

Risk-Sensitive Policy Optimization via Predictive CVaR Policy Gradient

Ju-Hyun Kim Seungki Min

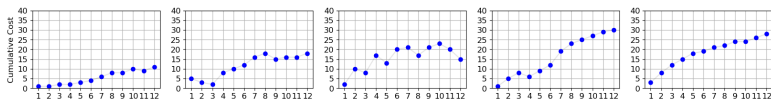
{vgb8111, skmin}@kaist.ac.kr

Department of Industrial and Systems Engineering
KAIST

ICML 2024, Jul 2024

Motivating Example

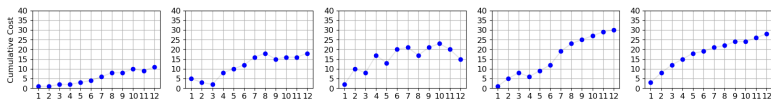
- Given a fixed policy π , consider 5 sample trajectories. For a risk-neutral RL, we can utilize the whole sample trajectories.



$$\mathbb{E}[C_{1:T}] \approx \frac{1}{N} \sum_{i \in [N]} C_{1:T}^{(i)}$$

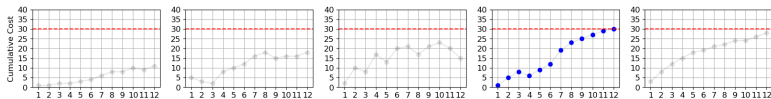
Motivating Example

- Given a fixed policy π , consider 5 sample trajectories. For a risk-neutral RL, we can utilize the whole sample trajectories.



$$\mathbb{E}[C_{1:T}] \approx \frac{1}{N} \sum_{i \in [N]} C_{1:T}^{(i)}$$

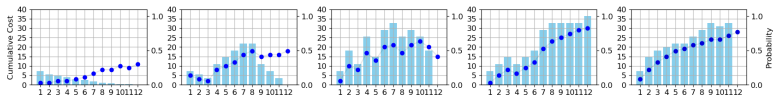
- For a CVaR RL, one may utilize **the worst q fraction** among the whole sample trajectories. \rightarrow High variance and low sample efficiency



$$\text{CVaR}_q[C_{1:T}] = \mathbb{E}[C_{1:T} | C_{1:T} \geq \text{VaR}_q[C_{1:T}]] \approx \frac{1}{qN} \sum_{qN \text{ worst}} C_{1:T}^{(i)}$$

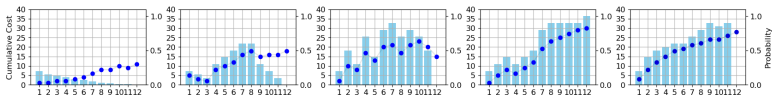
Main Idea

- We introduce a “*predictive tail probability process*” $Q^\pi = (Q_t^\pi)_{t \in [T]}$.
 - In each period, it predicts the probability that the current sample path ends up being one of the worst q fraction of outcomes.

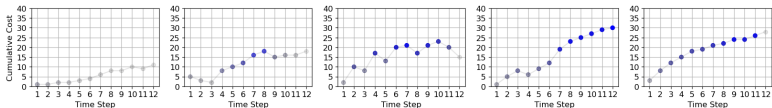


Main Idea

- We introduce a “predictive tail probability process” $Q^\pi = (Q_t^\pi)_{t \in [T]}$.
 - In each period, it predicts the probability that the current sample path ends up being one of the worst q fraction of outcomes.



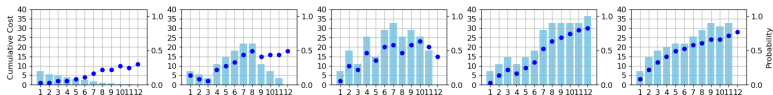
- We reformulate CVaR as an expectation of reweighted cost realization by the “predictive tail probability”.



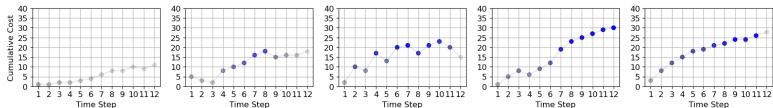
$$\text{CVaR}_q[C_{1:T}] = \frac{1}{q} \mathbb{E} \left[\sum_{t \in [T]} Q_t C_t \right] \approx \frac{1}{qN} \sum_{i \in [N]} \sum_{t \in [T]} Q_t^{(i)} C_t^{(i)}$$

Main Idea

- We introduce a “predictive tail probability process” $Q^\pi = (Q_t^\pi)_{t \in [T]}$.
 - In each period, it predicts the probability that the current sample path ends up being one of the worst q fraction of outcomes.



- We reformulate CVaR as an expectation of reweighted cost realization by the “predictive tail probability”.



$$\text{CVaR}_q[C_{1:T}] = \frac{1}{q} \mathbb{E} \left[\sum_{t \in [T]} Q_t C_t \right] \approx \frac{1}{qN} \sum_{i \in [N]} \sum_{t \in [T]} Q_t^{(i)} C_t^{(i)}$$

- For CVaR RL, we can utilize **all sample trajectories** via reweighting.
 - Low variance, high sample efficiency

Predictive CVaR Policy Gradient (PCVaR)

- Our goal is to find an optimal policy π^* solving

$$\min_{\pi \in \Pi^\Theta} \{J_q(\pi) := q \cdot \text{CVaR}_q^\pi [C_{1:T}]\}. \quad (*)$$

- Using **reformulated CVaR objective**, we change the optimization (*) into adjusted optimization problem with parameters (θ, η, ϕ) as

$$\min_{\theta \in \Theta} \{J(\theta, \eta, \phi) \mid \eta, \phi \text{ s.t. } \dots\},$$

where

$$J(\theta, \eta, \phi) := \mathbb{E} \left[\sum_{t \in [T]} \hat{Q}_t C_t \right], \hat{Q}_t = f^\phi(X_{t+1}, C_{1:t} - \eta).$$

- Compared to risk-neutral PG objective, we just replace C_t as $\hat{Q}_t C_t$ (even for policy learning process).

Predictive CVaR Policy Gradient (PCVaR)

- The objective $J_q(\pi)$ in (*) can be rewritten as

$$J_q(\pi) = \min_{\eta \in \mathbb{R}} \mathbb{E}^{\pi} [q\eta + (C_{1:T} - \eta)^+]. \quad (1)$$

- Given $\pi \in \Pi^{\mathcal{H}}$, optimal solution η^{π} of (1) is $\text{VaR}_q(C_{1:T})$.

Predictive CVaR Policy Gradient (PCVaR)

- The objective $J_q(\pi)$ in (*) can be rewritten as

$$J_q(\pi) = \min_{\eta \in \mathbb{R}} \mathbb{E}^{\pi} [q\eta + (C_{1:T} - \eta)^+]. \quad (1)$$

- Given $\pi \in \Pi^{\mathcal{H}}$, optimal solution η^{π} of (1) is $\text{VaR}_q(C_{1:T})$.

Definition: Predictive tail probability process $Q^{\pi, \eta}$

Given $\pi \in \Pi^{\mathcal{H}}$ and $\eta \in \mathbb{R}$, $Q^{\pi, \eta} = (Q_t^{\pi, \eta})_{t \in \{0, \dots, T\}}$ is defined as follow:

$$Q_t^{\pi, \eta} := \mathbb{P}(C_{1:T} \geq \eta | H_{t+1})$$

Predictive CVaR Policy Gradient (PCVaR)

- The objective $J_q(\pi)$ in (*) can be rewritten as

$$J_q(\pi) = \min_{\eta \in \mathbb{R}} \mathbb{E}^\pi [q\eta + (C_{1:T} - \eta)^+]. \quad (1)$$

- Given $\pi \in \Pi^{\mathcal{H}}$, optimal solution η^π of (1) is $\text{VaR}_q(C_{1:T})$.

Definition: Predictive tail probability process $Q^{\pi, \eta}$

Given $\pi \in \Pi^{\mathcal{H}}$ and $\eta \in \mathbb{R}$, $Q^{\pi, \eta} = (Q_t^{\pi, \eta})_{t \in \{0, \dots, T\}}$ is defined as follow:

$$Q_t^{\pi, \eta} := \mathbb{P}(C_{1:T} \geq \eta | H_{t+1})$$

Proposition 3.2: Reformulation of CVaR objective

If $\sum_{t \in [T]} C_t^\pi$ has no probability mass at $\eta^\pi = \text{VaR}_q^\pi[C_{1:T}]$,

$$J_q(\pi) = \mathbb{E}^\pi \left[\sum_{t \in [T]} Q_t^{\pi, \eta^\pi} C_t \right].$$

Predictive CVaR Policy Gradient (PCVaR)

- We decompose the CVaR policy optimization (*) into *three* optimization problems with **reformulated objective**.

$$\min_{\theta \in \Theta} \left\{ J(\theta, \eta, \phi) \mid \begin{array}{l} \eta \in \arg \min_{\eta' \in \mathbb{R}} L(\theta, \eta'), \\ \phi \in \arg \min_{\phi' \in \Phi} M(\theta, \eta, \phi') \end{array} \right\},$$

where

$$J(\theta, \eta, \phi) := \mathbb{E} \left[\sum_{t \in [T]} \hat{Q}_t C_t \right], \quad L(\theta, \eta) := \mathbb{E} [q\eta + (C_{1:T} - \eta)^+],$$

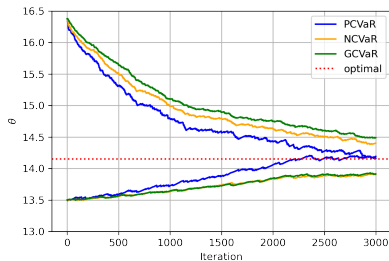
$$M(\theta, \eta, \phi) := \mathbb{E} \left[\sum_{t \in [T]} \left(\mathbb{I}\{C_{1:T} \geq \eta\} - \hat{Q}_t \right)^2 \right].$$

- Update θ (risk-neutral PG with $\hat{Q}_t C_t$), η (simple SGD), ϕ (typical supervised learning) in parallel.

- Consistency of the estimators related to ϕ and η (respectively **Proposition 4.1** and **Proposition 4.2**)
- Unbiasedness of the gradient estimators of $J(\theta, \eta, \phi)$ (**Theorem 4.3**)
- Variance reduction in the gradient estimation of PCVAR (**Proposition 4.5**)

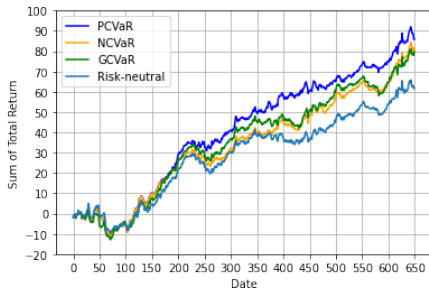
Numerical Experiments

- Continuous Blackjack (synthetic data)



Evaluation point	PCVaR	NCVaR	GCVaR
$\theta = \arg \max_{\theta} E[R_{1:T}]$	33.66	97.59	96.92
$\theta = \arg \max_{\theta} CVaR_q[R_{1:T}]$	6.71	24.01	24.67

- Pair trading (real-world data)



Contribution

- We suggest “*Predictive CVaR Policy Gradient (PCVAR)*”, relying on

$$q \cdot \text{CVaR}_q [C_{1:T}] = \mathbb{E} \left[\sum_{t \in [T]} Q_t C_t \right],$$

where $Q_t = \mathbb{P}(\text{current sample} \in \text{the worst } q \text{ fractions} \mid \text{history})$.
conditional $\mathbb{E} \rightarrow$ risk-neutral \mathbb{E}

- PCVAR utilizes **all sample trajectories**.
→ Improves sample efficiency and then accelerates learning.
- PCVAR can be applied on top of any risk-neutral policy gradient algorithm.
- Its effectiveness is demonstrated with theoretical analyses and numerical experiments.