

# Upper Bounds in RL 2

August 6, 2021

# Background: Markov Decision Processes

- $X$  - state space,  $(X_t)_{t \geq 0}$  - sequence of random states;
- $A$  - action space,  $(A_t)_{t \geq 0}$  - sequence of random actions;
- Policy  $\pi$  - distribution on  $A$ ;

$$\pi(a|x) = P(A_t = a | X_t = x)$$

- Markov kernel  $P^a(x'|x) = P(X_t = x' | X_{t-1} = x, A_{t-1} = a)$ ;
- Deterministic reward  $r : X \times A \rightarrow \mathbb{R}$ ;
- At step  $t$  in the state  $X_t = x$  the agent performs action  $A_t = a \sim \pi(\cdot|x)$ , transits to  $X_{t+1} = x' \sim P^a(\cdot|x)$  and obtains reward  $r_{t+1} = r(x, a)$ ;

# How to measure policy's quality?

Given two policies  $\pi_1$  and  $\pi_2$ , how to compare their quality?

## Definition

Let's fix some policy  $\pi$ . The state value function for a state  $x$  at time  $t$  is

$$V_t^\pi(x) = E_\pi\left[\sum_{k=t}^{\infty} \gamma^k r_{k+1} | x_t = x\right]$$

The ways of estimating  $V_t^\pi$ :

- Monte-Carlo: run a series of independent simulations;
- Temporal Difference-based methods (TD);

# Aim

- Approximation of optimal solution - finding optimal policies and optimal value functions. How good is the policy we learned from the approximation procedures? How far are we from the optimal solution?
- To give an answer, upper bounds of optimal state-value function can be evaluated.

# Theory

- For getting the upper bound of optimal Value Function we can use the algorithm based on the next principles.
- Let's  $P^a V^*(x) = \int_{\mathcal{X}} V^*(x') P^a(x'|x) dx'$ , then

## Definition

We call a function  $V^{\text{up}}$  the *upper solution* for  $V^*$  if it satisfies

$$V^{\text{up}} \geq \max_{a \in A} \{r(x, a) + \gamma P^a V^{\text{up}}\}.$$

# Theory

How to construct  $V^{\text{up}}$  ?

Denote  $Y^{x,a} \sim P(\cdot|x, a)$  and  $P^a\Phi(x) = 0$  for all  $a \in A$ .

- Then

$$V^{\text{up}}(x) = E[\max_a \{r(x, a) + \gamma(V^{\text{up}}(Y^{x,a}) - \Phi(Y^{x,a}))\}] \quad (1)$$

We may choose  $\Phi_{\pi}^{x,a}(y) = V^{\pi}(y) - (P^a V^{\pi})(x)$ . And due to Banach's fixed point theorem, we can find upper solutions via iterations (UVIP Algorithm):

$$V_{k+1}^{\text{up},\pi}(x) = E[\max_a \{r^a(x) + \gamma(V_k^{\text{up},\pi}(Y^{x,a}) - \Phi_{\pi}^{x,a}(Y^{x,a}))\}] \quad (2)$$

An important property:  $V_{k+1}^{\text{up},\pi}(x) \geq V^*(x)$  for any  $x \in X$  and  $k \in N$ , provided that  $V_0^{\text{up},\pi}(x) \geq V^*(x)$  and  $\gamma < 1$ ;

# Upper Value Iteration Procedure

For the fixed policy  $\pi$  and its value function we can construct upper solution with the following scheme

$$V_{k+1}^{\text{up}}(x) = \mathbb{E} \left[ \max_{a \in A} \{ r^a(x) + \gamma (V_k^{\text{up}}(Y^{x,a}) - V^{\pi}(Y^{x,a}) + (P_k^a V^{\pi})(Y^{x,a})) \} \right]$$

# Upper Value Iteration Procedure

We can introduce several improvements for the basic UVIP algorithm

- Incremental UVIP:

$$V_{k+1}^{\text{up}}(x) = V_k^{\text{up}}(x) + \alpha_k \delta_{k+1}(V_k^{\text{up}}) \mathbb{1}_{\{x=X_k\}},$$

where

$$\begin{aligned} \delta_{k+1}(V) = \max_a \{ & r^a(X_k) + \gamma(V(Y_{k+1}^a) \\ & - V^\pi(Y_{k+1}^a) + (P^a V^\pi(X_k))) \} - V(X_k) \end{aligned}$$



# Upper Value Iteration Procedure

We can introduce several improvements for the basic UVIP algorithm

- Incremental UVIP:

$$V_{k+1}^{\text{up}}(x) = V_k^{\text{up}}(x) + \alpha_k \delta_{k+1}(V_k^{\text{up}}) \mathbb{1}_{\{x=X_k\}},$$

where

$$\begin{aligned} \delta_{k+1}(V) = \max_a \{ & r^a(X_k) + \gamma(V(Y_{k+1}^a) \\ & - V^\pi(Y_{k+1}^a) + (P^a V^\pi(X_k))) \} - V(X_k) \end{aligned}$$

- Take previous iterations of the UVIP for computing  $\Phi_\pi^{x,a}$

# Upper Value Iteration Procedure

We can introduce several improvements for the basic UVIP algorithm

- Incremental UVIP:

$$V_{k+1}^{\text{up}}(x) = V_k^{\text{up}}(x) + \alpha_k \delta_{k+1}(V_k^{\text{up}}) \mathbb{1}_{\{x=X_k\}},$$

where

$$\begin{aligned} \delta_{k+1}(V) = \max_a \{ & r^a(X_k) + \gamma(V(Y_{k+1}^a) \\ & - V^\pi(Y_{k+1}^a) + (P^a V^\pi(X_k))) \} - V(X_k) \end{aligned}$$

- Take previous iterations of the UVIP for computing  $\Phi_\pi^{x,a}$
- Transition kernel approximation

# Upper Value Iteration Procedure

We can introduce several improvements for the basic UVIP algorithm

- Incremental UVIP:

$$V_{k+1}^{\text{up}}(x) = V_k^{\text{up}}(x) + \alpha_k \delta_{k+1}(V_k^{\text{up}}) \mathbb{1}_{\{x=X_k\}},$$

where

$$\delta_{k+1}(V) = \max_a \{ r^a(X_k) + \gamma(V(Y_{k+1}^a) - V^\pi(Y_{k+1}^a) + (P^a V^\pi(X_k))) \} - V(X_k)$$

- Take previous iterations of the UVIP for computing  $\Phi_\pi^{x,a}$
- Transition kernel approximation
- Real-Time Upper Solutions

# Infinite-horizon case

- Take previous iterations of the UVIP for computing  $\Phi_{\pi}^{x,a}$

---

## Algorithm 1: UVIPv2

---

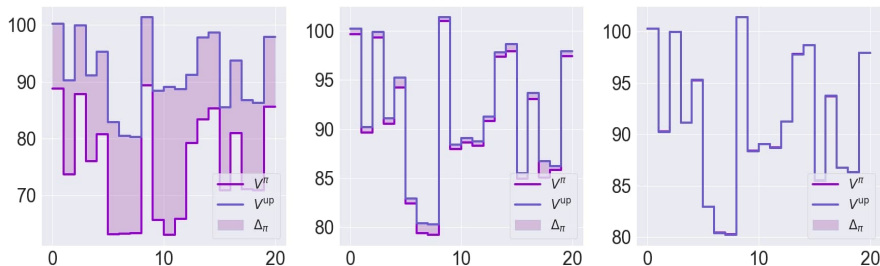
**Input:**  $V^{\pi}, V_0^{\text{up}}, V_1^{\text{up}}, \gamma, \varepsilon$

**Result:**  $V^{\text{up}}$

```
k = 1; while  $\|V_k^{\text{up}} - V_{k-1}^{\text{up}}\|_X > \varepsilon$  do
  for  $x \in X$  do
    for  $x \in X, a \in A$  do
      for  $y \in X$  do
         $\Phi_{\pi}^{x,a}(y) = V_{k-1}^{\text{up}}(y) - (P^a V_{k-1}^{\text{up}})(x);$ 
      end
    end
     $V_{k+1}^{\text{up}}(x) = E[\max_a \{r^a(x) + \gamma(V_k^{\text{up}}(Y^{x,a}) - \Phi_{\pi}^{x,a}(Y^{x,a}))\}], \quad Y^{x,a} \sim P^a(\cdot|x);$ 
  end
  k = k + 1;
end
 $V^{\text{up}} = V_k^{\text{up}}.$ 
```

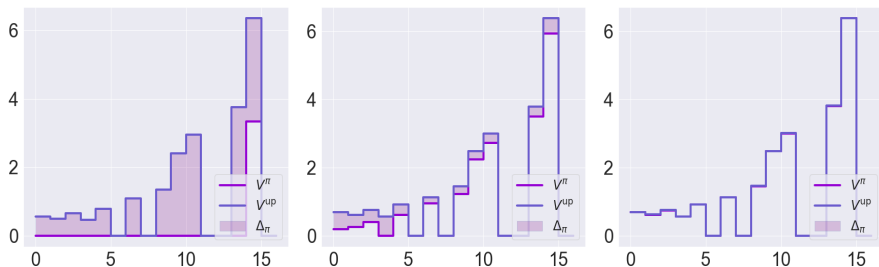
---

# Garnet + Monte-Carlo



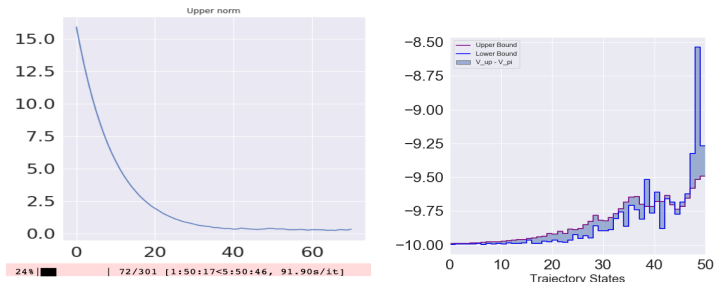
**Figure:** The difference between  $V_k^{\text{up}}$  and  $V^{\pi}$ . X-axis represents states of Garnet environment. All of the pairs represents steps of Value Iteration procedure.

# Frozen Lake + Monte-Carlo



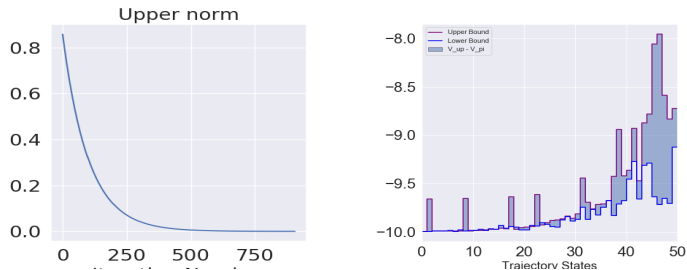
**Figure:** The difference between  $V_k^{\text{up}}$  and  $V^\pi$ . X-axis represents states of Frozen Lake environment. All of the pairs represents steps of Value Iteration procedure.

# Acrobot + Monte-Carlo



**Figure:** The difference between  $V^{\text{up},k}$  and  $V^{\pi}$ . X-axis represents states of Acrobot environment. The pictures are computed with DQN policy.

# Acrobot + Incremental



**Figure:** The difference between  $V^{\text{up},k}$  and  $V^{\pi}$ . X-axis represents states of Acrobot environment. The pictures are computed with DQN policy.



## UVIP with empirical P

- Empirically estimate the markov kernel  $P^a$
- Use  $\hat{P}_k^a(y|x)$  for sampling  $Y^{a,X}$

---

### Algorithm 2: (Incremental + Empirical P)

---

**Input:**  $V^\pi$ ,  $V_0^{\text{up}}$ ,  $\alpha$

Estimate  $\hat{P}_0^a(\cdot|x)$ ;

**for**  $k = 1, 2, \dots$  **do**

    Update  $\hat{P}_k^a(\cdot|x)$ ;

**for**  $X \in \mathcal{X}$  **do**

$$\delta_k = \max_{a \in \mathcal{A}} \{ r^a(X) + \gamma(V_{k-1}^{\text{up}}(Y^{a,X}) - V^\pi(Y^{a,X}) + (\hat{P}_k^a V^\pi)(X)) \} - V_{k-1}^{\text{up}}(X),$$

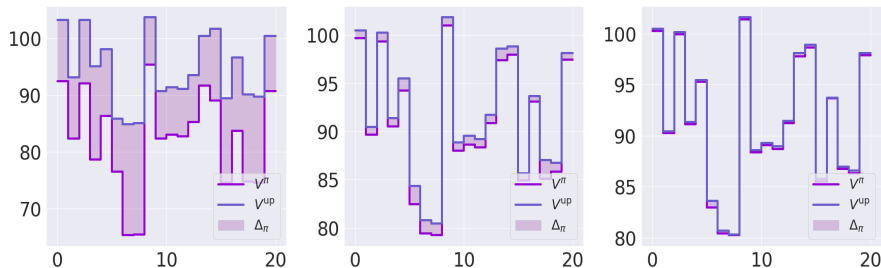
$$Y^{a,X} \sim \hat{P}_k^a(\cdot|x)$$

$$V_k^{\text{up}}(X) = V_{k-1}^{\text{up}}(X) + \alpha \delta_k$$

**end**

**end**

# Garnet



**Figure:** The difference between  $V^{\text{up},k}$  and  $V^{\pi}$ . X-axis represents states of Garnet environment. All of the pairs represents steps of Value Iteration procedure.

# Garnet

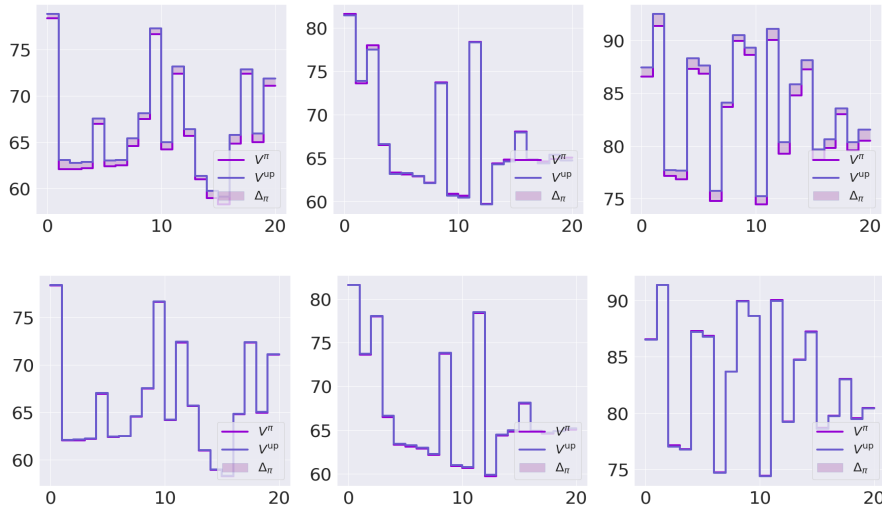
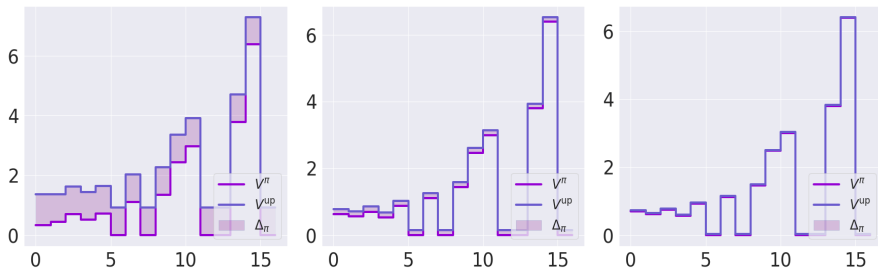


Figure: First row – sampling from  $\hat{P}_k^a(y|x)$ , second row – from  $P^a(y|x)$

# Frozen Lake



**Figure:** The difference between  $V^{up,k}$  and  $V^\pi$ . X-axis represents states of Frozen Lake environment. All of the pairs represents steps of Value Iteration procedure.

# Finite-horizon case

---

**Algorithm 3:** Real-Time Upper Solutions(RTUS)

---

**Input:**  $V^\pi, \forall x \in X, \forall t \in [H], V_t^{\text{up},0}, V_t^{\text{up},1}$

**for**  $k = 1, 2, \dots$  **do**

    Initialize  $X_1^k$

**for**  $t = 1, \dots, H$  **do**

$A_t^k \in \arg \max_{a \in A} \{r^a(X_t^k) + P^a V_{t+1}^{\text{up},k-1}(X_t^k)\}$

        Act with  $A_t^k$  and observe  $X_{t+1}^k$

$\delta_t^k(x) =$

$(\max_{a \in A} \{r^a(x) + V_{t+1}^{\text{up},k-1}(Y_{t+1}^{a,k}) - V_{t+1}^{\text{up},k-2}(Y_{t+1}^{a,k}) +$

$(P^a V_{t+1}^{\text{up},k-2})(x)\} - V_t^{\text{up},k-1}(x)) \cdot \mathbb{I}\{x = X_t^k\},$

$Y_{t+1}^{a,k} \sim P^a(\cdot | X_t^k)$

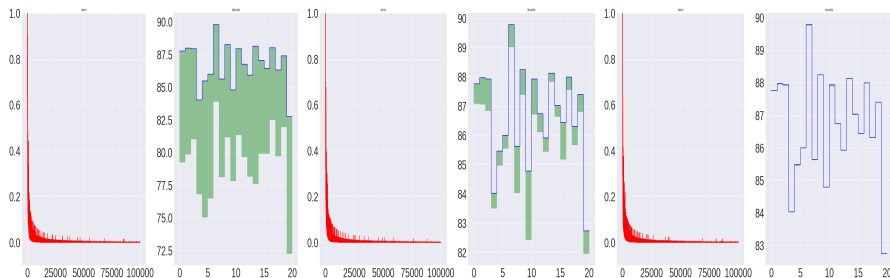
$V_t^{\text{up},k}(x) = V_t^{\text{up},k-1}(x) + \alpha_t \delta_t^k(x)$

**end**

**end**

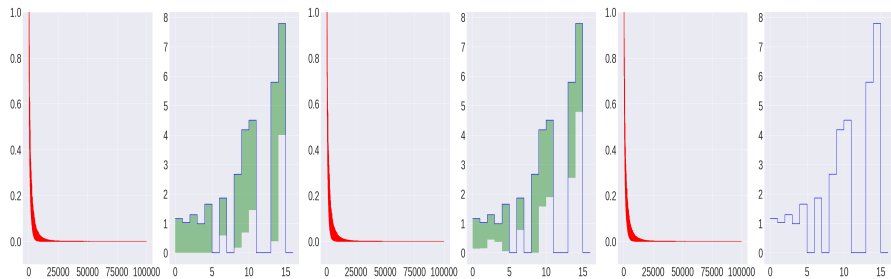
---

# Garnet



**Figure:** The difference between  $V_0^{\text{up},k}$  and  $V_0^{\pi}$ . X-axis represents states of Garnet environment. The first pair of pictures is computed with Q-learning policy on 1000 iterations, the second on 5000 iterations and the third group with RTDP algorithm.

# Frozen Lake



# Kernel-Based Reinforcement Learning

- The transition functions of KBRL's model,  $\hat{P}^a : \hat{\mathcal{S}} \times \hat{\mathcal{S}}_x \mapsto [0, 1]$ , are given by:

$$\hat{P}^a(\hat{s}_i^b | s) = \begin{cases} \kappa_\tau^a(s, s_i^b), & a = b \\ 0, & \text{else} \end{cases}$$

where  $\kappa_\tau^a(s, s_i^b)$  is normalized kernel function.



# Atari games

---

**Algorithm 2:** Approximate UVIP

---

**Input:** Sample  $(x_1, \dots, x_N)$ ;  $V^\pi, \tilde{V}_0^{\text{up}}, M_1, M_2, \gamma, \varepsilon$

**Result:**  $\hat{V}^{\text{up}}$

Generate  $r^a(X_i), Y_j^{x_i, a} \sim P^a(\cdot | x_i)$  for all  $i \in [N], j \in [M_1 + M_2], a \in \mathcal{A}$ ;

$k = 1$ ; **while**  $\|\tilde{V}_k^{\text{up}} - \tilde{V}_{k-1}^{\text{up}}\|_{x_N} > \varepsilon$  **do**

**for**  $a \in \mathcal{A}$  **do**

**for**  $i \in [N]$  **do**

**for**  $j \in [M_1 + M_2]$  **do**

$\tilde{V}_k^{\text{up}}(Y_j^{x_i, a}) = I[\tilde{V}_k^{\text{up}}](Y_j^{x_i, a})$  with  $I[\cdot](\cdot)$  defined in (6);

**end**

$$\bar{V}^{(i, a)} = M_1^{-1} \sum_{j=1}^{M_1} V^\pi(Y_j^{x_i, a});$$

**end**

**end**

**for**  $i \in [N]$  **do**

$$\tilde{V}_{k+1}^{\text{up}}(x_i) = M_2^{-1} \sum_{j=M_1+1}^{M_1+M_2} \max_{a \in \mathcal{A}} \{ r^a(x_i) + \gamma (\tilde{V}_k^{\text{up}}(Y_j^{x_i, a}) - V^\pi(Y_j^{x_i, a}) + \bar{V}^{(i, a)}) \};$$

**end**

$k = k + 1$ ;

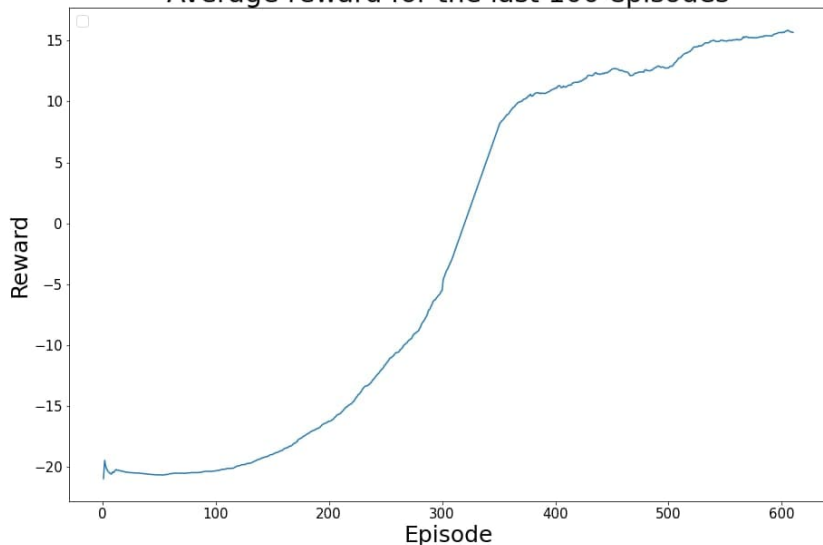
**end**

$\hat{V}^{\text{up}} = \tilde{V}_k^{\text{up}}$ .

---

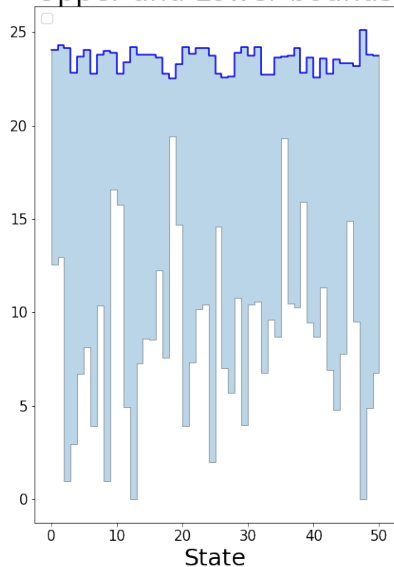
# Atari games

Average reward for the last 100 episodes

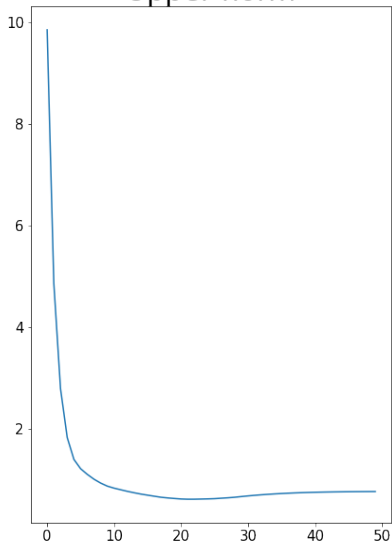


# Atari games

## Upper and Lower bounds



## Upper norm

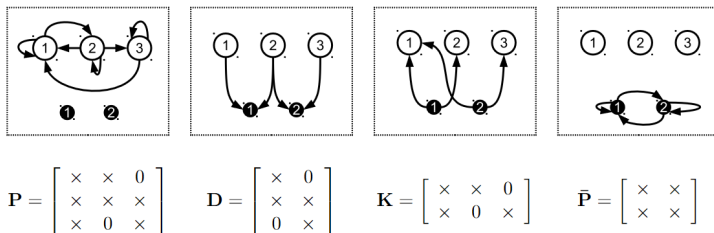


# Future work

- Merge modifications and run on more complicated environments (for example, Atari games)

Thank you for your attention!

# Stochastic-Factorization



**Figure:** Reducing the dimension of a transition model from  $n = 3$  states to  $m = 2$  artificial states. Original states  $s_i$  are represented as big white circles; small black circles depict artificial states  $\bar{s}_h$ . The symbol ‘ $\times$ ’ is used to represent nonzero elements.

# Kernel-Based Stochastic Factorization

---

**Algorithm 1** Batch KBSF

---

**Input:**  $S^a = \{(s_k^a, r_k^a, \hat{s}_k^a) | k = 1, 2, \dots, n_a\}$  for all  $a \in A$  ▷ Sample transitions  
 $\bar{S} = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_m\}$  ▷ Set of representative states  
**Output:**  $\tilde{\mathbf{v}} \approx \hat{\mathbf{v}}^*$   
**for each**  $a \in A$  **do**  
    Compute matrix  $\dot{\mathbf{D}}^a$ :  $\dot{d}_{ij}^a = \bar{\kappa}_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_j)$   
    Compute matrix  $\dot{\mathbf{K}}^a$ :  $\dot{k}_{ij}^a = \kappa_{\bar{\tau}}(\bar{s}_i, s_j^a)$   
    Compute vector  $\bar{\mathbf{r}}^a$ :  $\bar{r}_i^a = \sum_j \dot{k}_{ij}^a r_j^a$   
    Compute matrix  $\bar{\mathbf{P}}^a = \dot{\mathbf{K}}^a \dot{\mathbf{D}}^a$   
Solve  $\bar{M} \equiv (\bar{S}, A, \bar{\mathbf{P}}^a, \bar{\mathbf{r}}^a, \gamma)$  ▷ i.e., compute  $\bar{\mathbf{Q}}^*$   
Return  $\tilde{\mathbf{v}} = \Gamma \mathbf{D} \bar{\mathbf{Q}}^*$ , where  $\mathbf{D}^\top = [(\dot{\mathbf{D}}^1)^\top, (\dot{\mathbf{D}}^2)^\top, \dots, (\dot{\mathbf{D}}^{|A|})^\top]$

---