

# **Exploration in Tabular MDPs:**

## **Upper Confidence Bound Value Iteration (UCBVI)**

## **Announcements**

1. Scribing Lecture Notes (see Piazza for details)

2. Course Project Website  
<https://wensun.github.io/CS6789projects.html>

## Recap:

### Generative Model and Statistical Limits

Theorem: (Azar et al. '13) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \gamma \sqrt{\frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N}} + \frac{c\gamma}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N},$$

where  $c$  is an absolute constant.

Corollary: for  $\epsilon < 1$ , provided  $N \geq \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$  then

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \epsilon \text{ (with prob. greater than } 1 - \delta)$$

A.s.a

$N$  i.i.d samples

$\widehat{P}(s'|s,a)$

plan inside  $\widehat{P}$ .

## Recap:

### Generative Model and Statistical Limits

**Theorem:** (Azar et al. '13) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \gamma \sqrt{\frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N}} + \frac{c\gamma}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N},$$

where  $c$  is an absolute constant.

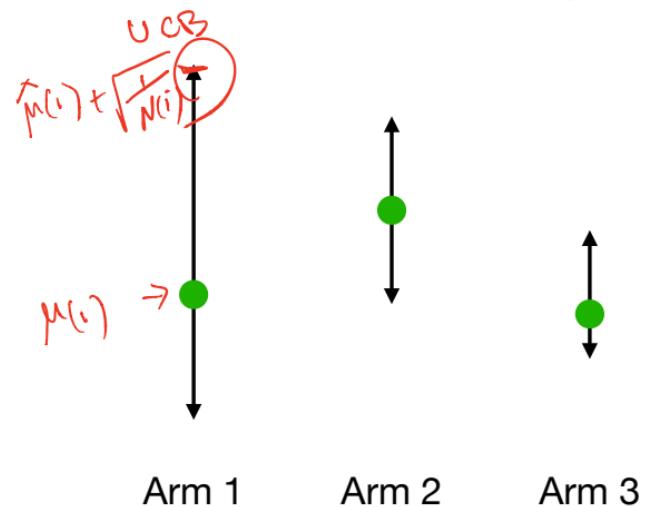
**Corollary:** for  $\epsilon < 1$ , provided  $N \geq \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$  then

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \epsilon \text{ (with prob. greater than } 1 - \delta)$$

Generative Model: we can reset to anywhere we want (may not be realistic)

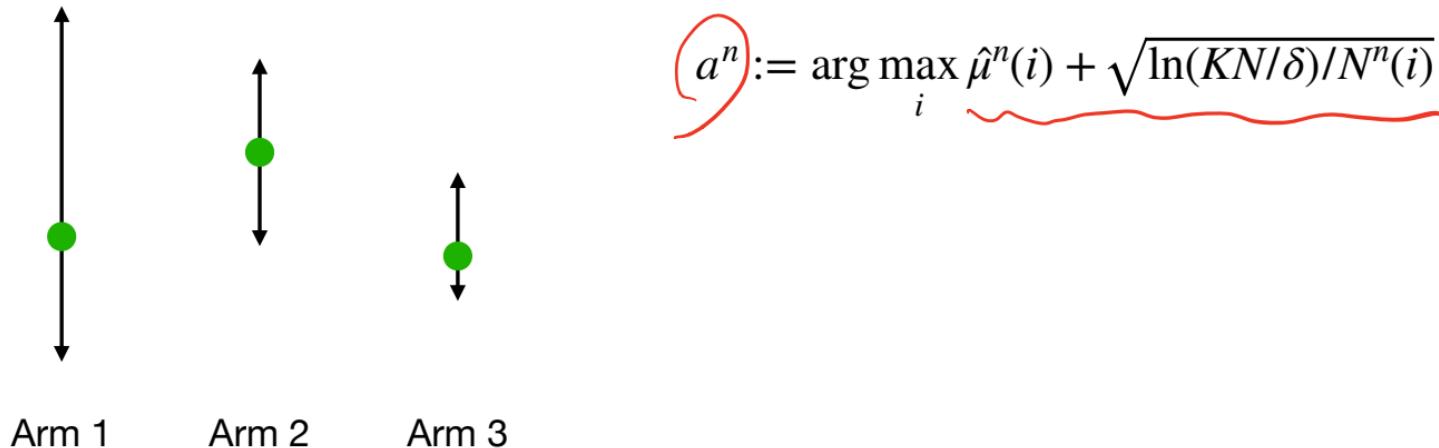
## Recap:

### Multi-armed Bandits and UCB Algorithm



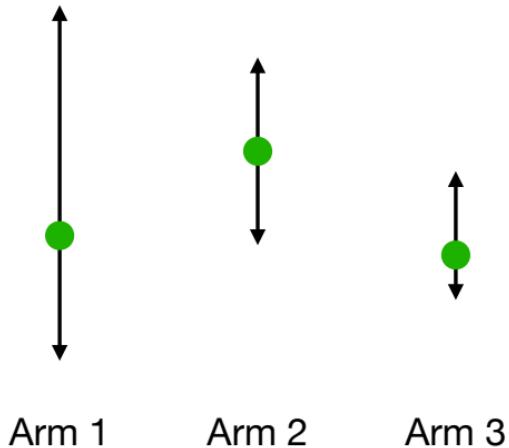
## Recap:

### Multi-armed Bandits and UCB Algorithm



## Recap:

### Multi-armed Bandits and UCB Algorithm

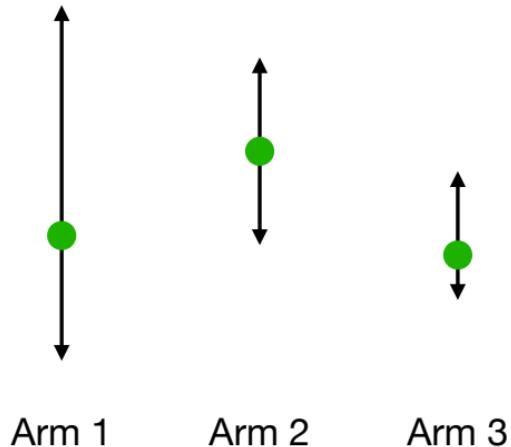


$$a^n := \arg \max_i \hat{\mu}^n(i) + \sqrt{\ln(KN/\delta)/N^n(i)}$$

$$\mathbb{E} \left[ N\mu(a^\star) - \sum_{n=1}^N \mu(a^n) \right] \leq \widetilde{O}(\sqrt{KN})$$

## Recap:

### Multi-armed Bandits and UCB Algorithm



$$a^n := \arg \max_i \hat{\mu}^n(i) + \sqrt{\ln(KN/\delta)/N^n(i)}$$

$$\mathbb{E} \left[ N\mu(a^\star) - \sum_{n=1}^N \mu(a^n) \right] \leq \widetilde{O}(\sqrt{KN})$$

Key step in the proof:  $\cancel{\text{Def}} \quad UCB(a^n) \geq UCB(a^\star)$   
 $\geq \mu(a^\star)$

$$\mu(a^\star) - \mu(a^n) \leq \hat{\mu}(a^n) + \sqrt{\frac{\ln(KN/\delta)}{N^n(a_n)}} - \mu(a^n)$$

$UCB(a^n)$

## Today: Efficient Learning in Finite Horizon tabular MDPs

Finite horizon episode (time-dependent) discrete MDP  $\mathcal{M} = \left\{ \{r_h\}_{h=0}^{H-1}, \{P_h\}_{h=0}^H, H, \mu, S, A \right\}$

## Today: Efficient Learning in Finite Horizon tabular MDPs

Finite horizon episode (time-dependent) discrete MDP  $\mathcal{M} = \left\{ \underbrace{\{r_h\}_{h=0}^{H-1}}_A, \underbrace{\{P_h\}_{h=0}^H}_B, H, \mu, S, A \right\}$

Only reset from  $\mu$ : we assume it's a delta distribution, all mass at a fixed  $s_0$   $\mu(s_0)=1$   
 $\mu(s_1)=0$

Unknown Transition  $P$  (for simplicity assume reward is known)

$s \neq s_0$

## Today: Efficient Learning in Finite Horizon tabular MDPs

Finite horizon episode (time-dependent) discrete MDP  $\mathcal{M} = \left\{ \{r_h\}_{h=0}^{H-1}, \{P_h\}_{h=0}^H, H, \mu, S, A \right\}$

Only reset from  $\mu$ : we assume it's a delta distribution, all mass at a fixed  $s_0$

Unknown Transition  $P$  (for simplicity assume reward is known)

Different from the Generative Model Setting!

## Today: Efficient Learning in Finite Horizon tabular MDPs

Finite horizon episode (time-dependent) discrete MDP  $\mathcal{M} = \left\{ \{r_h\}_{h=0}^{H-1}, \{P_h\}_{h=0}^H, H, \mu, S, A \right\}$

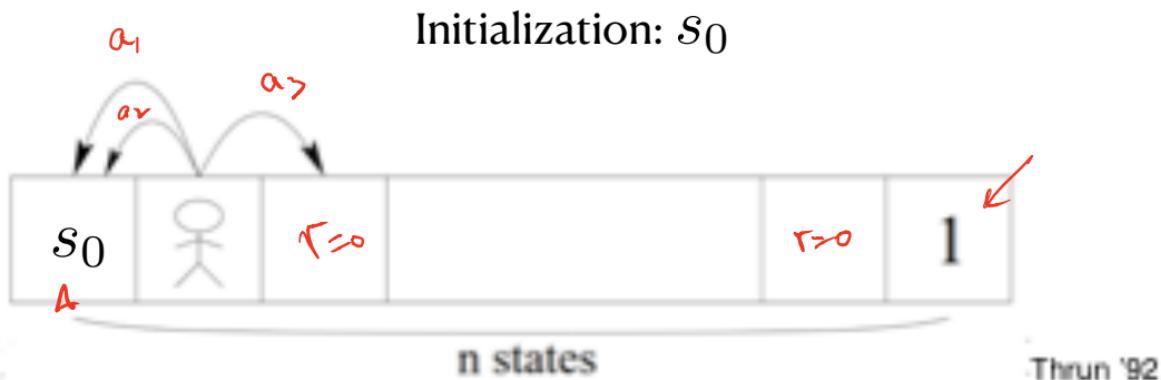
Only reset from  $\mu$ : we assume it's a delta distribution, all mass at a fixed  $s_0$

Unknown Transition  $P$  (for simplicity assume reward is known)

Different from the Generative Model Setting!

EXPLORATION!

## Why we need strategic exploration?

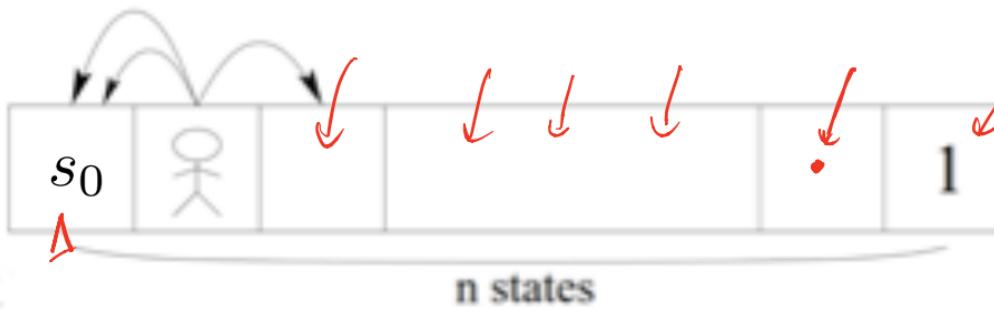


Length of chain is  $H$

$a$

## Why we need strategic exploration?

Initialization:  $s_0$



Thrun '92

Length of chain is  $H$

Probability of random walk hitting reward 1 is  $(1/3)^{-H}$

# **Learning Protocol**

## **Learning Protocol**

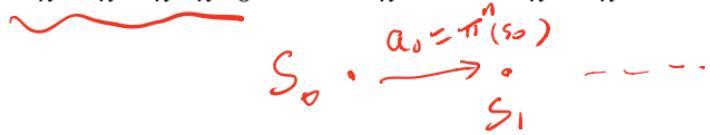
1. Learner initializes a policy  $\pi^1$

## Learning Protocol

1. Learner initializes a policy  $\pi^1$

2. At episode  $n$ , learner executes  $\pi^n$ :

$\{s_h^n, a_h^n, r_h^n\}_{h=0}^{H-1}$ , with  $a_h^n = \pi^n(s_h^n)$ ,  $r_h^n = r(s_h^n, a_h^n)$ ,  $s_{h+1}^n \sim P(\cdot | s_h^n, a_h^n)$



## Learning Protocol

1. Learner initializes a policy  $\pi^1$
2. At episode  $n$ , learner executes  $\pi^n$ :  
 $\{s_h^n, a_h^n, r_h^n\}_{h=0}^{H-1}$ , with  $a_h^n = \pi^n(s_h^n)$ ,  $r_h^n = r(s_h^n, a_h^n)$ ,  $s_{h+1}^n \sim P(\cdot | s_h^n, a_h^n)$
3. Learner updates policy to  $\pi^{n+1}$  using all prior information  
 $\underbrace{\{s_h^i, a_h^i\}_{h=0}^{H-1}}_{n}, i \geq 1$

## Learning Protocol

1. Learner initializes a policy  $\pi^1$
2. At episode  $n$ , learner executes  $\pi^n$ :  
 $\{s_h^n, a_h^n, r_h^n\}_{h=0}^{H-1}$ , with  $a_h^n = \pi^n(s_h^n)$ ,  $r_h^n = r(s_h^n, a_h^n)$ ,  $s_{h+1}^n \sim P(\cdot | s_h^n, a_h^n)$
3. Learner updates policy to  $\pi^{n+1}$  using all prior information

Performance measure: REGRET

$$\frac{1}{N} \mathbb{E} \left[ \sum_{n=1}^N (V^\star - V^{\pi^n}) \right] = \text{poly}(S, A, H) \sqrt{\frac{N}{N}} \quad \text{or } 0, \quad \frac{\sqrt{N}}{N} \rightarrow 0$$

## Notations for Today

$$= \sum_{s' \in S} p(s'|s,a) \cdot f(s')$$

$$\underbrace{\mathbb{E}_{s' \sim P(\cdot|s,a)} [f(s')]} := \underbrace{P(\cdot|s,a)} \cdot f$$

$d_h^\pi(s, a)$ : state-action distribution induced by  $\pi$  at time step  $h$   
(i.e., probability of  $\pi$  visiting  $(s, a)$  at time step  $h$  starting from  $s_0$ )

$$\pi = \{\pi_0, \dots, \pi_{H-1}\}$$

$\triangle$   
*Time-Dependent*

$$\pi_h : S \rightarrow A$$

$$s_0 \quad a_0 = \pi_0(s_0), \quad \dots \quad a_n = \pi_n(s_n) \quad \dots$$

## UCBVI: Optimistic Model-based Learning

*Episode*  
~~Inside iteration  $n$  :~~

## UCBVI: Optimistic Model-based Learning

**Inside iteration  $n$  :**

Use all previous data to estimate transitions  $\widehat{P}_{\emptyset}^n, \dots, \widehat{P}_{H-1}^n$

## UCBVI: Optimistic Model-based Learning

Inside iteration  $n$  :

Use all previous data to estimate transitions  $\widehat{P}_1^n, \dots, \widehat{P}_{H-1}^n$

Design reward bonus  $b_h^n(s, a), \forall s, a, h$

70

## UCBVI: Optimistic Model-based Learning

Inside iteration  $n$  :

Use all previous data to estimate transitions  $\widehat{P}_1^n, \dots, \widehat{P}_{H-1}^n$

Design reward bonus  $b_h^n(s, a), \forall s, a, h$

Optimistic planning with learned model:  $\pi^n = \underbrace{\text{Value-Iter}}_{\text{Value-Iter}} \left( \{ \widehat{P}_h^n, \underbrace{r_h + b_h^n}_{\text{Reward Bonus}} \}_{h=1}^{H-1} \right)$

## UCBVI: Optimistic Model-based Learning

Inside iteration  $n$  :

Use all previous data to estimate transitions  $\widehat{P}_1^n, \dots, \widehat{P}_{H-1}^n$

Design reward bonus  $b_h^n(s, a), \forall s, a, h$

Optimistic planning with learned model:  $\pi^n = \text{Value-Iter} \left( \{ \widehat{P}_h^n, r_h + b_h^n \}_{h=1}^{H-1} \right)$

Collect a new trajectory by executing  $\pi^n$  in the real world  $\{P_h\}_{h=0}^{H-1}$  starting from  $s_0$

$\{ s_0^n, a_0^n, \dots, s_{H-1}^n, a_{H-1}^n \}$

$s_{(n)}$

## UCBVI—Part 1: Model Estimation

Let us consider the **very beginning** of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h$$

## UCBVI—Part 1: Model Estimation

Let us consider the **very beginning** of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

## UCBVI—Part 1: Model Estimation

Let us consider the **very beginning** of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

$$\underbrace{N_h^n(s, a)}_{\text{sum}} = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h, \quad \underbrace{N_h^n(s, a, s')}_{\text{sum}} = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s, a, s')\}, \forall s, a, h$$

## UCBVI—Part 1: Model Estimation

Let us consider the **very beginning** of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

$$N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h, \quad N_h^n(s, a, s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s, a, s')\}, \forall s, a, h$$

Estimate model  $\widehat{P}_h^n(s' | s, a), \forall s, a, s', h :$

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}$$

## UCBVI—Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}} \approx cH \sqrt{\frac{1}{N_h^n(s, a)}}$$

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$$

Encourage to explore  
new state-actions

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$$

Encourage to explore  
new state-actions

**Value Iteration (aka DP) at episode n using  $\{\widehat{P}_h^n\}_h$  and  $\{r_h + b_h^n\}_h$**

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$$

Encourage to explore  
new state-actions

$h=0, \dots, h=H-1, \Delta$

**Value Iteration (aka DP) at episode n using  $\{\widehat{P}_h^n\}_h$  and  $\{r_h + b_h^n\}_h$**

$$\widehat{V}_H^n(s) = 0, \forall s$$

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$$

Encourage to explore  
new state-actions

**Value Iteration (aka DP) at episode n using  $\{\widehat{P}_h^n\}_h$  and  $\{r_h + b_h^n\}_h$**

$$\begin{aligned} & r(s, a) \in [0, 1] \\ & ||\varphi||_\infty \leq H \end{aligned}$$

$$\widehat{V}_H^n(s) = 0, \forall s \quad \widehat{Q}_h^n(s, a) = \min \left\{ \underbrace{r_h(s, a) + b_h^n(s, a)}_{\Delta} + \widehat{P}_h^n(\cdot | s, a) \cdot \underbrace{\widehat{V}_{h+1}^n}_{\Delta}, \underbrace{H}_{\Delta} \right\}, \forall s, a$$

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$$

Encourage to explore  
new state-actions

**Value Iteration (aka DP) at episode  $n$  using  $\{\widehat{P}_h^n\}_h$  and  $\{r_h + b_h^n\}_h$**

$$\widehat{V}_H^n(s) = 0, \forall s \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, \quad H \right\}, \forall s, a$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s \quad h=H, \dots, 0$$
$$\pi^n = \left\{ \pi_0^n, \dots, \pi_H^n \right\}$$

## UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode  $n$ :

$$\mathcal{D}_h^n = \{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}, \forall h, \quad N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h,$$

$$b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$$

Encourage to explore  
new state-actions

**Value Iteration (aka DP) at episode n using  $\{\widehat{P}_h^n\}_h$  and  $\{r_h + b_h^n\}_h$**

$$\widehat{V}_H^n(s) = 0, \forall s \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n(\cdot), \forall s, a \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

$$\|\widehat{V}_h^n\|_\infty \leq H, \forall h, n$$

[s, a]

## UCBVI: Put All Together

For  $n = 1 \rightarrow N$ :

1. Set  $N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h$  ✓

2. Set  $N_h^n(s, a, s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s, a, s')\}, \forall s, a, a', h$  ✓

3. Estimate  $\widehat{P}^n : \widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall s, a, s', h$  ✓

4. Plan:  $\pi^n = VI \left( \{ \widehat{P}_h^n, r_h + b_h^n \}_h \right)$ , with  $b_h^n(s, a) = cH \sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$   $\approx H \sqrt{\frac{1}{N(s,a)}}$

5. Execute  $\pi^n : \{s_0^n, a_0^n, r_0^n, \dots, s_{H-1}^n, a_{H-1}^n, r_{H-1}^n, s_H^n\}$

$$s_0^n \leftarrow s_0$$

## Theorem: UCBVI Regret Bound

$$\mathbb{E} \left[ \text{Regret}_N \right] := \mathbb{E} \left[ \sum_{n=1}^N (V^\star - V^{\pi^n}) \right] \leq \widetilde{O} \left( \underbrace{H^2 \sqrt{S^2 AN}} \right)$$

## Theorem: UCBVI Regret Bound

$$\mathbb{E} [\text{Regret}_N] := \mathbb{E}_{\Delta} \left[ \sum_{n=1}^N (V^\star - V^{\pi^n}) \right] \leq \widetilde{O} \left( H^2 \sqrt{S^2 AN} \right)$$

### Remarks:

Note that we consider expected regret here (policy  $\pi^n$  is a random quantity).  
High probability version is not hard to get (need to do a martingale argument)

## Theorem: UCBVI Regret Bound

$$\mathbb{E} [\text{Regret}_N] := \mathbb{E} \left[ \sum_{n=1}^N (V^\star - V^{\pi^n}) \right] \leq \widetilde{O} \left( H^2 \sqrt{S^2 A N} \right)$$

### Remarks:

Note that we consider expected regret here (policy  $\pi^n$  is a random quantity).  
High probability version is not hard to get (need to do a martingale argument)

Dependency on  $H$  and  $S$  are suboptimal; but the same algorithm can achieve  $H^2 \sqrt{S A N}$  in the leading term [Azar et.al 17 ICML]

## Outline of Proof

Bonus  $b_h^n(s, a)$  is related to  $\left( \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^{\star} \right)$

## Outline of Proof

Bonus  $b_h^n(s, a)$  is related to  $\left( \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^{\star} \right)$

VI with bonus inside the learned model gives optimism, i.e.,  $\widehat{V}_h^n(s) \geq V_h^{\star}(s), \forall h, n, s, a$

## Outline of Proof

Bonus  $b_h^n(s, a)$  is related to  $\left( \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^{\star} \right)$

VI with bonus inside the learned model gives optimism, i.e.,  $\widehat{V}_h^n(s) \geq V_h^{\star}(s), \forall h, n, s, a$

Upper bound per-episode regret:  $V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - \underbrace{V_0^{\pi^n}(s_0)}_{\in \pi^n(\mathcal{R}, \mathcal{C})}$

$\pi^n, \widehat{\rho}^n, r_n + b_n$

## Outline of Proof

Bonus  $b_h^n(s, a)$  is related to  $\left( \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^\star \right)$

VI with bonus inside the learned model gives optimism, i.e.,  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall h, n, s, a$

Upper bound per-episode regret:  $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

Apply simulation lemma:  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \underset{\{\widehat{P}^n, r + b_h^n\}}{\sim} \{\bar{P}, \bar{r}\}$

## High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret:  $V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

*optimism*

$\widehat{V}^n \geq V^*$

1. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$ ?  
*small*

## High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret:  $V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$  ↗ ↗

1. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$ ?

$$V^* - V^{\pi^n} \leq \epsilon$$

Then  $\pi^n$  is close to  $\pi^*$ , i.e., we are doing exploitation

## High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret:  $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$ ?

Then  $\pi^n$  is close to  $\pi^\star$ , i.e., we are doing exploitation

$$(\widehat{V}_0^n(s_0) > V_0^{\pi^n}(s_0))$$

2. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$  ?

$$\begin{aligned} UCB &= \widehat{V}_0^n \\ &\uparrow \\ &\rightarrow V^{\pi^n} \end{aligned}$$

## High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret:  $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$ ?

Then  $\pi^n$  is close to  $\pi^\star$ , i.e., we are doing exploitation

2. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$  ?

$$\epsilon \leq \underbrace{\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)}_{\sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ \underbrace{b_h^n(s,a)}_{\widehat{P}_h^n(\cdot | s,a) - P_h(\cdot | s,a)} + \underbrace{(\widehat{P}_h^n(\cdot | s,a) - P_h(\cdot | s,a)) \cdot \widehat{V}_{h+1}^n} \right]}$$

## High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret:  $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$ ?

Then  $\pi^n$  is close to  $\pi^\star$ , i.e., we are doing exploitation

2. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$  ?

$$\epsilon \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

We collect data at steps where bonus is large or model is wrong, i.e., exploration

# 1. Model Error using Hoeffding's inequality & Union Bound

$$b_n^{(s,a)} \approx (\hat{P}_n - P_n) \cdot V_{htl}$$
$$\widehat{P}_h^n(s'|s,a) = \frac{N_h^n(s,a,s')}{N_h^n(s,a)}, \forall h, s, a, s'$$

# 1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function  $f: S \mapsto [0, H]$ , w/ prob  $1 - \delta$  :

$$f = \nabla_{ht}^*$$

$$\left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^T f \right| \leq O(H \underbrace{\sqrt{\ln(SAHN/\delta)/N_h^n(s, a)}}_{\Delta}), \forall \underbrace{s, a, h, N}$$

# 1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function  $f: S \mapsto [0, H]$ , w/ prob  $1 - \delta$  :

$$\left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O(H \sqrt{\ln(SAHN/\delta)/N_h^n(s, a)}), \forall s, a, h, N$$

Bonus  $b_h^n(s, a)$

## 1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function  $f: S \mapsto [0, H]$ , w/ prob  $1 - \delta$  :

$$\left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O(H \sqrt{\ln(SAHN/\delta)/N_h^n(s, a)}), \forall s, a, h, N$$

Bonus  $b_h^n(s, a)$

From now on, assume this event being true

## 1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function  $f: S \mapsto [0, H]$ , w/ prob  $1 - \delta$  :

$$\left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O(H \sqrt{\ln(SAHN/\delta)/N_h^n(s, a)}), \forall s, a, h, N$$

Bonus  $b_h^n(s, a)$

From now on, assume this event being true

Intuition:

## 1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function  $f: S \mapsto [0, H]$ , w/ prob  $1 - \delta$ :

$$\left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O(H \sqrt{\ln(SAHN/\delta)/N_h^n(s, a)}), \forall s, a, h, N$$

$$\left| (\widehat{P}_{sa} - P_{sa}) \cdot V^* \right| \quad \text{Bonus } b_h^n(s, a)$$

From now on, assume this event being true

**Intuition:**

1. Assume for some  $i$ ,  $s_h^i = s, a_h^i = a$ , then  $f(s_{h+1}^i)$  is an unbiased estimate of  $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} f(s')$

# 1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

for H-improvement,

use Bernstein

Rather than  
Hoeffding's

Given a fixed function  $f: S \mapsto [0, H]$ , w/ prob  $1 - \delta$ :

$$\left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^T f \right| \leq O(H \sqrt{\ln(SAHN/\delta)/N_h^n(s, a)}), \forall s, a, h, N$$

Bonus  $b_h^n(s, a)$

$E \sup_{s \sim P_h} V(s')$

From now on, assume this event being true

Intuition:

1. Assume for some  $i$ ,  $s_h^i = s, a_h^i = a$ , then  $f(s_{h+1}^i)$  is an unbiased estimate of  $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} f(s')$

2. Note  $\widehat{P}_h^n(\cdot | s, a) \cdot f = \frac{1}{N_h^n(s, a)} \sum_{i=1}^{n-1} \mathbf{1}[(s_h^i, a_h^i) = (s, a)] f(s_{h+1}^i)$

## 2. Proving Optimism via Induction

**Lemma [Optimism]:**  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\begin{aligned}\widehat{V}_H^n(s) &= 0, & \widehat{Q}_h^n(s, a) &= \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\} \\ \widehat{V}_h^n(s) &= \max_a \widehat{Q}_h^n(s, a), & \pi_h^n(s) &= \arg \max_a \widehat{Q}_h^n(s, a), \forall s\end{aligned}$$

## 2. Proving Optimism via Induction

**Lemma [Optimism]:**  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\begin{aligned}\widehat{V}_H^n(s) &= 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\} \\ \widehat{V}_h^n(s) &= \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s\end{aligned}$$

Inductive hypothesis:  $\underbrace{\widehat{V}_{h+1}^n(s)}_{\text{red}} \geq V_{h+1}^\star(s), \quad \forall s$

## 2. Proving Optimism via Induction

**Lemma [Optimism]:**  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$
$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

①  $\widehat{Q} = H$

②  $\widehat{Q} =$

Inductive hypothesis:  $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^\star(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^\star(s, a) = \underbrace{r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n}_{\text{update}}, - \underbrace{r_h(s, a) + P_h(\cdot | s, a) \cdot V_{h+1}^\star}_{\text{Bell-Eqn}}$$

## 2. Proving Optimism via Induction

**Lemma [Optimism]:**  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\begin{aligned}\widehat{V}_H^n(s) &= 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\} \\ \widehat{V}_h^n(s) &= \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s\end{aligned}$$

Inductive hypothesis:  $\underbrace{\widehat{V}_{h+1}^n(s)}_{\geq V_{h+1}^\star(s)}, \quad \forall s$

$$\begin{aligned}\widehat{Q}_h^n(s, a) - Q_h^\star(s, a) &= r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - r_h(s, a) - P_h(\cdot | s, a) \cdot V_{h+1}^\star \\ &\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \underbrace{V_{h+1}^\star}_{\geq \widehat{V}_{h+1}^n} - P_h(\cdot | s, a) \cdot \underbrace{V_{h+1}^\star}_{\geq \widehat{V}_{h+1}^n}\end{aligned}$$

## 2. Proving Optimism via Induction

**Lemma [Optimism]:**  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\begin{aligned}\widehat{V}_H^n(s) &= 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\} \\ \widehat{V}_h^n(s) &= \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s\end{aligned}$$

Inductive hypothesis:  $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^\star(s), \forall s$

$$\begin{aligned}\widehat{Q}_h^n(s, a) - Q_h^\star(s, a) &= r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - r_h(s, a) - P_h(\cdot | s, a) \cdot V_{h+1}^\star \\ &\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot V_{h+1}^\star - P_h(\cdot | s, a) \cdot V_{h+1}^\star \\ &= b_h^n(s, a) + \underbrace{\left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^\star}_{\leq H \sqrt{\frac{1}{N_h^n(s, a)}}} \approx b_h^n(s, a)\end{aligned}$$

## 2. Proving Optimism via Induction

**Lemma [Optimism]:**  $\widehat{V}_h^n(s) \geq V_h^\star(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\begin{aligned}\widehat{V}_H^n(s) &= 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\} \\ \widehat{V}_h^n(s) &= \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s\end{aligned}$$

Inductive hypothesis:  $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^\star(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^\star(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - r_h(s, a) - P_h(\cdot | s, a) \cdot V_{h+1}^\star$$

$$\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot V_{h+1}^\star - P_h(\cdot | s, a) \cdot V_{h+1}^\star$$

$$\widehat{Q}_h^n \geq Q_h^\star$$

$$= b_h^n(s, a) + \left| \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^\star \right| \leq b_h^n(s, a)$$

$$\Rightarrow \widehat{V}_h^n(s) \geq V_h^\star(s), \forall s$$

$$\geq b_h^n(s, a) - b_h^n(s, a) = 0, \quad \forall s, a$$

for the case where reward is unknown:

### 3. Upper Bounding Regret using Optimism

$$\hat{r}_h^n(s,a) = \frac{1}{N_h^n(s,a)} \sum_{t=1}^{n^l} r_h^t \mathbb{1}(s_h^t, a_h^t = s,a)$$

$$|\hat{r}(s,a) - r(s,a)| \leq \sqrt{\frac{1}{N_h^n(s,a)}}$$

per-episode regret :=  $V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \hat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$

A  
unknown

$$\pi^n, \hat{p}, r+b_h$$

This is something  
we can control!  
And this is related  
to our policy  $\pi^n$

(Here, reward  
is known)

Recall simulation lemma — the lemma measures the difference of a policy under two MDPs

## 4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\left\{ \pi^n, \widehat{P}, r, b_h \right\} \quad \widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\left\{ \pi^n, P, r \right\}$$

## 4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) = \widehat{Q}_0^n(s_0, \underbrace{\pi^n(s_0)}_{\text{underlined}}) - \underbrace{Q_0^{\pi^n}(s_0, \pi^n(s_0))}_{\text{underlined}}$$

## 4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min_a \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) = \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0))$$

$$\leq \cancel{r_0(s_0, \pi^n(s_0))} + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - \cancel{r_0(s_0, \pi^n(s_0))} - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}$$

A  $\min\{a, b\} \leq a$

Bell-Eqn

## 4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\begin{aligned} \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) &= \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0)) \\ &\leq r_0(s_0, \pi^n(s_0)) + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - r_0(s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n} \\ &= \underbrace{b_h^n(s_0, \pi^n(s_0))}_{\text{Red}} + \underbrace{\widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n}_{\Delta} - \underbrace{P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}}_{\Delta} \end{aligned}$$

## 4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\begin{aligned}
 \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) &= \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0)) \\
 &\stackrel{\textcolor{red}{\Delta}}{=} r_0(s_0, \pi^n(s_0)) + b_0^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - r_0(s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n} \\
 &= b_0^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n} \\
 &= b_0^n(s_0, \pi^n(s_0)) + \left( \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \right) \cdot \widehat{V}_1^n + P_0(\cdot | s_0, \pi^n(s_0)) \cdot \left( \widehat{V}_1^n - V_1^{\pi^n} \right)
 \end{aligned}$$

## 4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\begin{aligned} \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) &= \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0)) \\ &\leq r_0(s_0, \pi^n(s_0)) + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - r_0(s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n} \\ &= b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n} \\ &= b_h^n(s_0, \pi^n(s_0)) + \underbrace{\left( \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \right)}_{(1)} \cdot \widehat{V}_1^n + \underbrace{P_0(\cdot | s_0, \pi^n(s_0)) \cdot \left( \widehat{V}_1^n - V_1^{\pi^n} \right)}_{(2) \text{ Recursion}} \\ &= \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + \underbrace{(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a))}_{(1)} \cdot \widehat{V}_{h+1}^n \right] \end{aligned}$$

## 4. Upper bounding Regret via Simulation Lemma

$$\begin{aligned} \text{per-episode regret} &:= V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0) \\ &\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right] \end{aligned}$$

#### 4. Upper bounding Regret via Simulation Lemma

$$\text{per-episode regret} := V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$$

But  $\widehat{V}_h^n$  is data-dependent  
(this is different from  $V_h^*$ ) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's  
inequality

$$|(\widehat{P}_h^n - P_h) \cdot \widehat{V}_{h+1}^n| \leq \sqrt{\frac{1}{N}}$$

## 4. Upper bounding Regret via Simulation Lemma

per-episode regret :=  $V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$

But  $\widehat{V}_h^n$  is data-dependent  
(this is different from  $V_h^\star$ ) !!!

Let's do Holder's inequality

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$a \cdot b \leq \|a\|_1 \|b\|_\infty$$

Truncate by  $H$

$$\left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \underbrace{\|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1}_{\leq H} \underbrace{\|\widehat{V}_{h+1}^n\|_\infty}_{\leq H}$$

## 4. Upper bounding Regret via Simulation Lemma

per-episode regret :=  $V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$

But  $\widehat{V}_h^n$  is data-dependent  
(this is different from  $V_h^\star$ ) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\|P_h - \widehat{P}_h(s, a)\|_1 \leq \sqrt{\frac{s \ln(1/\delta)}{N_h^n(s, a)}}$$

$$(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{with prob } 1 - \delta$$

## 4. Upper bounding Regret via Simulation Lemma

$$\text{per-episode regret} := V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$$

But  $\widehat{V}_h^n$  is data-dependent  
(this is different from  $V_h^\star$ ) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$b(s,a) = H \sqrt{\frac{\ln(1/\delta)}{N_h^n(s,a)}}$$

$$\left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{with prob } 1 - \delta$$

## 4. Upper bounding Regret via Simulation Lemma

$$\text{per-episode regret} := V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$$

But  $\widehat{V}_h^n$  is data-dependent  
(this is different from  $V_h^\star$ ) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$\leq 2 \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$\left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{with prob } 1 - \delta$$

## 4. Upper bounding Regret via Simulation Lemma

per-episode regret :=  $V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$

But  $\widehat{V}_h^n$  is data-dependent  
(this is different from  $V_h^\star$ ) !!!

Let's do Holder's inequality

$$\begin{aligned} &\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right] \\ &\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right] \\ &\quad \text{--- } (\widehat{P}-P) \cdot \widehat{V} \\ &\quad = (\widehat{P}-P) \cdot (V^*) + (\widehat{P}-P) \cdot (\widehat{V}-V^*) \\ &\leq 2 \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right] = 2H \sqrt{S \ln(SAHN/\delta)} \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ \sqrt{\frac{1}{N_h^n(s, a)}} \right] \\ &\quad \text{--- } \left( \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty \\ &\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{with prob } 1 - \delta \end{aligned}$$

## 5. Final Step

$$|(\hat{P} - P) \cdot V^*| \leq b_n \quad \|(\hat{P} \cdot (\cdot_{sa}) - P \cdot (\cdot_{sa}))\|_1 \leq \sqrt{s} b_n(s_{-a})$$

Remember we had two failure events for bounding transitions errors.

## 5. Final Step

Remember we had two failure events for bounding transitions errors.

$$\mathbb{E} [\text{Regret}_N] = \mathbb{E} \left[ \underbrace{\mathbf{1}\{\text{events hold}\}}_{\text{red underline}} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[ \underbrace{\mathbf{1}\{\text{events don't hold}\}}_{\text{red underline}} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right]$$

## 5. Final Step

Remember we had two failure events for bounding transitions errors.

$$\begin{aligned}\mathbb{E} [\text{Regret}_N] &= \mathbb{E} \left[ \mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[ \mathbf{1}\{\text{events don't hold}\} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] \\ &\leq \mathbb{E} \left[ \mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{P}(\text{events don't hold}) \cdot NH\end{aligned}$$

$\approx \delta$

$\leq NH$

## 5. Final Step

Remember we had two failure events for bounding transitions errors.

$$\begin{aligned}\mathbb{E} [\text{Regret}_N] &= \mathbb{E} \left[ \mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[ \mathbf{1}\{\text{events don't hold}\} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] \\ &\leq \mathbb{E} \left[ \mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^\star(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{P}(\text{events don't hold}) \cdot NH \\ &\leq H\sqrt{S \ln(SANH/\delta)} \mathbb{E} \left[ \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} \right] + 2\delta NH\end{aligned}$$

$\leq \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}}$

$N_h^n(s_h^n, a_h^n)$

$n=1, h=0, \dots, H-1$

$+ 2\delta NH$

## 5. Final Step

$$\sum_{n=1}^N \left( \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} \right)$$

## 5. Final Step

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} = \underbrace{\left( \sum_{h=0}^{H-1} \right)_{s,a}}_{\text{A}} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}}$$

## 5. Final Step

$$\sum_{i=1}^n \frac{1}{\sqrt{i}} \leq 2\sqrt{n} + 1$$

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} = \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)}$$

## 5. Final Step

CS-inequality

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} = \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \underbrace{\sqrt{N_h^N(s,a)}}_{\text{CS-inequality}} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)}$$

$$\begin{aligned} & \sum_{s,a} 1 \cdot \sqrt{N_h^N(s,a)} \\ & \leq \sqrt{\sum_{s,a} 1} \quad \sqrt{\sum_{s,a} N_h^N(s,a)} \\ & = SA \end{aligned}$$

## 5. Final Step

$$\begin{aligned} \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} &= \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{\sum_{s,a} N_h^N(s,a)} \\ &\quad \text{A red oval encloses } \sum_{s,a} N_h^N(s,a) \text{ with an equals sign below it.} \\ &\leq \sum_{h=0}^{H-1} \sqrt{SAN} = H\sqrt{SAN} \end{aligned}$$

## 5. Final Step

$$\begin{aligned} \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} &= \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)} \\ &\leq \sum_{h=0}^{H-1} \sqrt{SAN} = H\sqrt{SAN} \end{aligned}$$

$$\mathbb{E} [\text{Regret}_N] \leq \underbrace{2H^2 S \sqrt{AN \ln(SAHN/\delta)}} + \underbrace{2\delta NH}$$

## 5. Final Step

$$\begin{aligned}
 \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} &= \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)} \\
 &\leq \sum_{h=0}^{H-1} \sqrt{SAN} = H\sqrt{SAN}
 \end{aligned}$$

$$\mathbb{E} [\text{Regret}_N] \leq 2H^2 S \sqrt{AN \ln(SAHN/\delta)} + 2\delta NH$$

Set  $\delta = 1/(HN)$

$$\leq 2H^2 S \sqrt{AN \cdot \ln(SAH^2N^2)} = \widetilde{O}\left(H^2 S \sqrt{AN}\right)$$

$$S\sqrt{H^2 \sqrt{SAN}}$$

## High-level Idea: Exploration or Exploitation Tradeoff

$$\text{Upper bound per-episode regret: } V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

*optimism (Induction)*  
*≤ simulation lemma*

1. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$ ?

Then  $\pi^n$  is close to  $\pi^*$ , i.e., we are doing exploitation

Arg<sup>ue</sup>:  $N^{\frac{2}{3}}$

2. What if  $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$  ?

$$\epsilon \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[ b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$\sqrt{N}$

We collect data at steps where bonus is large or model is wrong, i.e., exploration