

Risk-Aware Transfer in Reinforcement Learning using Successor Features

Michael Gimelfarb^{1,3} André Barreto² Scott Sanner^{1,3} Chi-Guhn Lee¹

¹University of Toronto

²DeepMind

³Vector Institute (Affiliate Program)

Introduction

Motivation

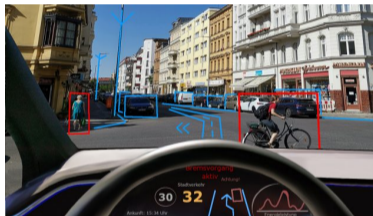
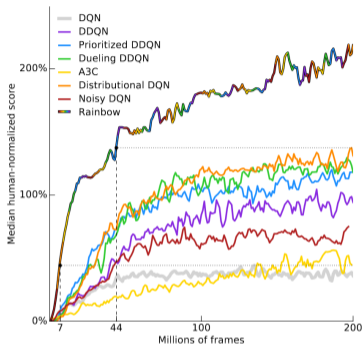


Figure 1: Sample Efficiency (Source: original paper on Rainbow DQN)

Figure 2: Risk-Awareness (Source: Wikimedia Commons)

Motivation

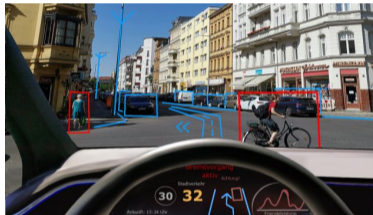
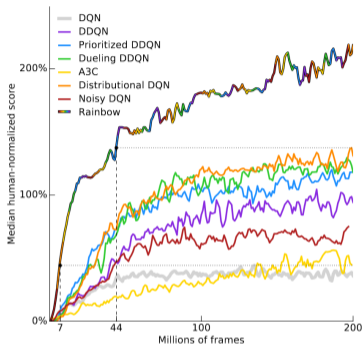


Figure 1: Sample Efficiency (Source: original paper on Rainbow DQN)

Figure 2: Risk-Awareness (Source: Wikimedia Commons)

- **transfer learning**

- replace $E[]$ by non-linear **utility** $U[]$

Motivation

Our goals:

Motivation

Our goals:

- transfer between tasks with **shared dynamics and different goals**

Motivation

Our goals:

- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting¹

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

Motivation

Our goals:

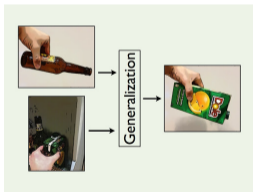
- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting¹
- provide **task generalization** by exploiting the structure of the task/reward space

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

Motivation

Our goals:

- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting¹
- provide **task generalization** by exploiting the structure of the task/reward space
- design a method suitable for offline RL²



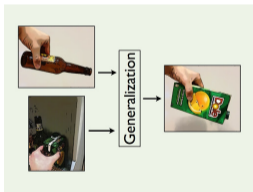
¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

²Levine, Sergey, et al. "Offline reinforcement learning..." arXiv. 2020.

Motivation

Our goals:

- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting¹
- provide **task generalization** by exploiting the structure of the task/reward space
- design a method suitable for offline RL²



- incorporate risk explicitly by e.g. penalizing the variance of returns

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

²Levine, Sergey, et al. "Offline reinforcement learning..." arXiv. 2020.

Introduce **Risk-Aware Successor Features** (RaSF)

Motivation

Introduce **Risk-Aware Successor Features (RaSF)**

	Transfers Skills	Exploits Task Structure	Risk-Sensitive
Risk-Aware RL	7	7	3
Risk-Aware Transfer	3	7	3
Successor Features	3	3	7
RaSF (Ours)	3	3	3

Preliminaries – Successor Features

- $E[R_0 + R_1 + \dots]$ requires averaging all possible future outcomes of the world – **curse of dimensionality**

- $E[R_0 + R_1 + \dots]$ requires averaging all possible future outcomes of the world – **curse of dimensionality**
- use cached $Q(s^\theta; a^\theta)$ in each successor state to **bootstrap** the estimated Q-values in state s

$$Q(s; a) = E_{S^\theta \sim P(\cdot|s;a)}[R_t + Q(S^\theta; (S^\theta))]$$

Policy Improvement

Suppose an initial policy π_0 is given:

Suppose an initial policy is given:

^ compute the value of , e.g.

$$Q(s; a) = E_{S^0 \sim P(\cdot|s;a)}[R_t + Q(S^0; (S^0))] \quad / \text{ policy evaluation}$$

Suppose an initial policy is given:

^ compute the value of Q , e.g.

$$Q(s; a) = E_{S^0 \sim P(\cdot|s;a)}[R_t + \gamma Q(S^0; (S^0))] \quad / \text{ policy evaluation}$$

^ construct a new policy π^0 according to

$$\pi^0(s) = \arg \max_{a \in A} Q(s; a) \quad / \text{ policy improvement}$$

Suppose an initial policy is given:

- ^ compute the value of Q , e.g.

$$Q(s; a) = E_{S^0 \sim P(\cdot|s;a)}[R_t + \gamma Q(S^0; (S^0))] \quad / \text{ policy evaluation}$$

- ^ construct a new policy π^0 according to

$$\pi^0(s) = \arg \max_{a \in A} Q(s; a) \quad / \text{ policy improvement}$$

Policy Improvement Theorem: π^0 is "better" than π , e.g. $Q^{\pi^0}(s; a) \geq Q^{\pi}(s; a)$

Generalized Policy Iteration

Generalized Policy Iteration

Alternating between evaluation and improvement leads to an optimal policy

Generalized Policy Iteration

Alternating between evaluation and improvement leads to an optimal policy

Key Idea: Replace π with multiple source policies $\pi_1 :: \pi_n$.

Generalized Policy Improvement 2.0

Suppose $\pi_1; \dots; \pi_n$ are given¹:

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

Generalized Policy Improvement 2.0

Suppose Q_1, \dots, Q_n are given¹:

^ compute the values

$Q_1(s; a), \dots, Q_n(s; a)$ / generalized policy evaluation (GPE)

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

Generalized Policy Improvement 2.0

Suppose $Q^1; \dots; Q^n$ are given¹:

^ compute the values

$Q^1(s; a); \dots; Q^n(s; a)$ / generalized policy evaluation (GPE)

^ pick the policy with the best return

$$i^* = \arg \max_i Q^i(s; a)$$

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

Generalized Policy Improvement 2.0

Suppose $Q^1; \dots; Q^n$ are given¹:

^ compute the values

$$Q^1(s; a); \dots; Q^n(s; a) \quad / \text{ generalized policy evaluation (GPE)}$$

^ pick the policy with the best return

$$i^* \in \arg \max_i Q^i(s; a)$$

^ construct Q^0 as usual but w.r.t. the "best" policy i^*

$$Q^0(s) \in \arg \max_{a \in A} Q^{i^*}(s; a) = \arg \max_{a \in A} \max_{i=1 \dots n} Q^i(s; a) \quad / \text{ generalized policy improvement (GPI)}$$

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

Generalized Policy Improvement 2.0

Suppose $Q^1; \dots; Q^n$ are given¹:

^ compute the values

$$Q^1(s; a); \dots; Q^n(s; a) \quad / \text{ generalized policy evaluation (GPE)}$$

^ pick the policy with the best return

$$i \in \arg \max_i Q^i(s; a)$$

^ construct Q^0 as usual but w.r.t. the "best" policy π_i

$$Q^0(s) \in \arg \max_{a \in A} Q^i(s; a) = \arg \max_{a \in A} \max_{i=1 \dots n} Q^i(s; a) \quad / \text{ generalized policy improvement (GPI)}$$

Key result: Q^0 is better than π_i .

¹Barreto, Andre, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

Preliminaries { Risk-Aversion in MDPs using Entropic Utility Functions

Optimizing risk measures in sequential problems is \mathbb{R} ard

²Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

Optimizing risk measures in sequential problems is \mathbb{R} hard
^ considering optimizing a \mathbb{R} nal cost:

²Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

Time-Consistency

Optimizing risk measures in sequential problems is hard

^ considering optimizing a final cost:

^ consider the dynamic risk measure

$$\rho_{k;N}(Z) = \max_{p \in \mathcal{P}_{k;N}} E_p[Z | \mathcal{F}_k];$$

for $k = 0; 1; 2$

²Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

Time-Consistency

Optimizing risk measures in sequential problems is hard

^ considering optimizing a final cost Z :

^ it is easy to check that $\rho_1(Z) = 60$ for all \mathbb{P} , so Z is riskier than a deterministic cash flow of $W = 50$ at time 1

^ consider the dynamic risk measure

$$\rho_{k;N}(Z) = \max_{\mathbb{P} \in \mathcal{P}_{0:T}} E_{\mathbb{P}}[Z | \mathcal{F}_k];$$

for $k = 0; 1; 2$

²Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

Optimizing risk measures in sequential problems is hard

^ considering optimizing a final cost:

^ it is easy to check that $v_1(Z)(\omega) = 60$ for all ω , so Z is riskier than a deterministic cash flow of $W = 50$ at time 1

^ consider the dynamic risk measure

^ yet, $v_0(Z)(\omega) = 40$ and so Z is less risky than W at time 0!

$$v_{k;N}(Z) = \max_{p \in \mathcal{P}_k} E_p[Z | \mathcal{F}_k];$$

for $k = 0; 1; 2$

²Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

Optimizing risk measures in sequential problems is hard

^ considering optimizing a final cost:

^ it is easy to check that $v_1(Z) = 60$ for all ω , so Z is riskier than a deterministic cash flow of $W = 50$ at time 1

^ consider the dynamic risk measure

^ yet, $v_0(Z) = 40$ and so Z is less risky than W at time 0!

$$v_{k;N}(Z) = \max_{p \in \mathcal{P}_k} E_p[Z | \mathcal{F}_k];$$

Key idea: Z has become riskier just because time has passed!

for $k = 0; 1; 2$

²Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

Variance Reduction in MDPs with Entropic Utilities

We use the entropic utility to measure risk:

Variance Reduction in MDPs with Entropic Utilities

We use the entropic utility to measure risk:

^ defined in terms of the moment-generating function

$$U[R] = -\log E e^{h R^i}$$

Variance Reduction in MDPs with Entropic Utilities

We use the entropic utility to measure risk:

^ defined in terms of the moment-generating function

$$U[R] = -\log E[e^{hR}]$$

^ has the Taylor expansion

$$U[R] = E[R] + \frac{1}{2}V[R] + O(\sigma^2) \quad / \quad \text{is a level of risk-aversion}$$

Variance Reduction in MDPs with Entropic Utilities

We use the entropic utility to measure risk:

^ defined in terms of the moment-generating function

$$U[R] = -\log E e^{h R}$$

^ has the Taylor expansion

$$U[R] = E[R] + \frac{1}{2} V[R] + O(\epsilon^2) \quad / \quad \epsilon \text{ is a level of risk-aversion}$$

^ connected to the mean-variance optimization in MDPs

¹Mannor, Shie, and John N. Tsitsiklis. "Mean-variance optimization in Markov decision processes." ICML. 2011.

Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- ^ dynamic programming: has a Bellman equation formulation

Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- ^ dynamic programming: has a Bellman equation formulation

- ^ recursive property: behaves similar to expectation in total-reward episodic MDPs¹ (and discounted MDPs with simple modifications)

¹Osogami, Takayuki. "Robustness and risk-sensitivity in Markov decision processes." NeurIPS 2012.

Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- ^ dynamic programming: has a Bellman equation formulation
- ^ recursive property: behaves similar to expectation in total-reward episodic MDPs¹ (and discounted MDPs with simple modifications)
- ^ convex/concave: satisfies properties that can be seen as rational decision making

¹Osogami, Takayuki. "Robustness and risk-sensitivity in Markov decision processes." NeurIPS 2012.

Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- ^ dynamic programming: has a Bellman equation formulation

- ^ recursive property: behaves similar to expectation in total-reward episodic MDPs¹ (and discounted MDPs with simple modifications)
- ^ convex/concave: satisfies properties that can be seen as rational decision making
- ^ time consistency: can focus on Markov policies

¹Osogami, Takayuki. "Robustness and risk-sensitivity in Markov decision processes." NeurIPS 2012.

Theory

Motivating Example

Why is the problem non-trivial?

Motivating Example

Why is the problem non-trivial?

Motivating Example

1. Define a family of tasks:

^ two source tasks:

low failure cost + high failure cost

^ one target task:

only X has high failure cost

Motivating Example

1. Define a family of tasks:

^ two source tasks:

low failure cost + high failure cost

^ one target task:

only X has high failure cost

2. Solve them with VI for xed

Motivating Example

1. Define a family of tasks:

^ two source tasks:

low failure cost + high failure cost

^ one target task:

only X has high failure cost

3. Apply risk-aware and risk-neutral GPI:

2. Solve them with VI for xed

Motivating Example

1. Define a family of tasks:

^ two source tasks:

low failure cost + high failure cost

^ one target task:

only X has high failure cost

3. Apply risk-aware and risk-neutral GPI:

2. Solve them with VI for fixed

Conclusion: only risk-aware GPI results in the correct target policy

Key Theoretical Results

Key Theoretical Results

Armed with this knowledge, we prove that risk-aware GPI:

Key Theoretical Results

Armed with this knowledge, we prove that risk-aware GPI:

$\hat{\pi}$ is a strict policy improvement operator

Key Theoretical Results

Armed with this knowledge, we prove that risk-aware GPI:

$\hat{\pi}$ is a strict policy improvement operator

$\hat{\pi}$ is optimal up to an irreducible task discrepancy gap

Generalized Policy Evaluation

Assume linear reward:

$$r(s; a; s^0) = \phi(s; a; s^0)^T w$$

Generalized Policy Evaluation

Assume linear reward:

$$r(s; a; s^0) = \phi(s; a; s^0)^T w$$

Now:

$$\begin{aligned}
 Q_w(s; a) &= E \sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s; A_0 = a; \pi \\
 &= E \sum_{t=0}^{\infty} \gamma^t \phi_t^T w \mid S_0 = s; A_0 = a; \pi \\
 &= E \sum_{t=0}^{\infty} \gamma^t \underbrace{\phi_t}_{z_t(s; a)} \mid S_0 = s; A_0 = a; \pi
 \end{aligned}$$

Can generalize GPE tdistributions of return :

$$Q_{h;}(s; a) = U \sum_{t=h}^{\infty} \gamma^t r(s_t; a_t(s_t); s_{t+1}) = U [\sum_{t=h}^{\infty} \gamma^t Q_{h;}(s; a) | w]$$

Can generalize GPE distributions of return :

$$Q_{h;}(s; a) = U \left[\sum_{t=h}^T \gamma^t r(s_t; a_t(s_t); s_{t+1}) \right] = U \left[\sum_{t=h}^T r_t(s; a)^t w \right]$$

One simple trick is to Taylor expand to the second moment:

$$U \left[\sum_{t=h}^T r_t(s; a)^t w \right] = E_P \left[\sum_{t=h}^T r_t(s; a)^t w \right] + \frac{1}{2} \text{Var}_P \left[\sum_{t=h}^T r_t(s; a)^t w \right] + O(\epsilon^2)$$

$$\sum_{t=h}^T r_t(s; a)^t w + \frac{1}{2} w^t \text{Var}_P \left[\sum_{t=h}^T r_t(s; a)^t w \right] = Q_{h;}(s; a)$$

Can generalize GPE to distributions of return :

$$Q_{h;}(s; a) = U \left[\sum_{t=h}^T \gamma^t r(s_t; a_t(s_t); s_{t+1}) \right] = U \left[\sum_{h(s; a)}^T w \right]$$

One simple trick is to Taylor expand to the second moment:

$$U \left[\sum_{h(s; a)}^T w \right] = E_P \left[\sum_{h(s; a)}^T w \right] + \frac{1}{2} \text{Var}_P \left[\sum_{h(s; a)}^T w \right] + O(\epsilon^2)$$

$$\sum_{h(s; a)}^T w + \frac{1}{2} w^T \text{Var}_P \left[\sum_{h(s; a)}^T w \right] = Q_{h;}(s; a)$$

Reduces to a (simpler) problem of estimating sufficient statistics of the feature occupancy

Experiments

Two domains from Barreto et al., 2017:

Two domains from Barreto et al., 2017:

Introduce reward volatility :

- ^ traps **X** for four-room
- ^ action noise + danger zones for reacher

Four-Room

Train on a sequence of 128 random task instances, for 1000 steps each

Four-Room

Train on a sequence of 128 random task instances, 1000 steps each

Four-Room

Sensitivity to parameter:

Four-Room

Sensitivity to parameter:

Train on four source tasks, test periodically on 8 unseen test tasks:

Train on four source tasks, test periodically on 8 unseen test tasks:

Does the agent learn risk-sensitive behavior?

Does the agent learn risk-sensitive behavior?

How sensitive is the agent to ϵ ? Does the C51 method help in learning SFs?

How sensitive is the agent to ϵ ? Does the C51 method help in learning SFs?

Conclusion

Conclusion

In conclusion:

Conclusion

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals

Conclusion

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals
- we extended generalized policy improvement to the risk-aware setting with entropic utilities

Conclusion

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals
- we extended generalized policy improvement to the risk-aware setting with entropic utilities
- we then extended the notion of generalized policy evaluation via the Taylor expansion of the entropic utility

Conclusion

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals
- we extended generalized policy improvement to the risk-aware setting with entropic utilities
- we then extended the notion of generalized policy evaluation via the Taylor expansion of the entropic utility
- together, risk-aware GPI and GPE are shown to inherit the superior task generalization abilities of successor features, while also learning to avoid risky situations

Thank you.