) = \sum_{i=1}^m \alpha_m^i (\underbrace{V_{h+1}^{ti}(s_{h+1}^{ti})}_{(a)} - \underbrace{\langle P(\cdot | s, a), V_{h+1}^x \rangle}_{(b)})$$

$$+ \sum_{i=1}^m \alpha_m^i b_i + I(m=0) (\underbrace{H}_{(d)} - \underbrace{Q_h^x(s, a)}_{\substack{\leq H \\ \geq 0}})$$

$$(a) = \sum_{i=1}^m \alpha_m^i \phi_{h+1}^{ti} \in [0, H] \quad \in [0, H]$$

$$\textcircled{b} = \sum_{i=1}^m d_m^i \left(V_{n+1}^* (S_{h+1}^{t_i}) - \langle P_h \cdot |S_n\rangle, V_{n+1}^* \right) \in [-H, H]$$

Azuma's inequality: Martingale seq. $X_1 \sim X_n$.

$$|X_i| \leq B_i$$

$$\left| \sum_{i=1}^n X_i \right| \leq \sqrt{\sum_{i=1}^n B_i^2 \cdot \ln \frac{1}{\delta}}$$

w.p. $1-\delta$.

$$B_i = H \cdot d_m^i$$

$$\leq \sqrt{H^2 \cdot \sum_{i=1}^m (d_m^i)^2 \ln \frac{1}{\delta}}$$

union bound

$$\frac{H}{m}$$

$$\textcircled{b} \in \left[-\frac{c}{\sqrt{m}}, \frac{c}{\sqrt{m}} \right]$$

$$\textcircled{c} \in \left[\frac{c}{\sqrt{m}}, \frac{2c}{\sqrt{m}} \right]$$

$$\textcircled{d} \in [0, I(m=0) \cdot H]$$

$$Q_h^t(s, a) - Q_h^*(s, a)$$

$$\geq \sum_{i=1}^m d_m^i \left(V_{h+1}^{t_i}(s_{h+1}^{t_i}) - V_{h+1}^*(s_{h+1}^{t_i}) \right)$$

~~$$+ \left(-\frac{c}{\sqrt{m}} \right) + \frac{c}{\sqrt{m}} + 0$$~~

$$V_{h+1} \geq V_{h+1}^*$$

" " " "

Pf of inequality ①:

backward induction on h .

base case: $h = H, V_{H+1} \equiv V_{H+1}^* \equiv 0$.

$$Q_H^t \geq Q_H^*$$

inductive case: suppose for $h+1, \forall t$.

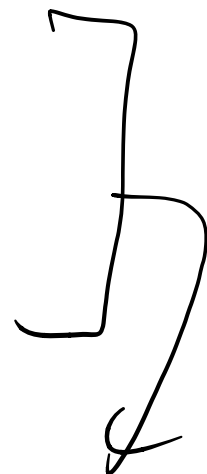
$$Q_{h+1}^t(s, a) \geq Q_{h+1}^*(s, a)$$

max over a on both

$$\max_a Q_{h+1}^t(s, a) \geq \underbrace{V_{h+1}^*(s)}$$

$$H \geq V_{h+1}^*(s)$$

$$V_{h+1}^t(s) \geq V_{h+1}^*(s)$$



$$\Rightarrow Q_n^t \geq Q_n \quad \text{.} \quad \text{.} \quad \text{.}$$

pf of ②

$$Q_n^t(s,a) - Q_n^*(s,a) \leq I(m=0)H +$$

$$\frac{3\epsilon}{\sqrt{m}} I(m>0) + \sum_{i=1}^m \alpha_m^i \phi_{n+1}^{t_i}$$

pf of thm:

$$\sum_{t=1}^T \left(V_1^*(s_1^t) - V_1^{\pi^t}(s_1^t) \right)$$

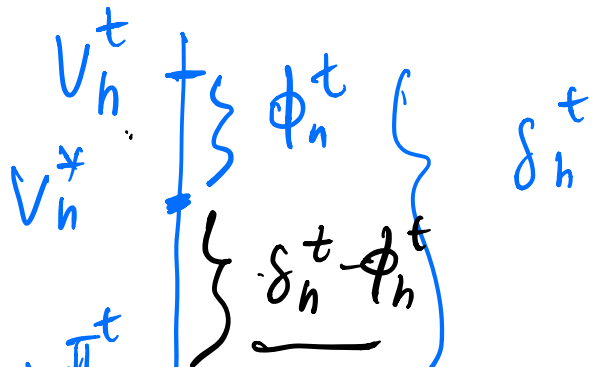
Lemma 1

$$\leq \sum_{t=1}^T \left(\underbrace{V_1^t(s_1^t) - V_1^{\pi^t}(s_1^t)}_{\min(H, Q_n^t(s_n^t, a_n^t))} \right)$$

define

$$\delta_n^t = \underbrace{V_n^t(s_n^t) - V_n^{\pi^t}(s_n^t)}$$

$$\text{Recall } \phi_n^t = V_n^*(s_n^t) - V_n^*(s_n^t)$$



Bounding

$$\sum_{t=1}^T \delta_1^t \quad ?$$

$$\sum_{t=1}^T \delta_{H+1}^t \equiv 0.$$

rekte

$$\sum_{t=1}^T \delta_n^t$$

to

$$\sum_{t=1}^T \delta_{n+1}^t$$

$$\delta_n^t \leq Q_n^t (S_n^t \cdot a_n^t) - Q_n^{\pi^t} (S_n^t \cdot a_n^t)$$

$$= \underbrace{(Q_n^t - Q_n^y)}_{(a_2)} (S_n^t \cdot a_n^t) + \underbrace{(Q_n^y - Q_n^{\pi^t})}_{(b_t)} (S_n^t \cdot a_n^t)$$

$$\sum_{t=1}^T (a_t) \leq \sum_{t=1}^T \underbrace{I(m_h^t = 0)}_{m_h^t (S_n^t \cdot a_n^t)} \cdot H \quad (a_1)$$

$$+ \sum_{t=1}^T \frac{3c}{\sqrt{m_h^t}} I(m_h^t > 0) \quad (a_2)$$

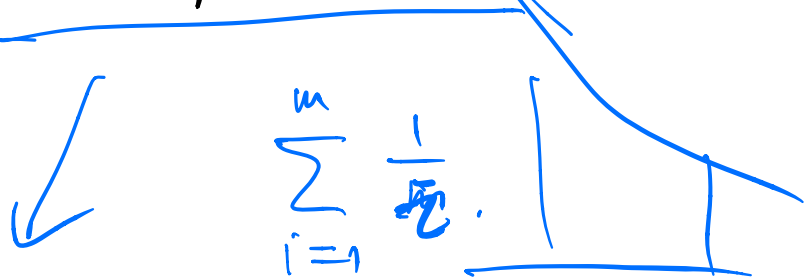
$$+ \sum_{t=1}^T \sum_{\tilde{j}=1}^{m_h^t} \alpha_{m_h^t}^i \cdot \phi_{n+1}^{ti} \quad (a_3)$$

$$\alpha_{1,1} \leq H \times S \times A \times H.$$

$$\alpha_{1,2} \leq \sum_{S.A} \sum_{t=(S_h^t \cdot a_h^t) = (S.A)} \frac{3c}{\sqrt{m_h^t}} I(m_h^t > 0)$$

$m_h^{T+1}(S.A)$ terms

$$= \sum_{S.A} 3c \cdot \sum_{m=1}^{m_h^{T+1}(S.A)} \sqrt{\frac{1}{m}}$$



$$\leq 2 \sqrt{m_h^{T+1}(S.A)}$$

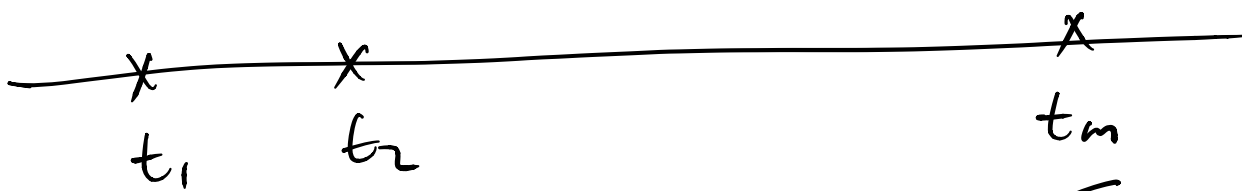
Cauchy-Schwarz

$$\leq 6c \cdot \sum_{S.A} \sqrt{m_h^{T+1}(S.A)}$$

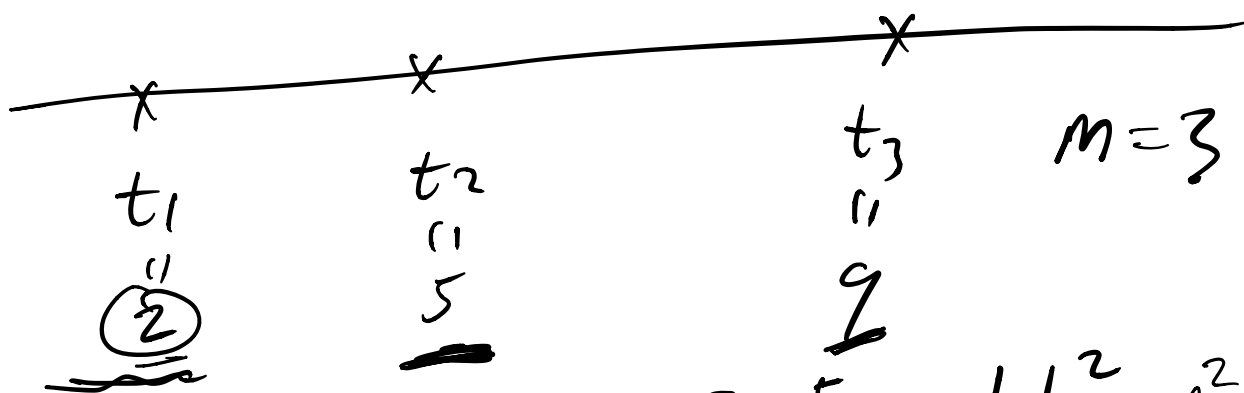
$$\leq 6c \cdot \sqrt{\frac{\sum_{S.A} 1}{S \cdot A} \cdot \frac{\sum_{S.A} m_h^{T+1}(S.A)}{T}}$$

(a.3) : regroup according to

$$(S_n^t \cdot a_n^t) = \underline{\underline{(S \cdot a)}}$$



$$m = m_h^{T+1} (S \cdot a)$$



$$d_1^1 \cdot \phi_{n+1}^2$$

$$d_2^1 \phi_{n+1}^2 + d_2^2 \phi_{n+1}^5$$

$$d_3^1 \phi_{n+1}^2 + d_3^2 \phi_{n+1}^5$$

$$+ d_3^3 \phi_{n+1}^9$$

$$m_h^t (S \cdot a)$$

$$\sum_{S \cdot a} \sum_{t_i (S_n^t \cdot a_n^t) = (S \cdot a)}$$

$$\sum_{i=1}^m d_m^i \phi_{n+1}^i$$

coeff of ϕ_{n+1}^2 : $d_1^1 + d_2^1 + d_3^1 \leq 1 + \frac{1}{H}$

ϕ_{n+1}^5 : $d_2^2 + d_3^2 \leq 1 + \frac{1}{H}$

$$\phi_{n+1}^q : \alpha_3^3 \leq 1 + \frac{1}{H}$$

$$\leq \left(1 + \frac{1}{H}\right) \cdot \sum_{t=1}^T \phi_{n+1}^t$$

$$\sum_t \textcircled{a_t} \leq H^2 SA + b \sqrt{SAT} + \left(1 + \frac{1}{H}\right) \sum_{t=1}^T \phi_{n+1}^t$$

$$\sum_t \textcircled{b_t} = \sum_t Q_h^t(S_h^t, a_h^t) - Q_h^{\pi_t}(S_h^t, a_h^t)$$

$$= \sum_t \left\langle P_h(\cdot | S_h^t, a_h^t), V_{h+1}^* - V_{h+1}^{\pi_t} \right\rangle$$

$\in [-H, H]$ $\leq H \cdot \sqrt{T}$

$$\leq \sum_t \left\langle P_h(\cdot | S_h^t, a_h^t) - P_{S_{h+1}^t}, V_{h+1}^* - V_{h+1}^{\pi_t} \right\rangle$$

$$+ \underbrace{V_{h+1}^*(S_{h+1}^t) - V_{h+1}^{\pi_t}(S_{h+1}^t)}$$

$$\leq \underbrace{\left(1 + \frac{1}{H}\right)} \cdot \underbrace{(\delta_{h+1}^t - \phi_{n+1}^t)}$$

combining:

$$\sum \delta_h^t \leq \left(1 + \frac{1}{H}\right) \sum \delta_{h+1}^t + H^2 SA$$

$$+ \underbrace{6c\sqrt{SAT}}_{\sqrt{H^3}} + \underbrace{H\sqrt{T}}.$$

$$\sum_t \delta_t \approx \underbrace{\left(1 + \frac{1}{H}\right)^H}_{\approx \sqrt{H^3}} \sum_t \delta_{H+1}$$

$$+ \sum_{h=1}^H \underbrace{\left(1 + \frac{1}{H}\right)^{H-h}}_{\leq H} \left(H^2 SA + \underbrace{c\sqrt{SAT}}_{\approx \sqrt{H^3}} \right)$$

$$\left(1 + \frac{1}{H}\right)^H \leq e.$$

$$\leq O\left(\underbrace{H^3 SA} + \sqrt{H^5 SAT} \right)$$

$$\Rightarrow \text{expected regret} \leq O\left(\sqrt{H^5 SAT}\right).$$

~~✗~~