

# Bridging Distributional and Risk-Sensitive Reinforcement Learning: Balancing Statistical, Computational, and Risk Considerations

**Hao Liang**

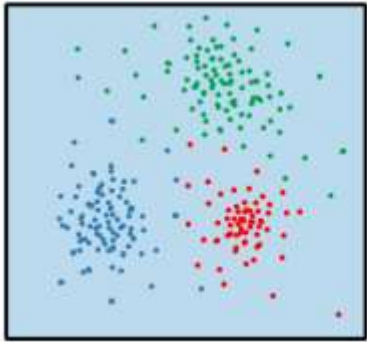
The Chinese University of Hong Kong, Shenzhen

- Hao Liang, and Zhi-Quan Luo. "**Bridging distributional and risk-sensitive reinforcement learning with provable regret bounds.**" arXiv preprint arXiv:2210.14051v3 (2022). Under review at *Journal of Machine Learning Research*.
- Hao Liang, and Zhi-Quan Luo. "**A distribution optimization framework for confidence bounds of risk measures.**" *International Conference on Machine Learning*. PMLR, 2023.

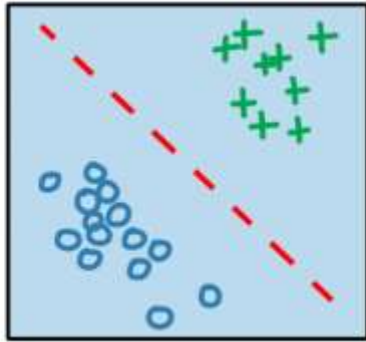
# What Is Reinforcement Learning (RL)?

machine learning

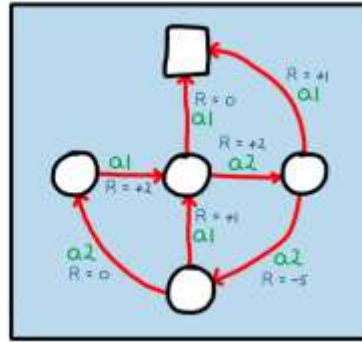
unsupervised learning



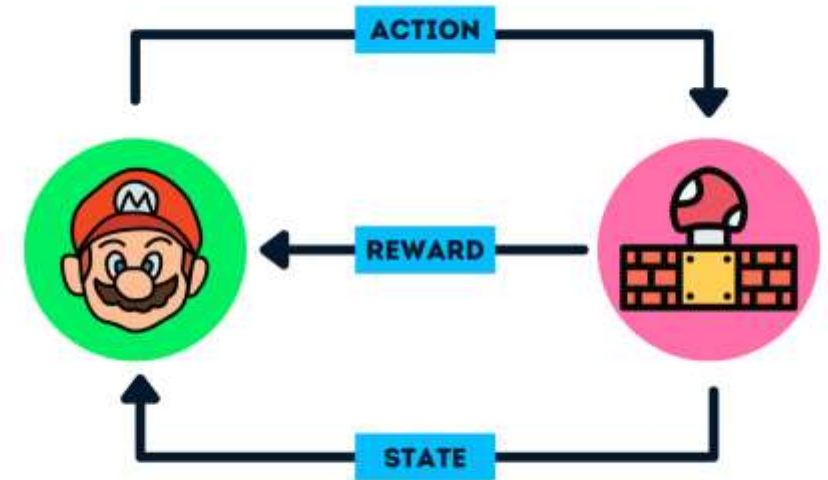
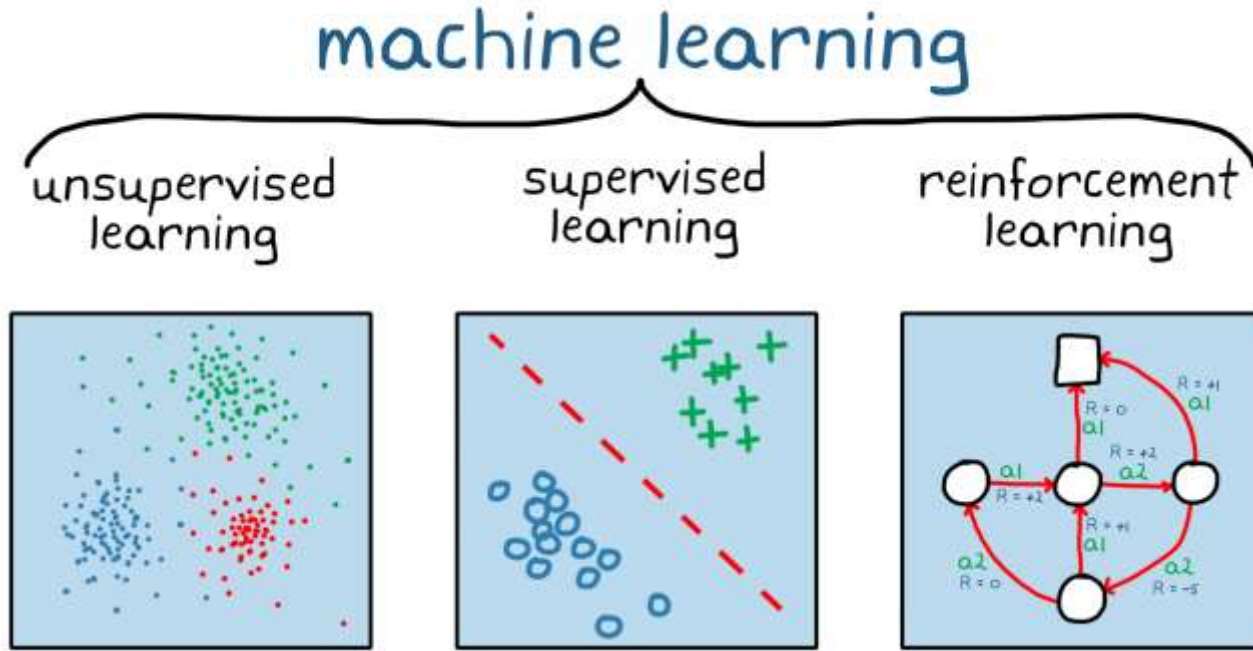
supervised learning



reinforcement learning



# What Is Reinforcement Learning (RL)?

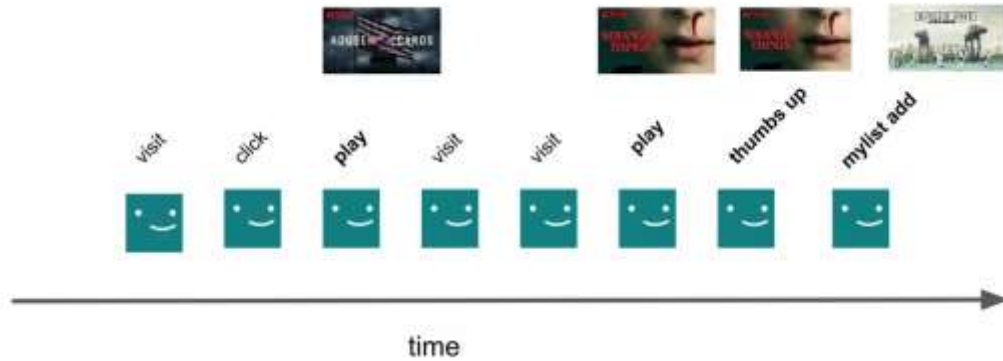


- **Sequential** decision-making under uncertainty
- The goal is to maximize **expected** cumulative rewards

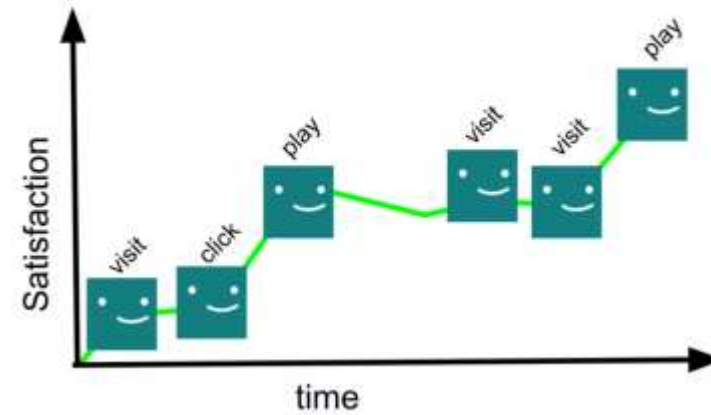
$$E \left[ \sum_t \text{Reward}_t \right]$$

# Standard RL is Risk-neutral

## Online recommendation



Maximize **average** “satisfaction”

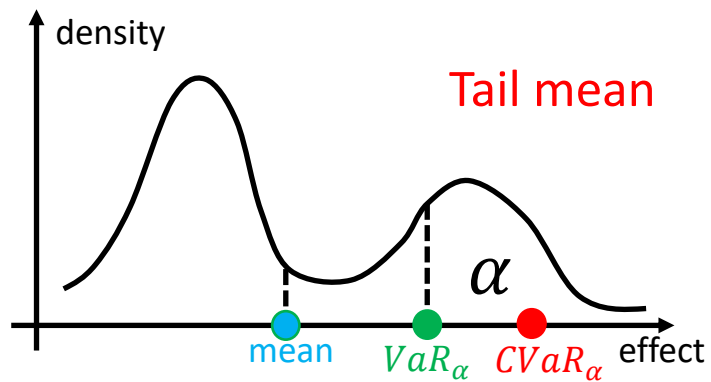
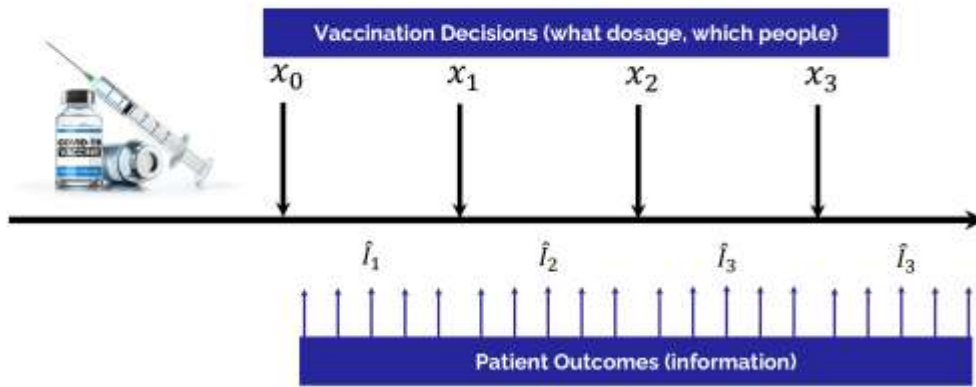


Standard RL aims to maximize the **expected** return

# Risk-sensitive RL (RSRL)

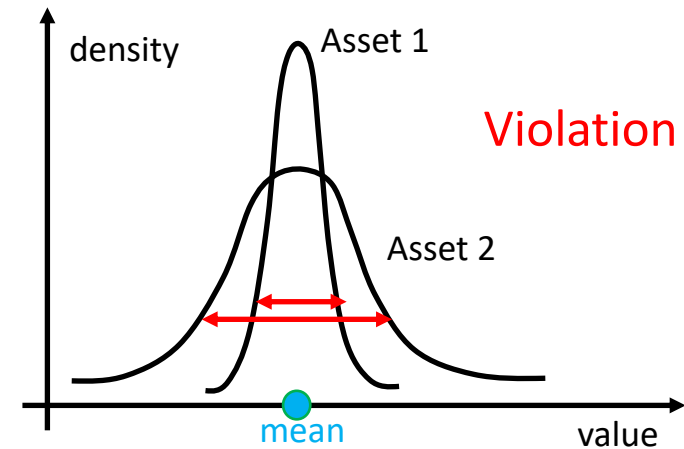
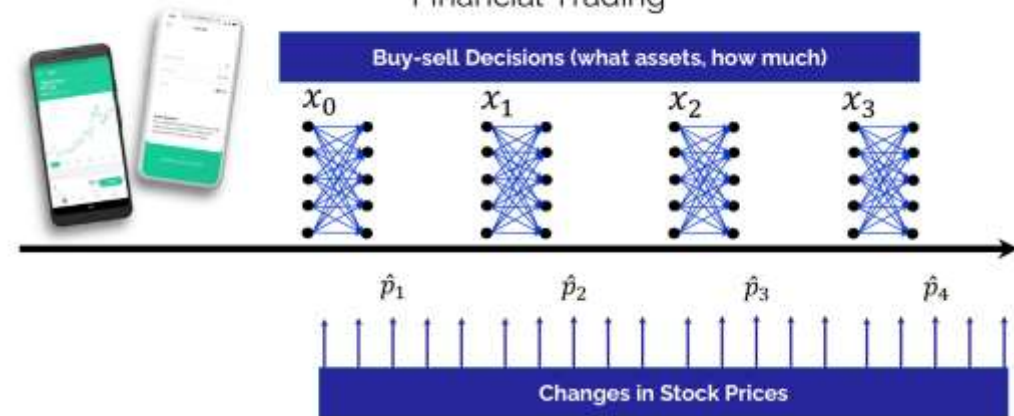
## Healthcare/Clinical trial

Testing new vaccines



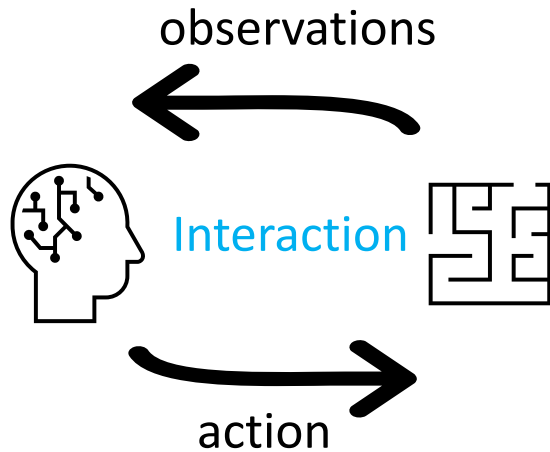
## Finance

Financial Trading



RSRL captures other **distributional** characteristics of return

# Towards **Sample-efficient** RSRL



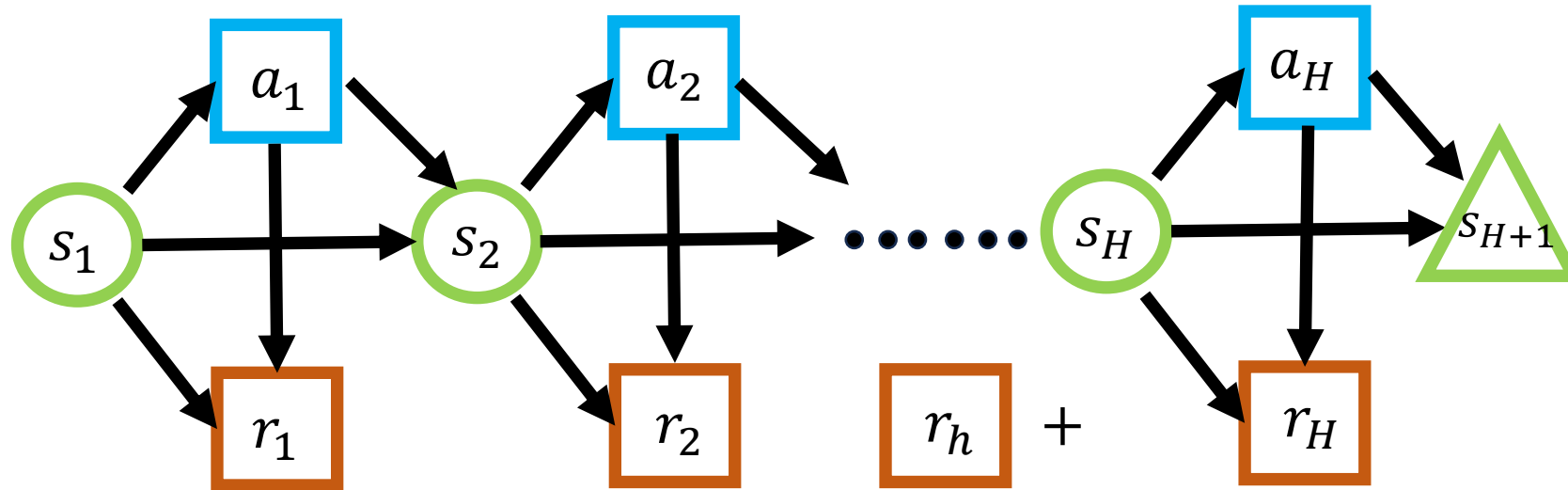
**Sample efficiency** is critical in RSRL!

- Healthcare monitoring systems
- Financial trading
- Online advertising

**Question: How to attain **sample-efficient** RSRL?**

**Principle: **Distributional Perspective****

# Markov Decision Process (MDP)



Tabular MDP  $\mathbf{M} = (S, A, P, r, H)$

■ Finite state space  $S$ , action space  $A$

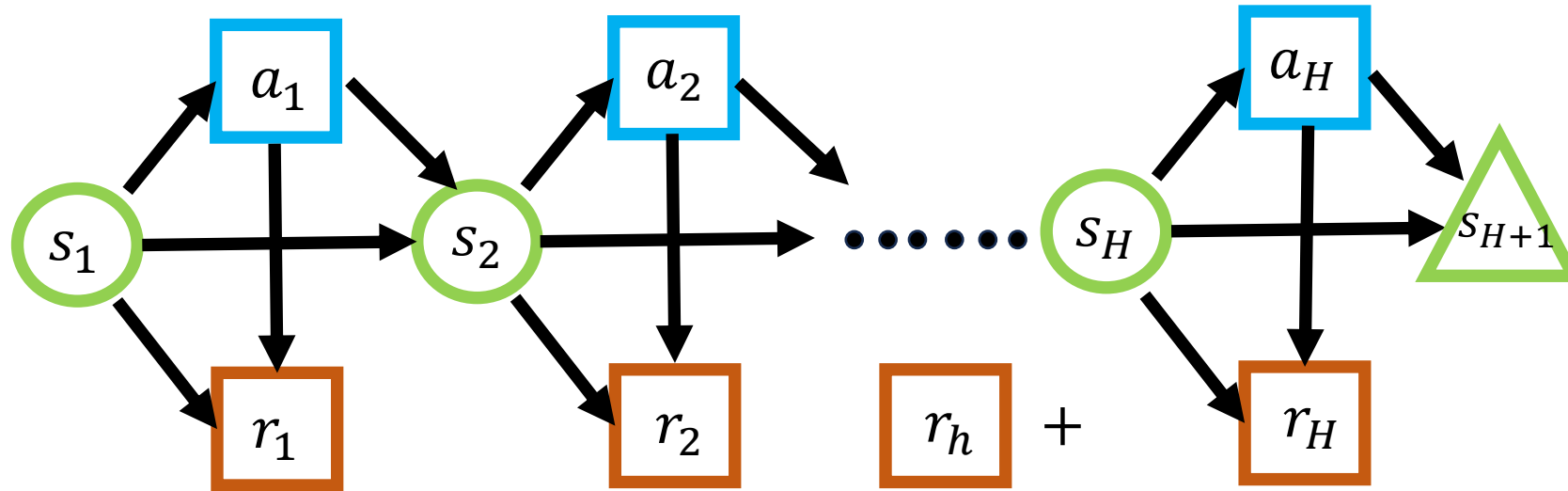
■ Transition kernel  $P_h(s, a)$

$$s' \sim P_h(s, a)$$

■ Reward function  $r_h(s, a)$

■ Horizon  $H$

# Markov Decision Process (MDP)



Tabular MDP  $\mathbf{M} = (S, A, P, r, H)$

■ Finite state space  $S$ , action space  $A$

■ Transition kernel  $P_h(s, a)$

$$s' \sim P_h(s, a)$$

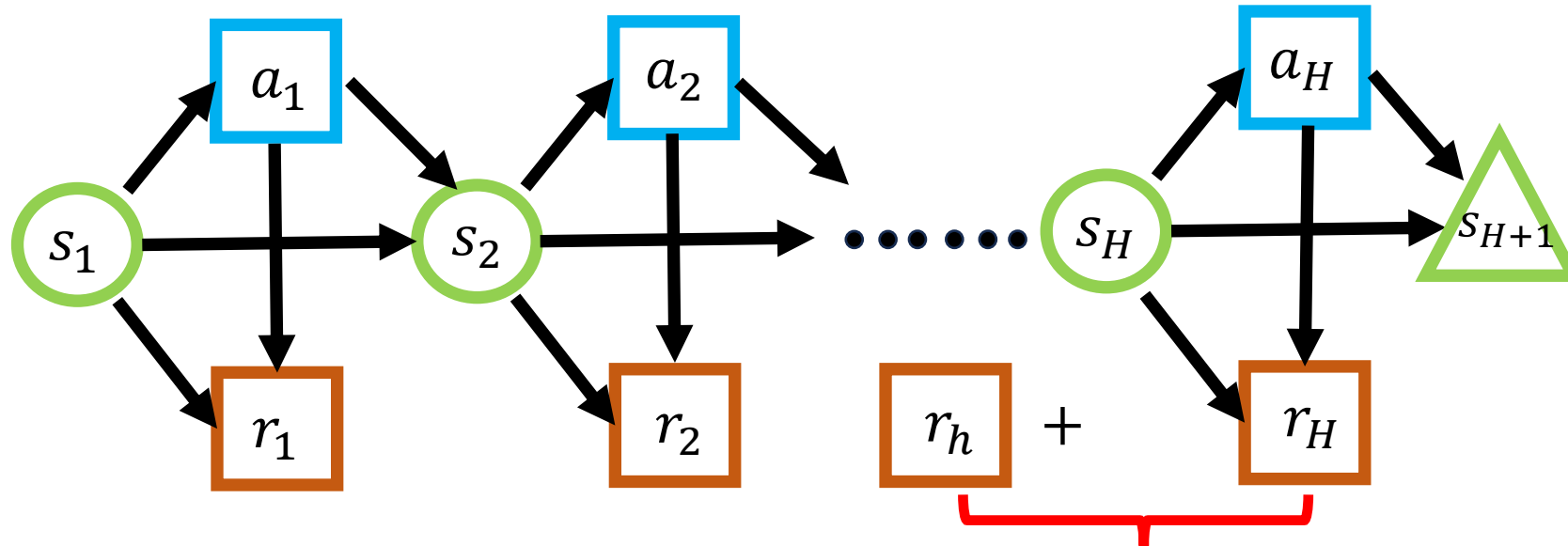
■ Reward function  $r_h(s, a)$

■ Horizon  $H$

■ Policy  $\pi = (\pi_h)_{h \in [H]}$   
 $\pi_h: S \rightarrow A \in \Pi$



# Markov Decision Process (MDP)



$Z_h^\pi$  Random Variable

Tabular MDP  $\mathbf{M} = (S, A, P, r, H)$

■ Finite state space  $S$ , action space  $A$

■ Transition kernel  $P_h(s, a)$

$$s' \sim P_h(s, a)$$

■ Reward function  $r_h(s, a)$

■ Horizon  $H$

■ Policy  $\pi = (\pi_h)_{h \in [H]}$

$$\Pi \ni \pi_h: S \rightarrow A$$

Return = cumulative reward

$$Z_h^\pi = r_h(s_h, a_h) + \dots + r_H(s_H, a_H),$$

$$a_h = \pi_h(s_h), s_{h+1} \sim P_h(s_h, a_h)$$

# Risk-neutral MDP vs. Risk-sensitive MDP

## Risk-neutral MDP

$$\max \mathbb{E}[Z_1^\pi]$$

## Risk-sensitive MDP

$$\max \rho(Z_1^\pi)$$

Risk measure  $\rho$ : R.V./distribution  $\rightarrow \mathbb{R}$   
reflects the risk preference towards the uncertainty

# Risk-neutral MDP vs. Risk-sensitive MDP

Risk-neutral MDP

$$\max \mathbf{E}[Z_1^\pi]$$

MANAGEMENT SCIENCE  
Vol. 18, No. 7, March, 1972  
*Printed in U.S.A.*

Risk-sensitive MDP

$$\max \rho(Z_1^\pi)$$

## RISK-SENSITIVE MARKOV DECISION PROCESSES\*

RONALD A. HOWARD† AND JAMES E. MATHESON‡§

Entropic risk measure (ERM) [1]

$$U_\beta(X) := \frac{1}{\beta} \log \mathbf{E}[\exp(\beta X)] = \mathbf{E}[X] + \frac{\beta}{2} \mathbf{V}[X] + O(|\beta|^2)$$

[1] Howard, Ronald A., and James E. Matheson. "Risk-sensitive Markov decision processes." *Management science* 18.7 (1972): 356-369.

# Risk-neutral MDP vs. Risk-sensitive MDP

Risk-neutral MDP

$$\max \mathbf{E}[Z_1^\pi]$$

Risk-sensitive MDP

$$\max \rho(Z_1^\pi)$$

MANAGEMENT SCIENCE  
Vol. 18, No. 7, March, 1972  
*Printed in U.S.A.*

## RISK-SENSITIVE MARKOV DECISION PROCESSES\*

RONALD A. HOWARD† AND JAMES E. MATHESON‡§

### Entropic risk measure (ERM) [1]

$$U_\beta(X) := \frac{1}{\beta} \log \mathbf{E}[\exp(\beta X)] = \mathbf{E}[X] + \frac{\beta}{2} \mathbf{V}[X] + O(|\beta|^2)$$

$\beta$  controls risk preference

- Risk-seeking  $\beta > 0$ : favoring high uncertainty
- Risk-averse  $\beta < 0$ : favoring low uncertainty
- Risk-neutral  $\beta \rightarrow 0$

[1] Howard, Ronald A., and James E. Matheson. "Risk-sensitive Markov decision processes." *Management science* 18.7 (1972): 356-369.

# Risk-sensitive MDP: **Optimality**

- Find the **optimal** policy  $\pi^* = (\pi_1^*, \dots, \pi_H^*)$

$$\pi^* := \operatorname{argmax}_{(\pi_1, \dots, \pi_H)} V_1^{\pi_1 \dots \pi_H}(s_1) := U_\beta(Z_1^{\pi_1 \dots \pi_H})$$

- A **multi-stage** maximization problem

$$\pi = (\pi_1, \dots, \pi_H) \in \Pi^H$$

- Direct search suffers **exponential** computational complexity [2]

$$|\Pi^H| = |\Pi|^H$$

# Risk-sensitive MDP: **Optimality**

- Find the **optimal** policy

$$\pi^* := \operatorname{argmax}_{(\pi_1, \dots, \pi_H)} V_1^{\pi_1 \dots \pi_H}(s_1) = \mathbf{U}_\beta(Z_1^{\pi_1 \dots \pi_H})$$

- A **multi-stage** maximization problem

$$\pi = (\pi_1, \dots, \pi_H) \in \Pi^H$$

- Direct search suffers **exponential** computational complexity

$$|\Pi^H| = |\Pi|^H$$

## **Risk-neutral optimality equation**

$$Q_h^*(s, a) = r_h(s, a) + \sum P_h(s' | s, a) V_{h+1}^*(s')$$
$$V_h^*(s) = \max_a Q_h^*(s, a), V_{H+1}^*(s) = 0$$

# Risk-sensitive MDP: **Optimality**

- Find the **optimal** policy

$$\pi^* := \operatorname{argmax}_{(\pi_1, \dots, \pi_H)} V_1^{\pi_1 \dots \pi_H}(s_1) = U_{\beta}(Z_1^{\pi_1 \dots \pi_H})$$

- A **multi-stage** maximization problem

$$\pi = (\pi_1, \dots, \pi_H) \in \Pi^H$$

- Direct search suffers **exponential** computational complexity

$$|\Pi^H| = |\Pi|^H$$

## **Risk-neutral optimality equation**

$$Q_h^*(s, a) = r_h(s, a) + \sum P_h(s' | s, a) V_{h+1}^*(s')$$
$$V_h^*(s) = \max_a Q_h^*(s, a), V_{H+1}^*(s) = 0$$



## **Optimal substructure**

- Break into **multiple single-stage** problems
- Recursion of value functions

# Risk-sensitive MDP: **Optimality**

- Find the **optimal** policy

$$\pi^* := \operatorname{argmax}_{(\pi_1, \dots, \pi_H)} V_1^{\pi_1 \dots \pi_H}(s_1) = \mathbf{U}_\beta(Z_1^{\pi_1 \dots \pi_H})$$

- A **multi-stage** maximization problem

$$\pi = (\pi_1, \dots, \pi_H) \in \Pi^H$$

- Direct search suffers **exponential** computational complexity

$$|\Pi^H| = |\Pi|^H$$

**Risk-neutral optimality equation**



**Optimal substructure**

$$Q_h^*(s, a) = r_h(s, a) + \sum P_h(s' | s, a) V_{h+1}^*(s')$$
$$V_h^*(s) = \max_a Q_h^*(s, a), V_{H+1}^*(s) = 0$$

- Break into **multiple single-stage** problems
- Recursion of value functions

**Optimal substructure for risk-sensitive MDP?**



# Risk-sensitive MDP: **Optimality**

- Find the **optimal** policy

$$\pi^* := \operatorname{argmax}_{(\pi_1, \dots, \pi_H)} V_1^{\pi_1 \dots \pi_H}(s_1) = \mathbf{U}_\beta(Z_1^{\pi_1 \dots \pi_H})$$

- A **multi-stage** maximization problem

$$\pi = (\pi_1, \dots, \pi_H) \in \Pi^H$$

- Direct search suffers **exponential** computational complexity

$$|\Pi^H| = |\Pi|^H$$

**Risk-neutral optimality equation**



**Optimal substructure**

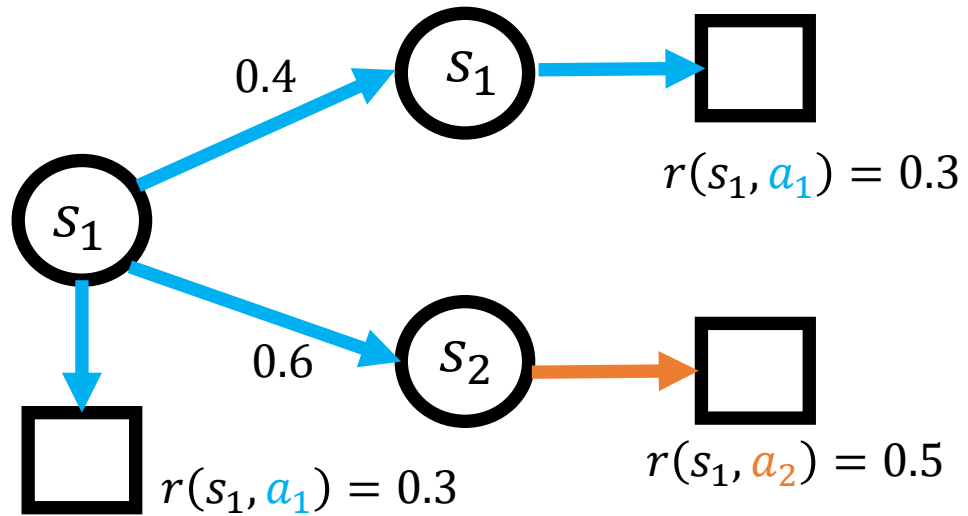
$$Q_h^*(s, a) = r_h(s, a) + \sum P_h(s' | s, a) V_{h+1}^*(s')$$
$$V_h^*(s) = \max_a Q_h^*(s, a), V_{H+1}^*(s) = 0$$

- Break into **multiple single-stage** problems
- Recursion of value functions

**Optimal substructure for risk-sensitive MDP?**

**Yes. Distributional dynamic programming**

# First Attempt: **Distributional** Policy Evaluation



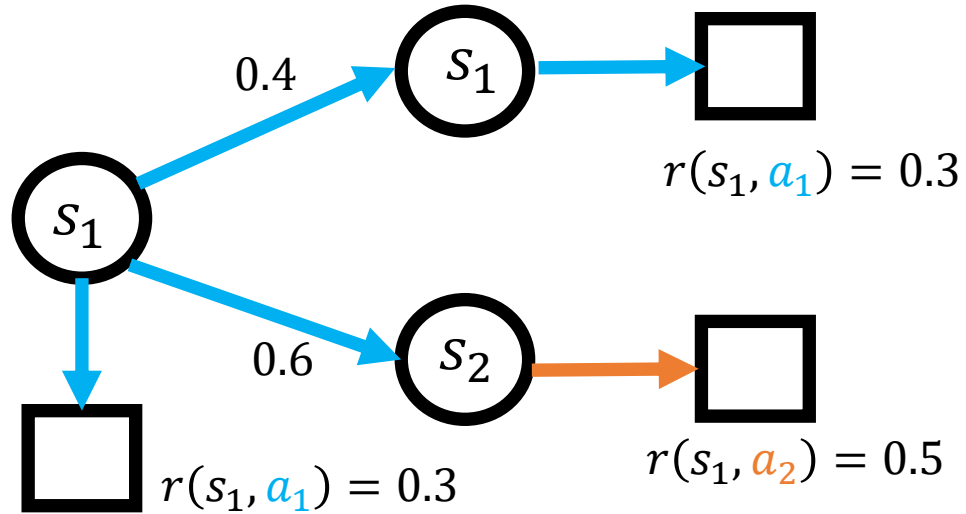
Given policy  $\pi$

$$\begin{aligned}\pi_1(s_1) &= a_1, \\ \pi_2(s_1) &= a_1, \pi_2(s_2) = a_2\end{aligned}$$

**Task**

Determine the r.v.  $Z_1^\pi(s_1)$

# First Attempt: **Distributional** Policy Evaluation



Given policy  $\pi$

$$\pi_1(s_1) = a_1, \\ \pi_2(s_1) = a_1, \pi_2(s_2) = a_2$$

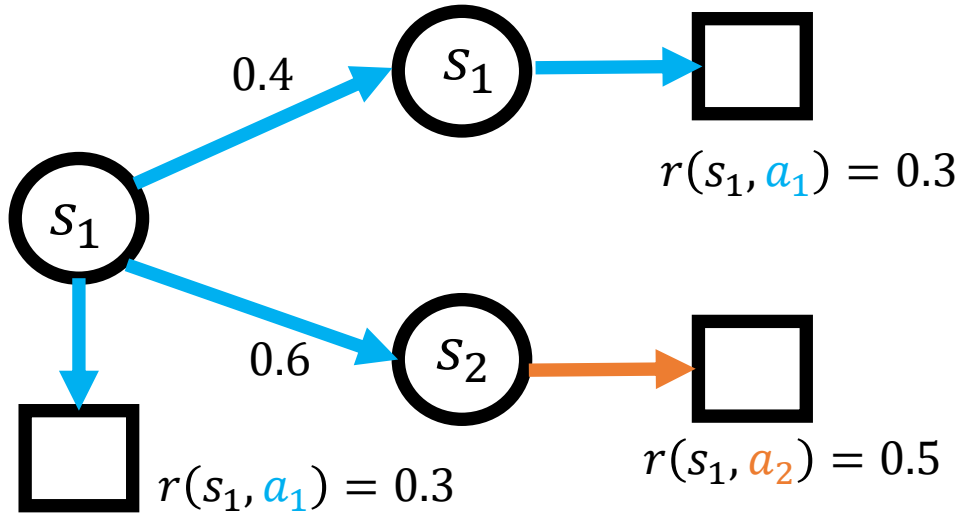
**Task**

Determine the r.v.  $Z_1^\pi(s_1)$

## Recursion of random variables

Denote by  $Y_h(s) := Z_h(s, \pi_h(s))$

# First Attempt: **Distributional** Policy Evaluation



Given policy  $\pi$   
 $\pi_1(s_1) = a_1,$   
 $\pi_2(s_1) = a_1, \pi_2(s_2) = a_2$

**Task**  
Determine the r.v.  $Z_1^\pi(s_1)$

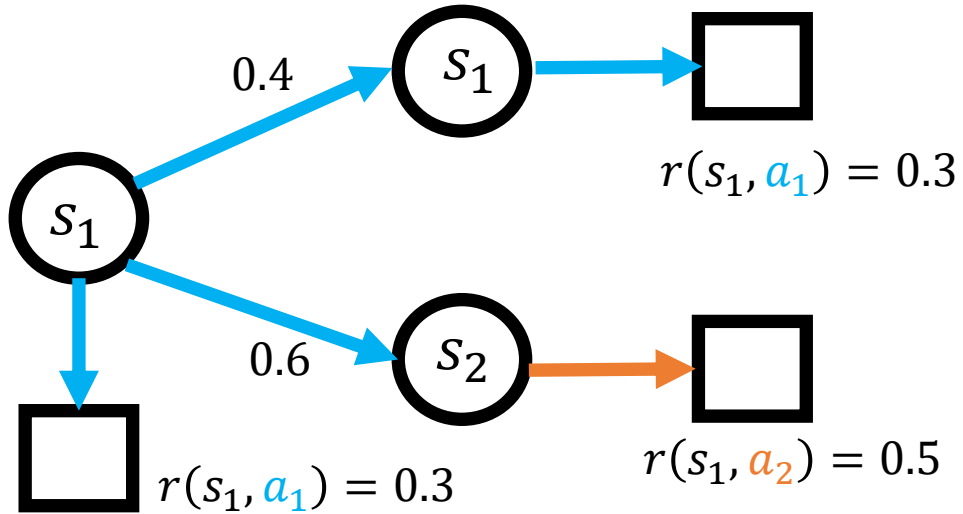
## Recursion of random variables

$$\begin{aligned} Y_2(s_1) &= 0.3 \\ Y_2(s_2) &= 0.5 \end{aligned}$$

Deterministic final reward

$$s' | s_1, a_1 \sim \begin{cases} s_1, w.p. 0.4 \\ s_2, w.p. 0.6 \end{cases}$$

# First Attempt: **Distributional** Policy Evaluation



Given policy  $\pi$   
 $\pi_1(s_1) = a_1,$   
 $\pi_2(s_1) = a_1, \pi_2(s_2) = a_2$

**Task**  
Determine the r.v.  $Z_1^\pi(s_1)$

## Recursion of random variables

$$Y_2(s_1) = 0.3$$
$$Y_2(s_2) = 0.5$$

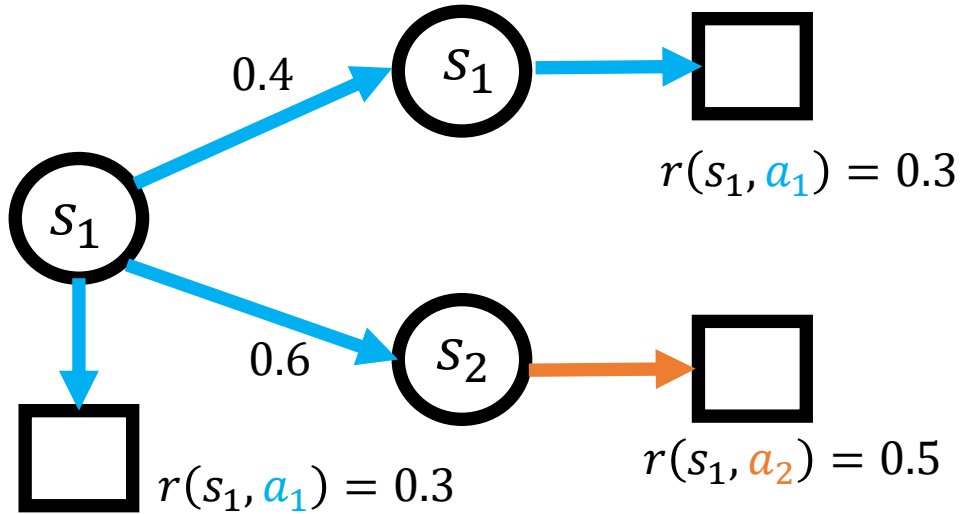


$$Y_2(S')|s_1, a_1 \sim \begin{cases} 0.3, w.p. 0.4 \\ 0.5, w.p. 0.6 \end{cases}$$

$$S'|s_1, a_1 \sim \begin{cases} s_1, w.p. 0.4 \\ s_2, w.p. 0.6 \end{cases}$$

Random future return

# First Attempt: **Distributional** Policy Evaluation



Given policy  $\pi$   
 $\pi_1(s_1) = a_1,$   
 $\pi_2(s_1) = a_1, \pi_2(s_2) = a_2$

**Task**  
 Determine the r.v.  $Z_1^\pi(s_1)$

## Recursion of random variables

$$\begin{aligned} Y_2(s_1) &= 0.3 \\ Y_2(s_2) &= 0.5 \end{aligned}$$

$$S' | s_1, a_1 \sim \begin{cases} s_1, w.p. 0.4 \\ s_2, w.p. 0.6 \end{cases}$$



$$Y_2(S') | s_1, a_1 \sim \begin{cases} 0.3, w.p. 0.4 \\ 0.5, w.p. 0.6 \end{cases}$$

$$\begin{aligned} Y_1(s_1) &= r_1(s_1, a_1) + Y_2(S') \\ &\sim \begin{cases} 0.3 + 0.3, w.p. 0.4 \\ 0.3 + 0.5, w.p. 0.6 \end{cases} \end{aligned}$$

Random initial return

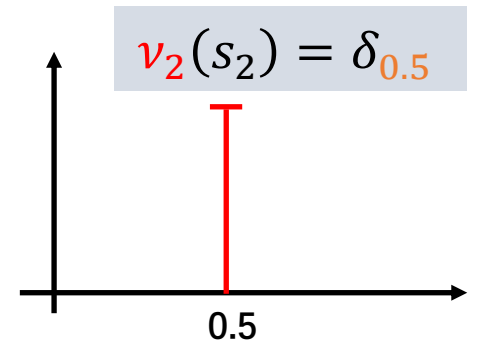
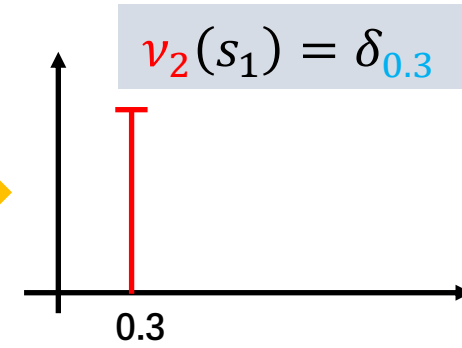
# Recursion of **Distributions**

## Recursion of r.v.s

$$Y_2(s_1) = 0.3$$

$$Y_2(s_2) = 0.5$$

Dirac



Denote by  $\nu_h(s)$  the distribution of  $Y_h(s)$

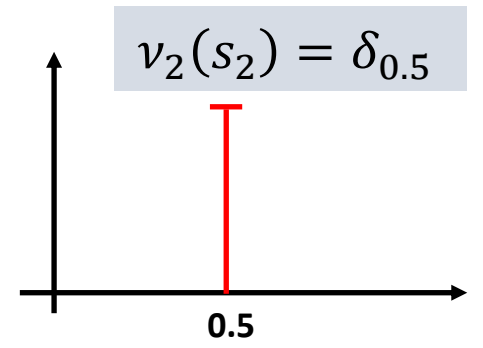
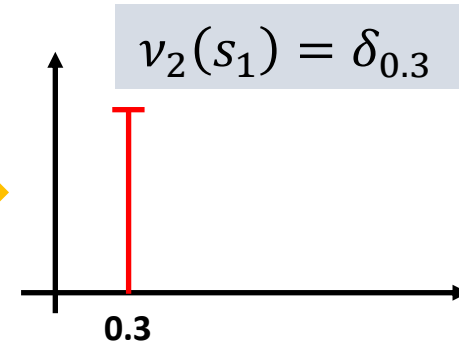
$$Y_h(s) \sim \nu_h(s)$$

# Recursion of Distributions

## Recursion of r.v.s

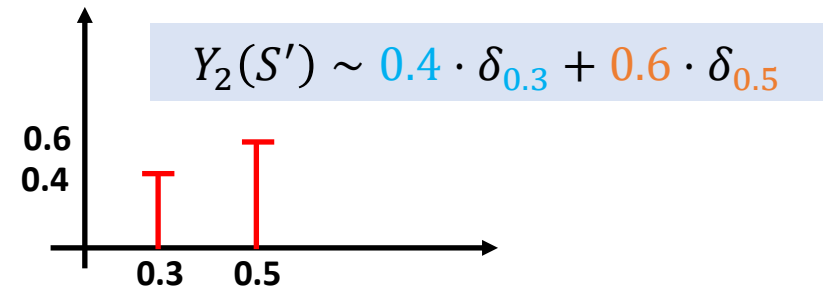
$$\begin{aligned} Y_2(s_1) &= 0.3 \\ Y_2(s_2) &= 0.5 \end{aligned}$$

Dirac



$$Y_2(S')|s_1, a_1 \sim \begin{cases} 0.3, w.p. 0.4 \\ 0.5, w.p. 0.6 \end{cases}$$

Mixture



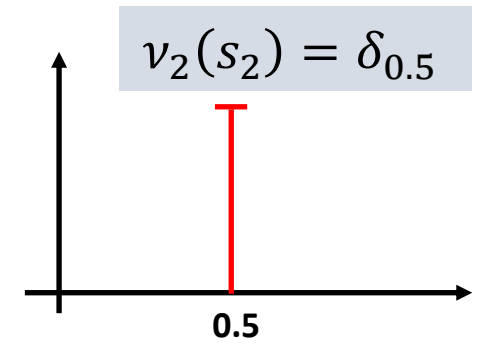
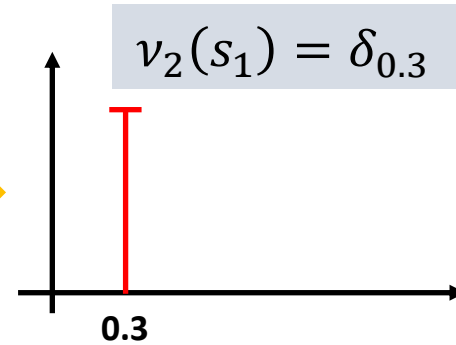


# Recursion of Distributions

## Recursion of r.v.s

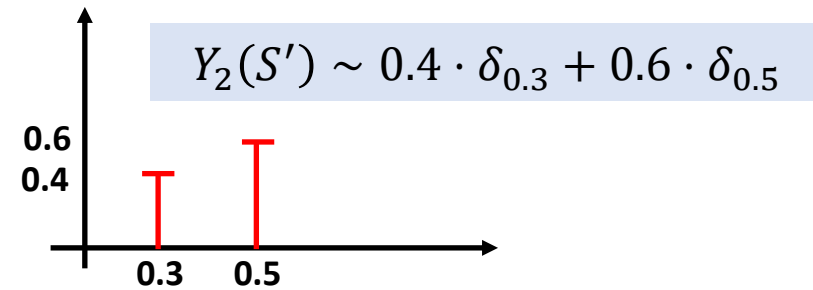
$$\begin{aligned} Y_2(s_1) &= 0.3 \\ Y_2(s_2) &= 0.5 \end{aligned}$$

Dirac



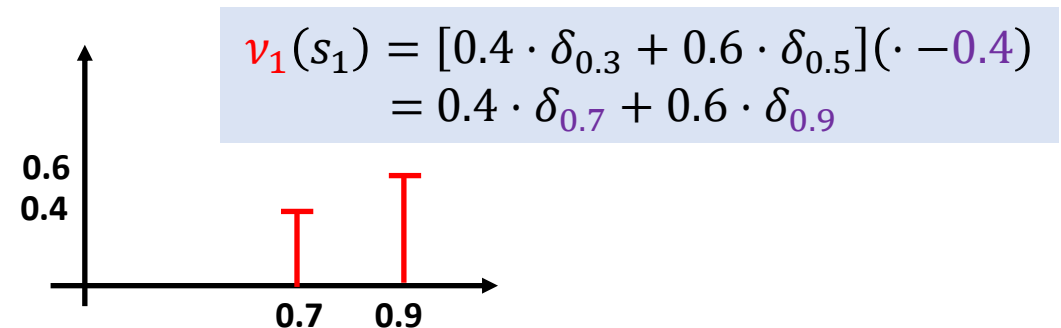
$$Y_2(S') | s_1, a_1 \sim \begin{cases} 0.3, w.p. 0.4 \\ 0.5, w.p. 0.6 \end{cases}$$

Mixture



$$\begin{aligned} Y_1(s_1) &= r_1(s_1, a_1) + Y_2(S') \\ &\sim \begin{cases} 0.4 + 0.3, w.p. 0.4 \\ 0.4 + 0.5, w.p. 0.6 \end{cases} \end{aligned}$$

Shift



# Distributional Dynamic Programming: Policy Evaluation

Fact 1: mixture distribution

$$X_i \sim F_i, P(I = i) = p_i \implies X_I \sim \sum_i p_i F_i$$

Fact 2: shift

$$X \sim F \implies X + c \sim F(\cdot - c)$$

# Distributional Dynamic Programming: Policy Evaluation

Fact 1: mixture distribution

$$X_i \sim F_i, P(I = i) = p_i \implies X_I \sim \sum_i p_i F_i$$

Fact 2: shift

$$X \sim F \implies X + c \sim F(\cdot - c)$$

Recursion of R.V.s

$$Z_h(s, a) = r_h(s, a) + Y_{h+1}(S')$$

$$S' \sim P_h(\cdot | s, a)$$

$$Y_h(s) = Z_h(s, \pi_h(s))$$

# Distributional Dynamic Programming: Policy Evaluation

Fact 1: mixture distribution

$$X_i \sim F_i, P(I = i) = p_i \Rightarrow X_I \sim \sum_i p_i F_i$$

Fact 2: shift

$$X \sim F \Rightarrow X + c \sim F(\cdot - c)$$

Recursion of R.V.s

$$\begin{aligned} Z_h(s, a) &= r_h(s, a) + Y_{h+1}(S') \\ S' &\sim P_h(\cdot | s, a) \\ Y_h(s) &= Z_h(s, \pi_h(s)) \end{aligned}$$

Two facts



Recursion of distributions

$$\begin{aligned} \eta_h(s, a) &= \sum P_h(s' | s, a) v_{h+1}(s') (\cdot - r_h(s, a)) \\ v_h(s) &= \eta_h(s, \pi_h(s)) \end{aligned}$$

mixture

shift

# Distributional Dynamic Programming: Policy Evaluation

**Fact 1:** mixture distribution

$$X_i \sim F_i, P(I = i) = p_i \implies X_I \sim \sum_i p_i F_i$$

**Fact 2:** shift

$$X \sim F \implies X + c \sim F(\cdot - c)$$

Recursion of R.V.s

$$\begin{aligned} Z_h(s, a) &= r_h(s, a) + Y_{h+1}(S') \\ S' &\sim P_h(\cdot | s, a) \\ Y_h(s) &= Z_h(s, \pi_h(s)) \end{aligned}$$

Two facts



Recursion of distributions

$$\begin{aligned} \eta_h(s, a) &= \sum P_h(s' | s, a) v_{h+1}(s') (\cdot - r_h(s, a)) \\ v_h(s) &= \eta_h(s, \pi_h(s)) \end{aligned}$$

**Distributional Bellman Operator**  $\mathbf{T}_d: P(R)^S \rightarrow P(R)^{S \times A}$

$$\begin{aligned} \eta_h(s, a) &= [\mathbf{T}_d v_{h+1}](s, a) \\ &:= \sum P_h(s' | s, a) v_{h+1}(s') (\cdot - r_h(s, a)) \end{aligned}$$

# Distributional Dynamic Programming: Risk-sensitive Control

## Policy Evaluation

Given  $\pi$ , determine  $Z_1^\pi$

## Risk-sensitive Control

$$\max_{\pi} U_{\beta} (Z_1^{\pi})$$

# Distributional Dynamic Programming: Risk-sensitive Control

## Risk-sensitive Control

$$\max_{\pi} U_{\beta}(Z_1^{\pi})$$

Key property 1: Additivity

$$U_{\beta}(X + c) = U_{\beta}(X) + c$$

Key property 2: Independence

$$U_{\beta}(F_2) \leq U_{\beta}(F_1) \implies U_{\beta}((1 - \theta)F_2 + \theta G) \leq U_{\beta}((1 - \theta)F_1 + \theta G)$$

# Distributional Dynamic Programming: Risk-sensitive Control

## Risk-sensitive Control

$$\max_{\pi} U_{\beta}(Z_1^{\pi})$$

Key property 1: Additivity

$$U_{\beta}(X + c) = U_{\beta}(X) + c$$

Key property 2: Independence

$$U_{\beta}(F_2) \leq U_{\beta}(F_1) \implies U_{\beta}((1 - \theta)F_2 + \theta G) \leq U_{\beta}((1 - \theta)F_1 + \theta G)$$

## Distributional Bellman Optimality Equation

$$\begin{aligned} \eta_h^*(s, a) &= [\mathbf{T}_d v_{h+1}^*](s, a) \\ \pi_h^*(s) &= \operatorname{argmax}_a U_{\beta}(\eta_h^*(s, a)) \\ v_h^*(s) &= \eta_h^*(s, \pi_h^*(s)) \end{aligned}$$

backward recursion



# Distributional Dynamic Programming: Risk-sensitive Control

## Risk-sensitive Control

$$\max_{\pi} U_{\beta}(Z_1^{\pi})$$

Key property 1: Additivity

$$U_{\beta}(X + c) = U_{\beta}(X) + c$$

Key property 2: Independence

$$U_{\beta}(F_2) \leq U_{\beta}(F_1) \implies U_{\beta}((1 - \theta)F_2 + \theta G) \leq U_{\beta}((1 - \theta)F_1 + \theta G)$$

## Distributional Bellman Optimality Equation

$$\eta_h^*(s, a) = [\mathbf{T}_d v_{h+1}^*](s, a)$$

$$\pi_h^*(s) = \operatorname{argmax}_a U_{\beta}(\eta_h^*(s, a))$$

$$v_h^*(s) = \eta_h^*(s, \pi_h^*(s))$$

backward recursion

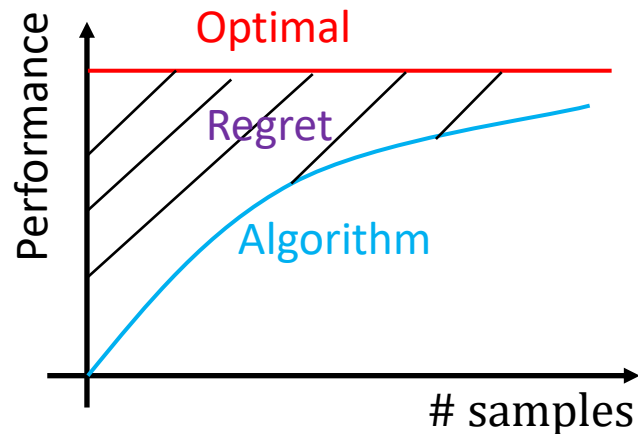
greedy is optimal

# From DP to RL

- DP requires the knowledge of transition
- RL deals with **unknown** transition
- Learn from interactions/**samples**
- Expensive samples



**Fundamental problem of RL**  
**Statistical/sample efficiency**

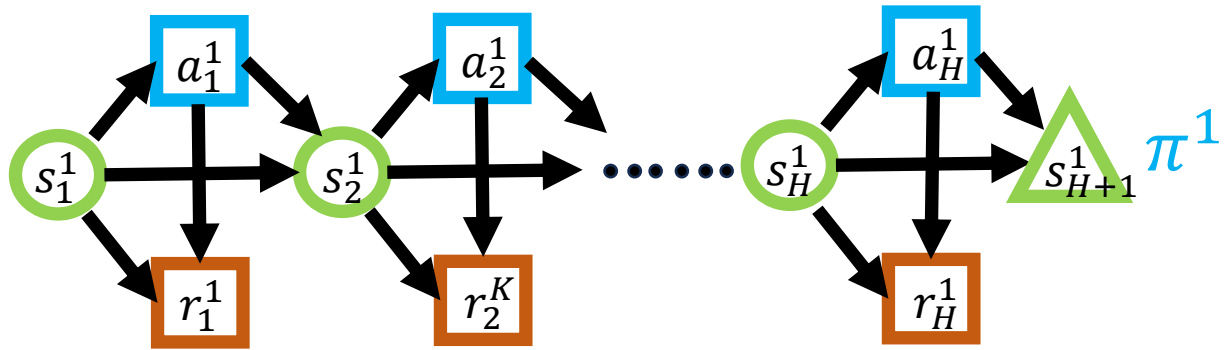


**Low regret**



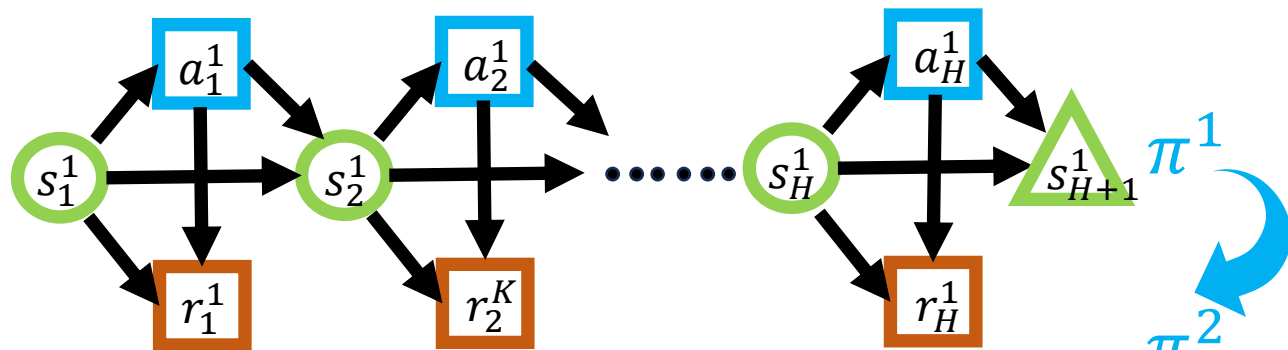
**High sample efficiency**

# Episodic MDP

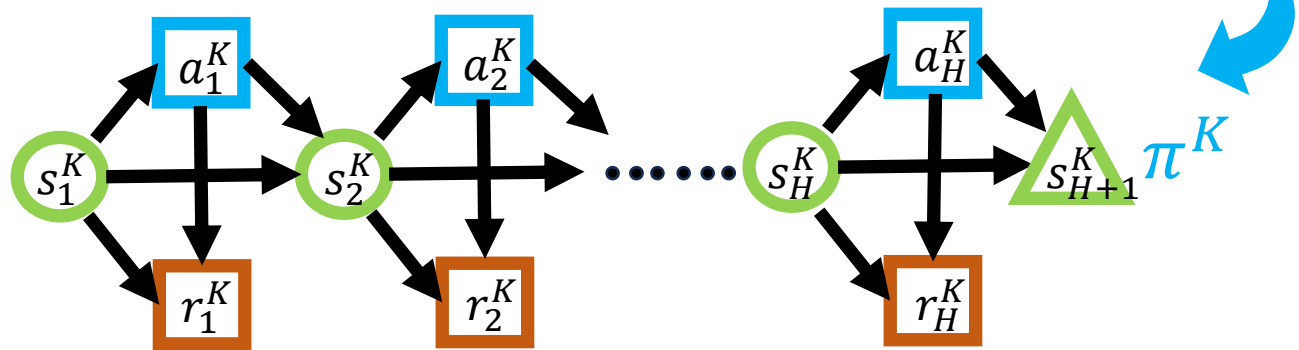


First episode

# Episodic MDP

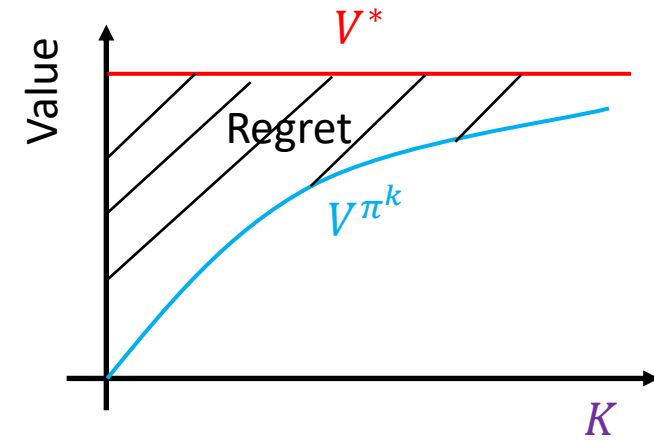
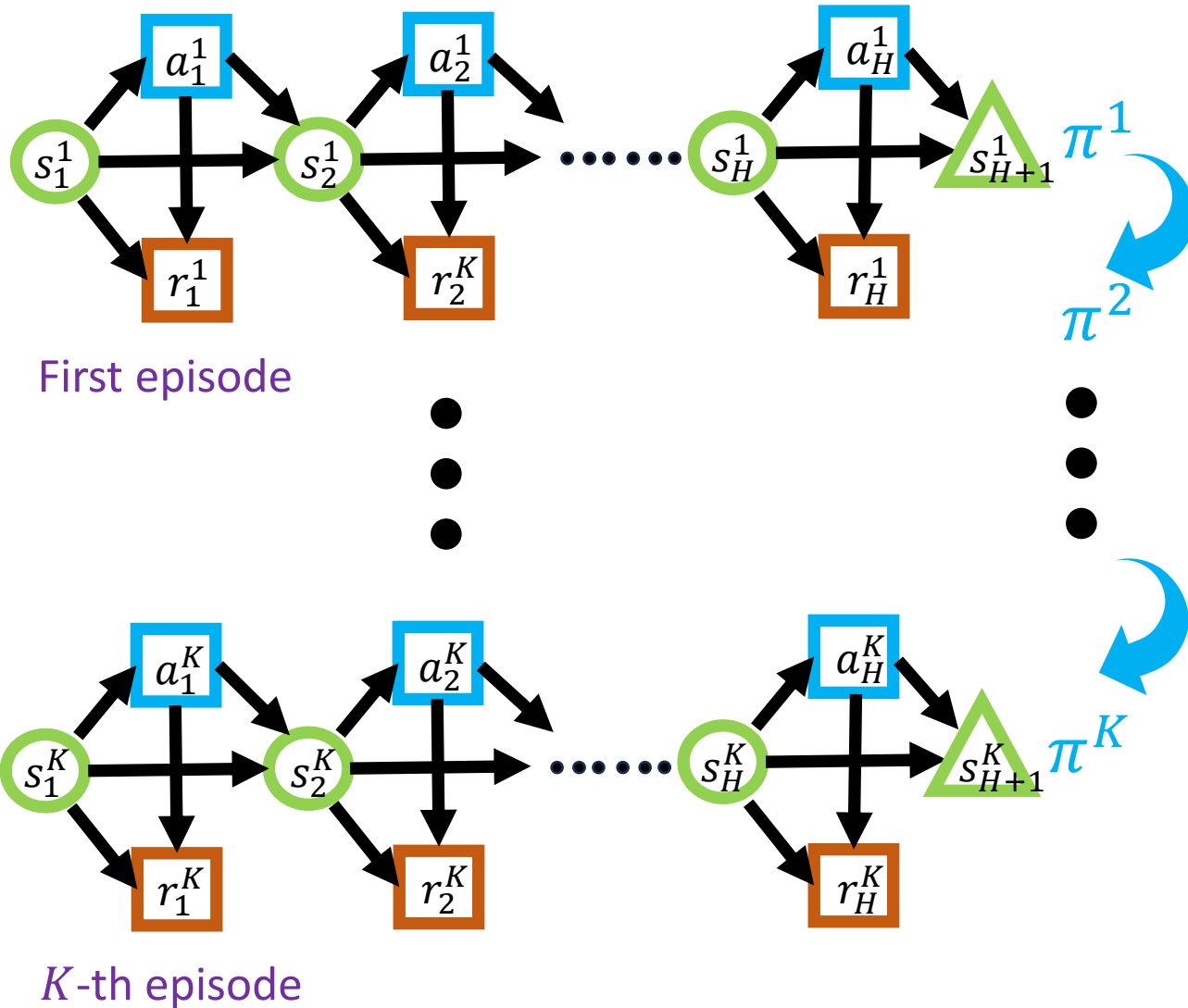


First episode



$K$ -th episode

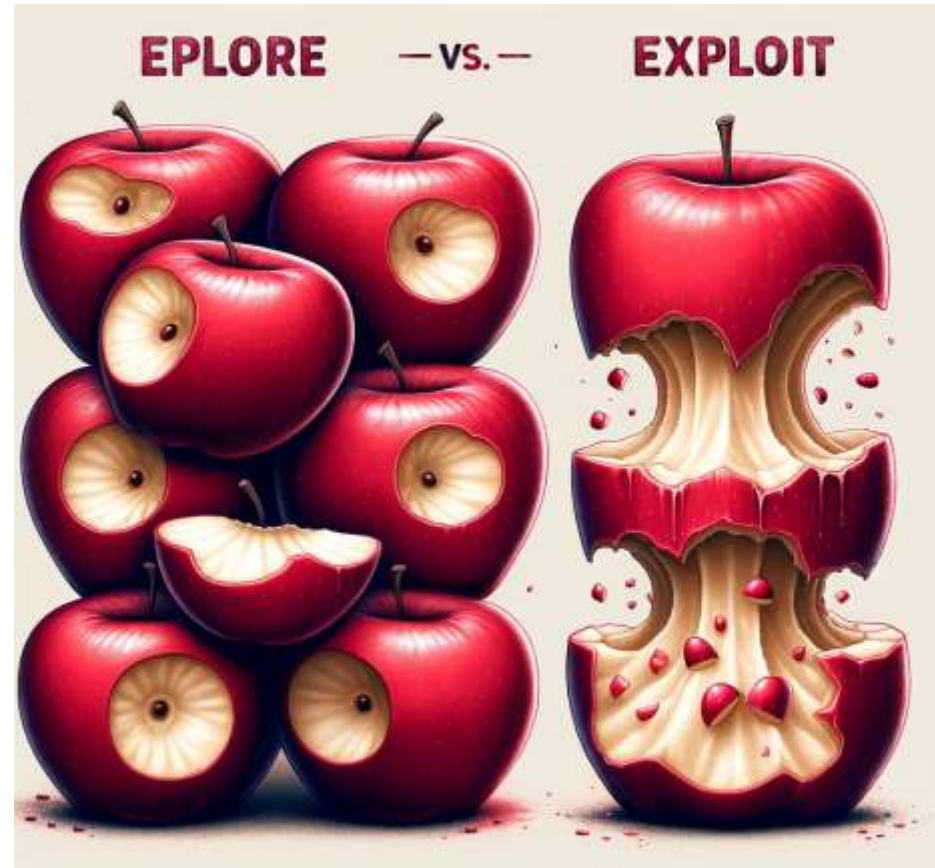
# Episodic MDP



Goal : minimize the regret over  $K$  episodes

$$\text{Regret}(\text{alg}, K) := \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k)$$

# Exploration vs. Exploitation Dilemma



## EXPLORATION

Try different actions  
to gather info

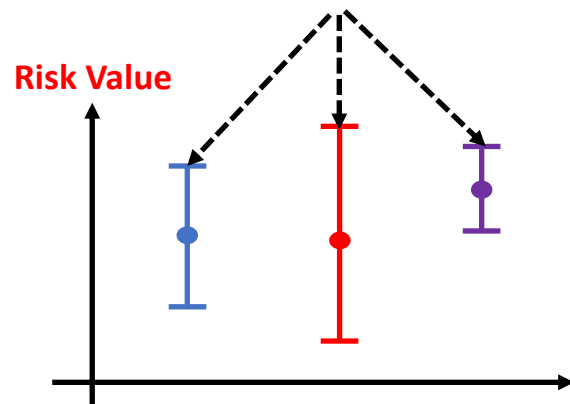
## EXPLOITATION

Use the best-known action

**Sample efficiency** requires balance between exploration and exploitation!

# Risk-sensitive Optimism in Face of Uncertainty

- An effective principle for efficient exploration  
Optimism in Face of Uncertainty (OFU)
- Act greedily w.r.t. **Upper Confidence Bound** of **Risk Value**
- Encourages exploring actions with the **best possible** outcomes



**Tighter UCB**



**Higher sample efficiency**

# Risk-sensitive **O**ptimistic **D**istribution **I**teration (**RODI**)

**Empirical model**

$$\hat{P}_h^k(s'|s, a) = N_h^k(s, a, s') / N_h^k(s, a)$$

**Approximate Distributional Bellman operator**

$$\left[ \hat{T}_d^k v_{h+1}^k \right] (s, a) := \sum \hat{P}_h^k(s'|s, a) v_{h+1}^k(s') (\cdot - r_h(s, a))$$

**Approximate Bellman recursion**

$$\eta_h^k \leftarrow \hat{T}_d^k v_{h+1}^k$$

**Distributional Optimism Operator**

$$\eta_h^k \leftarrow \mathbf{O}_{c^k} \eta_h^k$$

**Policy in episode  $k$**

$$\pi_h^k(s) \leftarrow \operatorname{argmax}_a U_\beta(\eta_h^k(s, a))$$



# Risk-sensitive **O**ptimistic **D**istribution **I**teration (**RODI**)

Empirical model

$$\hat{P}_h^k(s'|s, a) = N_h^k(s, a, s') / N_h^k(s, a)$$

Approximate Distributional Bellman operator

$$\left[ \hat{T}_d^k v_{h+1}^k \right] (s, a) := \sum \hat{P}_h^k(s'|s, a) v_{h+1}^k(s') (\cdot - r_h(s, a))$$

Approximate Bellman recursion

$$\eta_h^k \leftarrow \hat{T}_d^k v_{h+1}^k$$

Distributional Optimism Operator

$$\eta_h^k \leftarrow \mathbf{O}_{c^k} \eta_h^k$$

Policy in episode  $k$

$$\pi_h^k(s) \leftarrow \operatorname{argmax}_a U_\beta(\eta_h^k(s, a))$$

**RODI** in one line

$$\eta_h^k \leftarrow \mathbf{O}_{c^k} \hat{T}_d^k v_{h+1}^k$$

**Optimism**

$$U_\beta(\eta_h^k(s, a)) \geq U_\beta(\eta_h^*(s, a)) \\ \forall (s, a, k, h)$$

# Distributional Optimism

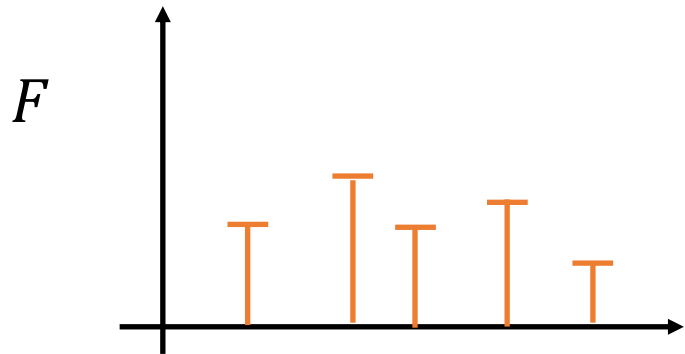
**Def.** A distribution  $F$  is **optimistic** over  $G$  if

$$U_{\beta}(F) \geq U_{\beta}(G)$$

Denote by  $F \geq G$

Given a confidence ball of distributions

$$B(F, c) := \{G \mid \|G - F\| \leq c\}$$



# Distributional Optimism

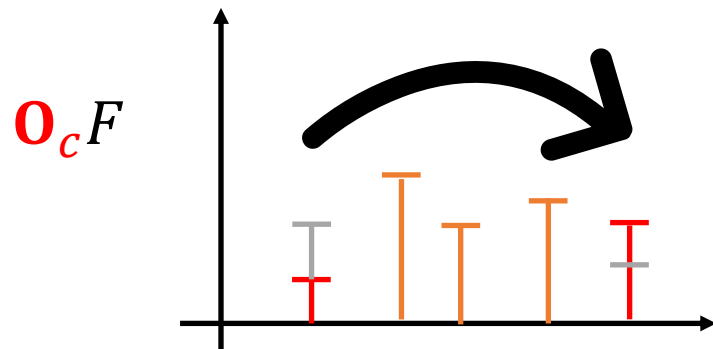
**Def.** A distribution  $F$  is **optimistic** over  $G$  if

$$U_{\beta}(F) \geq U_{\beta}(G)$$

Denote by  $F \geq G$

Given confidence ball of distributions

$$B(F, c) := \{G \mid \|G - F\| \leq c\}$$



**Distributional optimism operator**

$$O_c F \geq G, \forall G \in B(F, c)$$

**Intuition** move probability mass from the lower tail to the maximum value

- Keramati, Ramtin, et al. "Being optimistic to be conservative: Quickly learning a CVaR policy." *Proceedings of the AAAI conference on artificial intelligence*. Vol. 34. No. 04. 2020.
- Hao Liang, and Zhi-Quan Luo. "A distribution optimization framework for confidence bounds of risk measures." *International Conference on Machine Learning*. PMLR, 2023.

# Optimism of RODI

**RODI**

$$\eta_h \leftarrow \mathbf{O}_c \widehat{\mathbf{T}}_d v_{h+1}$$

**Induction**  $v_{h+1} \geq v_{h+1}^*$

$$\|\widehat{\mathbf{T}}_d v_{h+1} - \mathbf{T}_d v_{h+1}\| \leq c_h \implies \mathbf{T}_d v_{h+1} \in B(\widehat{\mathbf{T}}_d v_{h+1}, c_h)$$

**Distributional optimism operator**

$$\eta_h = \mathbf{O}_c \widehat{\mathbf{T}}_d v_{h+1} \geq \mathbf{T}_d v_{h+1}$$

**Induction**



$$\eta_h \geq \mathbf{T}_d v_{h+1} \geq \mathbf{T}_d v_{h+1}^* = \eta_h^* \longrightarrow v_h \geq v_h^*$$

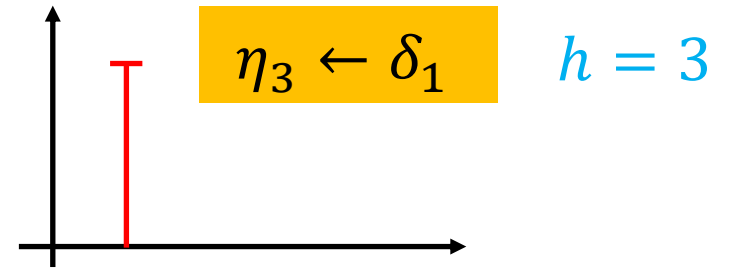
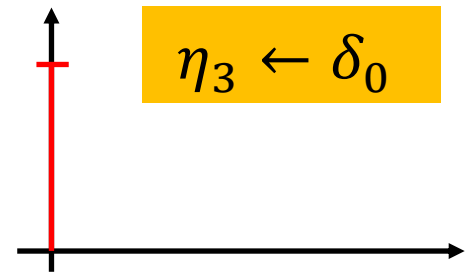
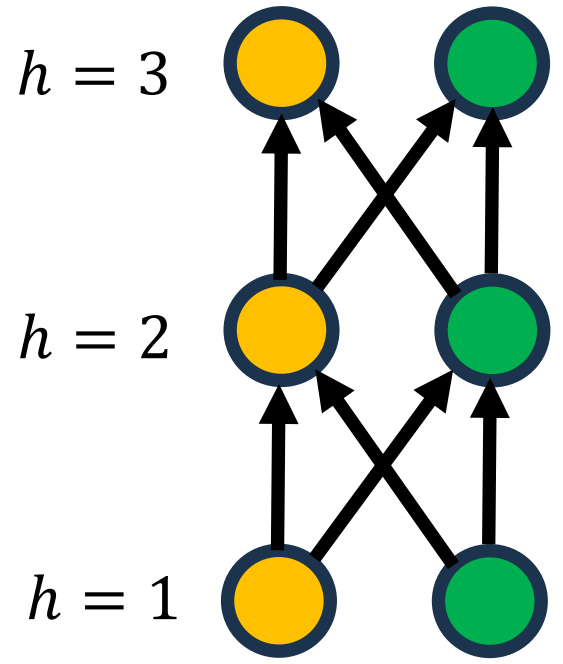
**Def.**

$$U_\beta(\eta_h(s, a)) \geq U_\beta(\eta_h^*(s, a)) \longrightarrow Q_h(s, a) \geq Q_h^*(s, a)$$



**Optimism of risk value**

# Computational Inefficiency of RODI

State space = {, }  
Uniform transition  
 $r(\text{green}) = 0, r(\text{yellow}) = 1$

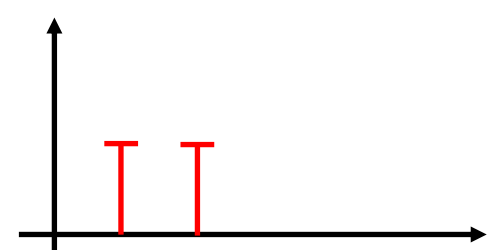
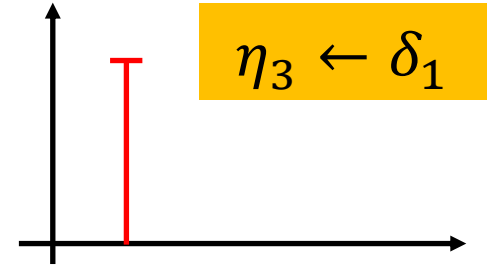
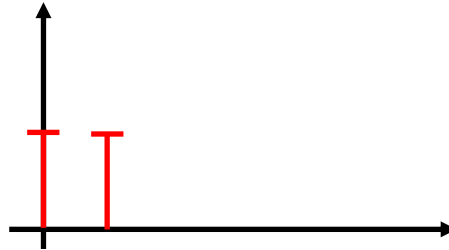
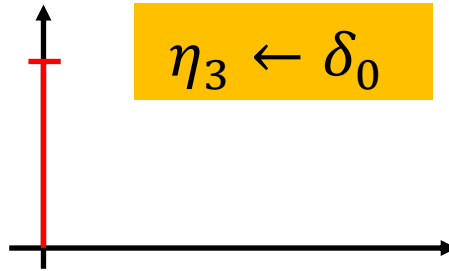
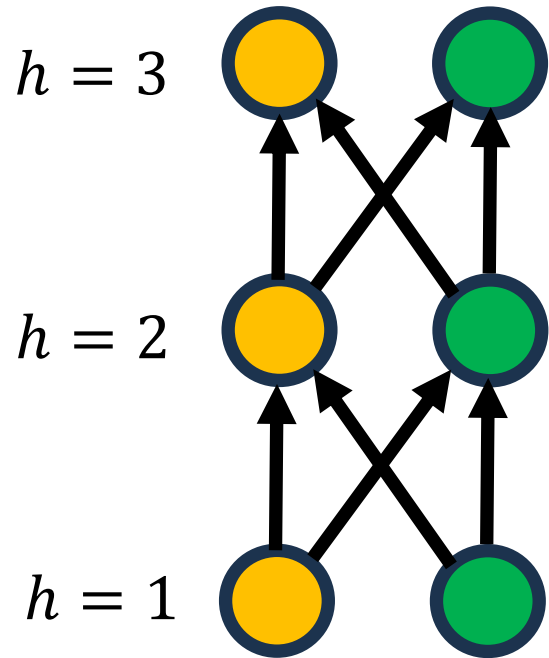


# Computational Inefficiency of RODI

State space = {, 

Uniform transition



$r(\text{green}) = 0, r(\text{yellow}) = 1$

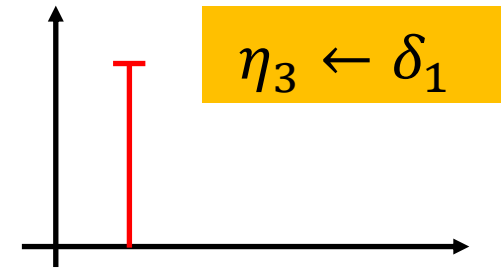
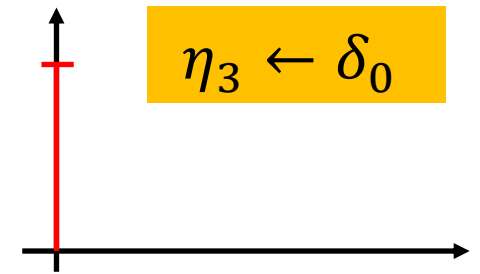
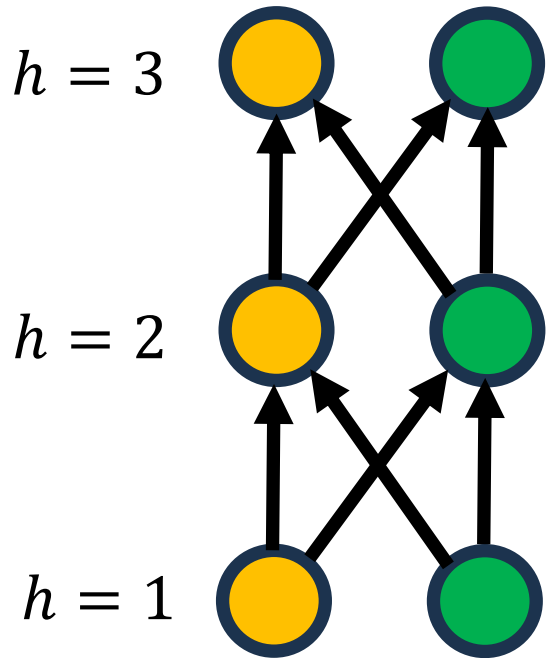


$h = 3$

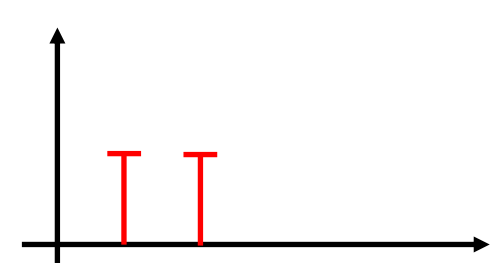
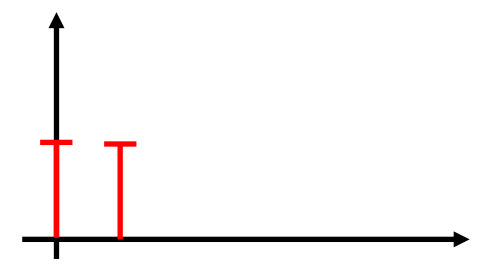
$h = 2$   
 $\eta_2 \leftarrow \mathbf{T}_d \nu_3$

# Computational Inefficiency of RODI

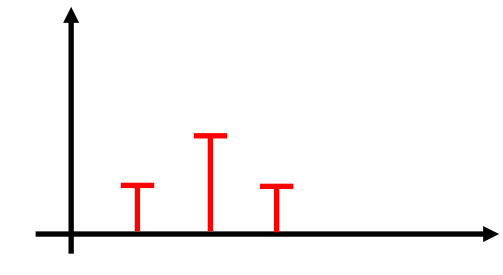
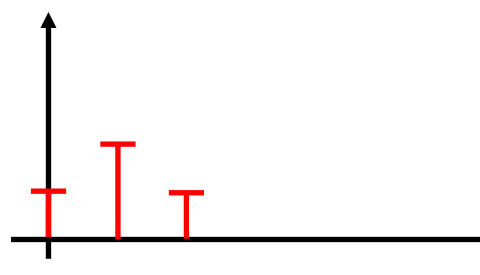
State space = {, }  
 Uniform transition  
 $r(\text{green}) = 0, r(\text{yellow}) = 1$



$h = 3$



$h = 2$   
 $\eta_2 \leftarrow \mathbf{T}_d \nu_3$



$h = 1$   
 $\eta_1 \leftarrow \mathbf{T}_d \nu_2$

Operator  $\mathbf{T}_d$  increases support size exponentially!

# RODI with Distribution Representation (RODI-Rep)

Operator  $\mathbf{T}_d$  increases support size exponentially

$$\eta_h \leftarrow \mathbf{T}_d \nu_{h+1} \quad \longrightarrow \quad |\eta_h| = S \cdot |\nu_{h+1}|$$



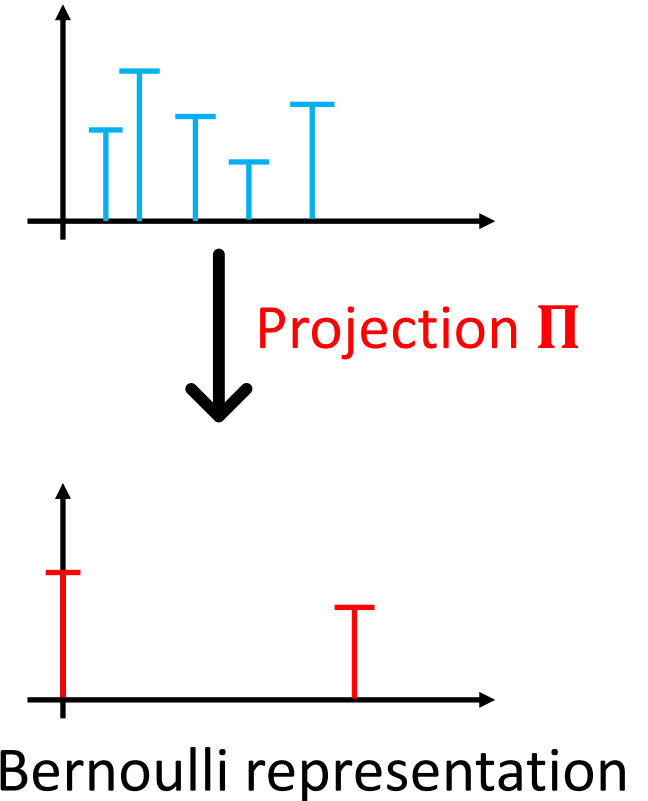
# RODI with Distribution Representation (RODI-Rep)

Operator  $\mathbf{T}_d$  increases support size exponentially

$$\eta_h \leftarrow \mathbf{T}_d \nu_{h+1} \quad \longrightarrow \quad |\eta_h| = S \cdot |\nu_{h+1}|$$

Represent distribution with **fixed** support via **projection**

$$\eta_h \leftarrow \mathbf{\Pi} \mathbf{T}_d \nu_{h+1} \quad \longrightarrow \quad |\eta_h| = |\nu_{h+1}| = |\eta_{h+1}|$$



# RODI with Distribution Representation (RODI-Rep)

Operator  $\mathbf{T}_d$  increases support size exponentially

$$\eta_h \leftarrow \mathbf{T}_d \nu_{h+1} \quad \longrightarrow \quad |\eta_h| = S \cdot |\nu_{h+1}|$$

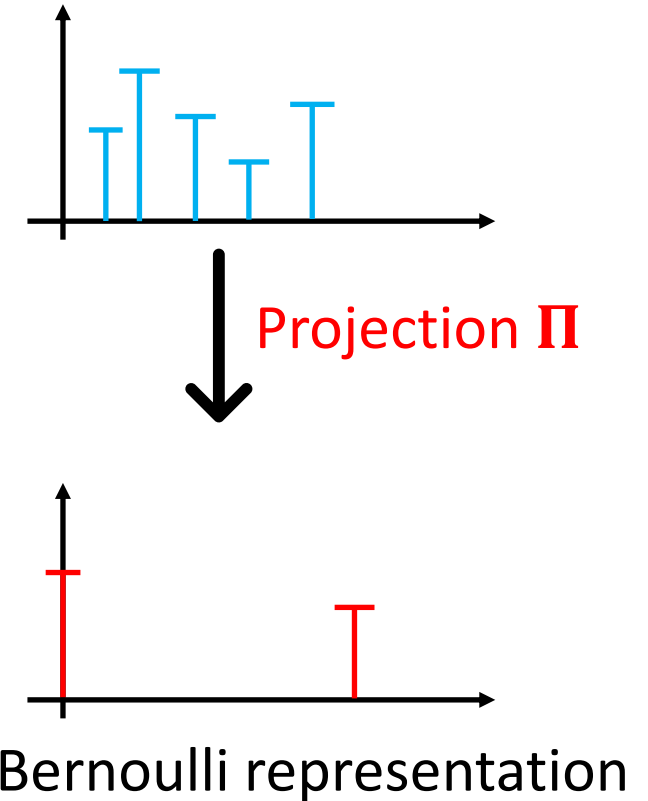
Represent distribution with **fixed** support via **projection**

$$\eta_h \leftarrow \mathbf{\Pi} \mathbf{T}_d \nu_{h+1} \quad \longrightarrow \quad |\eta_h| = |\nu_{h+1}| = |\eta_{h+1}|$$

**RODI-Rep**

$$\eta_h \leftarrow \mathbf{\Pi} \mathbf{O}_c \widehat{\mathbf{T}}_d \nu_{h+1} \quad \longrightarrow \quad |\eta_h| = |\nu_{h+1}|$$

Ensure optimism while maintaining **computational efficiency**



# Regret **Lower** Bound

$T := KH$  total time steps

For any algorithm, exists an MDP instance

$$\mathbb{E}[\text{Regret}(K)] \geq \Omega\left(\frac{\exp(\beta H/6) - 1}{\beta} \sqrt{SAT}\right)$$

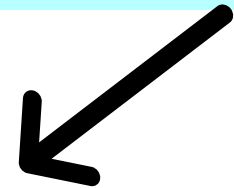


Fundamental trade-off between **risk sensitivity** and **sample complexity**

- Fix and **tighten** the previous lower bound
- Recover **tight** lower bound under risk-neutral setting
- Hold for  $\beta > 0$

Corrected **result** in [3]

$$\mathbb{E}[\text{Regret}(K)] \geq \Omega\left(\frac{\exp(|\beta|H/2) - 1}{|\beta|} \sqrt{K \log K}\right)$$



**Missing**  $S, A$   
**Loose dependency on**  $H$

- Reduction to **2-armed bandit**
- Proof has errors

[3] Fei, Yingjie, et al. "Risk-sensitive reinforcement learning: Near-optimal risk-sample tradeoff in regret." *Advances in Neural Information Processing Systems* 33 (2020): 22384-22395.

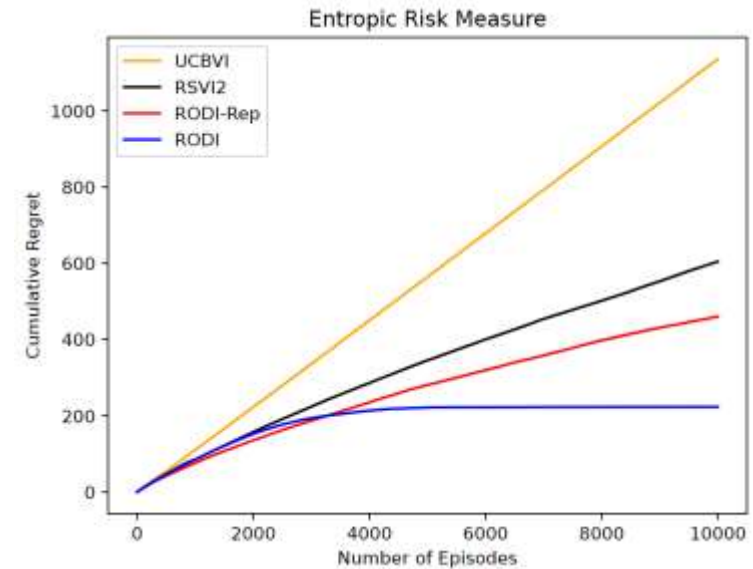
# Regret **Upper** Bound

Algorithm	Regret bound	Time	Space
RSVI	$\tilde{\mathcal{O}} \left( \frac{\exp( \beta H) - 1}{ \beta } \sqrt{HS^2AT} \right)$	$\mathcal{O}(TS^2A)$	$\mathcal{O}(HSA + T)$
RSVI2			
<b>RODI-Rep</b> <b>RODI</b>		$\mathcal{O}(KS^H)$	$\mathcal{O}(S^H)$
lower bound	$\Omega \left( \frac{\exp(\beta H/6) - 1}{\beta} \sqrt{SAT} \right)$	-	-

**RODI** and **RODI-Rep** satisfies

$$\text{Regret}(K) \leq \mathcal{O} \left( \frac{\exp(|\beta|H) - 1}{|\beta|} \sqrt{HS^2AT} \right)$$

- Matching the **best known result** in [4]
- **Clean** and interpretable analysis
- Outperform **RSVI2** [4] empirically



[4] Fei, Yingjie, et al. "Exponential bellman equation and improved regret bounds for risk-sensitive reinforcement learning." *Advances in Neural Information Processing Systems* 34 (2021): 20436-20446.

# Summary

## Main Contributions

- Risk-Sensitive **Distributional** Dynamic Programming framework
- **Computationally efficient** DRL algorithm with near-optimal **regret** guarantee
- Tight regret lower bound

# Summary

## Main Contributions

- Risk-Sensitive **Distributional** Dynamic Programming framework
- **Computationally efficient** DRL algorithm with near-optimal **regret** guarantee
- Tight regret lower bound

## Future Directions

- **Scalability** issues with large state-action spaces: **approximation** techniques
- **Multi-Agent** RSRL: coordination and cooperation
- Theoretical foundations: connections to **robust** control
- **Application** domains: cyber-physical systems

**Thank You**