

Conditionally Elicitable Dynamic Risk Measures for Deep Reinforcement Learning

Anthony Coache (University of Toronto)

Joint work with

Sebastian Jaimungal (University of Toronto & Oxford-Man Institute)

and

Álvaro Cartea (Oxford-Man Institute & University of Oxford)

SIAG/FME Conference Paper Prize Session ★ June 9, 2023 ★ Philadelphia, USA



Statistical Sciences
UNIVERSITY OF TORONTO



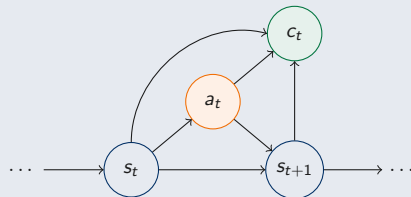
Conference on
Financial Mathematics and
Engineering

Reinforcement Learning (RL)

- **Model-agnostic** framework for **learning-based control**
- Learning optimal behaviours from interactions to minimise a cost signal

Markov Decision Process ($\mathcal{S}, \mathcal{A}, \pi, \mathbb{P}, c$)

- \mathcal{S} – State space
- \mathcal{A} – Action space
- $\pi^\theta(a_t|s_t)$ – Randomised policy characterised by θ
- $\mathbb{P}(s_0), \mathbb{P}(s_{t+1}|s_t, a_t)$ – Transition probabilities
- $c_t(s_t, a_t, s_{t+1})$ – Cost function



Risk-Aware RL

Standard RL aims at minimising problems of the form: $\min_{\theta} \mathbb{E}[Y^{\theta}]$, where $Y^{\theta} = \sum_t c_t^{\theta}$

- **Ignores the risk** of the costs!

Risk-aware RL with static risk measures, e.g. expected utility [Nass et al., 2019], risk-constrained \mathbb{E} [Di Castro et al., 2019], coherent risk [Tamar et al., 2016], etc.

- Optimising **static risk** measures leads to optimal **precommitment policies**!

Recent approaches to overcome the time-consistency issue, e.g.:

- *Dynamic risk measures* [Marzban et al., 2021; Coache and Jaimungal, 2021], *conditional risk mappings* [Cheng and Jaimungal, 2022], *recursive risk filters* [Bielecki et al., 2022]...

In this paper, we:

- develop a computational approach to solve RL problems with dynamic risk
- devise an efficient deep estimation method for elicitable dynamic risk measures
- prove that these dynamic risk measures may be approximated to an arbitrary accuracy using NNs

Risk-Aware RL

Standard RL aims at minimising problems of the form: $\min_{\theta} \mathbb{E}[Y^{\theta}]$, where $Y^{\theta} = \sum_t c_t^{\theta}$

- Ignores the risk of the costs!

Risk-aware RL with static risk measures, e.g. expected utility [Nass et al., 2019], risk-constrained \mathbb{E} [Di Castro et al., 2019], coherent risk [Tamar et al., 2016], etc.

- Optimising static risk measures leads to optimal precommitment policies!

Recent approaches to overcome the time-consistency issue, e.g.:

- *Dynamic risk measures* [Marzban et al., 2021; Coache and Jaimungal, 2021], *conditional risk mappings* [Cheng and Jaimungal, 2022], *recursive risk filters* [Bielecki et al., 2022]...

In this paper, we:

- develop a **computational approach** to solve RL problems with **dynamic risk**
- devise an **efficient deep estimation method** for elicitable dynamic risk measures
- prove that these dynamic risk measures may be **approximated to an arbitrary accuracy** using NNs

Time-Consistent Dynamic Risk

- $\mathcal{Y}_{t_1, t_2} := \mathcal{Y}_{t_1} \times \cdots \times \mathcal{Y}_{t_2}$ – Sequence of \mathcal{F}_t -measurable random costs on $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_t, \mathbb{P})$
- $\rho_{t, T} : \mathcal{Y}_{t, T} \rightarrow \mathcal{Y}_t$ – Dynamic risk measure $\{\rho_{t, T}\}_t$
- **Strong time-consistency** – For any $Y, Z \in \mathcal{Y}_{t_1, T}$ and $0 \leq t_1 < t_2 \leq T$, we have

$$Y_k = Z_k, \forall k = t_1, \dots, t_2 - 1 \text{ and } \rho_{t_2, T}(Y_{t_2}, \dots, Y_T) \leq \rho_{t_2, T}(Z_{t_2}, \dots, Z_T)$$

implies that $\rho_{t_1, T}(Y_{t_1}, \dots, Y_T) \leq \rho_{t_1, T}(Z_{t_1}, \dots, Z_T)$.

[Thm. 1, Ruszczyński, 2010]

Let $\{\rho_{t, T}\}_t$ be a monotone, cash-additive, and normalised dynamic risk measure. Then $\{\rho_{t, T}\}_t$ is strongly time-consistent iff it may be expressed with one-step conditional risk measures $\rho_t : \mathcal{Y}_{t+1} \rightarrow \mathcal{Y}_t$ as

$$\rho_{t, T}(Y_t, \dots, Y_T) = Y_t + \rho_t \left(Y_{t+1} + \rho_{t+1} \left(Y_{t+2} + \cdots + \rho_{T-1}(Y_T) \cdots \right) \right).$$

Time-Consistent Dynamic Risk

- $\mathcal{Y}_{t_1, t_2} := \mathcal{Y}_{t_1} \times \cdots \times \mathcal{Y}_{t_2}$ – Sequence of \mathcal{F}_t -measurable random costs on $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_t, \mathbb{P})$
- $\rho_{t, T} : \mathcal{Y}_{t, T} \rightarrow \mathcal{Y}_t$ – Dynamic risk measure $\{\rho_{t, T}\}_t$
- *Strong time-consistency* – For any $Y, Z \in \mathcal{Y}_{t_1, T}$ and $0 \leq t_1 < t_2 \leq T$, we have

$$Y_k = Z_k, \forall k = t_1, \dots, t_2 - 1 \text{ and } \rho_{t_2, T}(Y_{t_2}, \dots, Y_T) \leq \rho_{t_2, T}(Z_{t_2}, \dots, Z_T)$$

implies that $\rho_{t_1, T}(Y_{t_1}, \dots, Y_T) \leq \rho_{t_1, T}(Z_{t_1}, \dots, Z_T)$.

[Thm. 1, [Ruszczyński, 2010](#)]

Let $\{\rho_{t, T}\}_t$ be a monotone, cash-additive, and normalised dynamic risk measure.

Then $\{\rho_{t, T}\}_t$ is strongly time-consistent iff it may be expressed with **one-step conditional risk measures** $\rho_t : \mathcal{Y}_{t+1} \rightarrow \mathcal{Y}_t$ as

$$\rho_{t, T}(Y_t, \dots, Y_T) = Y_t + \rho_t \left(Y_{t+1} + \rho_{t+1} \left(Y_{t+2} + \cdots + \rho_{T-1}(Y_T) \cdots \right) \right).$$

Problem Setup

Problems of the form

$$\min_{\theta} \rho_{0,T}(\{c_t^\theta\}_t) = \min_{\theta} \rho_0 \left(c_0^\theta + \rho_1 \left(c_1^\theta + \cdots + \rho_{T-1} \left(c_{T-1}^\theta + \rho_T(c_T^\theta) \right) \cdots \right) \right)$$

where c_t^θ are \mathcal{F}_{t+1} -measurable random costs and ρ_t 's are **static risk measures**.

DP equations for the *value function*, i.e. running risk-to-go, for $s \in \mathcal{S}$:

$$V_t(s; \theta) = \rho_t \left(\underbrace{c_t^\theta}_{\text{current cost}} + \underbrace{V_{t+1}(s_{t+1}^\theta; \theta)}_{\text{one-step ahead risk-to-go}} \mid s_t = s \right),$$

under transition probabilities $\mathbb{P}^\theta(a, s' | s_t = s) = \mathbb{P}(s' | s, a) \pi^\theta(a | s_t = s)$

Problem Setup

Problems of the form

$$\min_{\theta} \rho_{0,T}(\{c_t^{\theta}\}_t) = \min_{\theta} \rho_0 \left(c_0^{\theta} + \rho_1 \left(c_1^{\theta} + \cdots + \rho_{T-1} \left(c_{T-1}^{\theta} + \rho_T(c_T^{\theta}) \right) \cdots \right) \right)$$

where c_t^{θ} are \mathcal{F}_{t+1} -measurable random costs and ρ_t 's are static risk measures.

DP equations for the *value function*, i.e. running risk-to-go, for $s \in \mathcal{S}$:

$$V_t(s; \theta) = \rho_t \left(\underbrace{c_t^{\theta}}_{\text{current cost}} + \underbrace{V_{t+1}(s_{t+1}^{\theta}; \theta)}_{\text{one-step ahead risk-to-go}} \mid s_t = s \right),$$

under transition probabilities $\mathbb{P}^{\theta}(a, s' | s_t = s) = \mathbb{P}(s' | s, a) \pi^{\theta}(a | s_t = s)$

Policy Gradient

- We wish to optimise the value function over policies θ via a policy gradient method:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} V(\cdot; \theta)$$

[Gradient of V , Coache et al., 2022]

Under some regularity assumptions, the gradient of the value function at any period $t \in \mathcal{T}$ and any state $s \in \mathcal{S}$ for dynamic spectral risk measures with finite support spectrum is

$$\begin{aligned} \nabla_{\theta} V_t(s; \theta) = & \sum_{m=1}^{k-1} \frac{p_m}{1 - \alpha_m} \mathbb{E}_{\mathbb{P}^{\theta}(\cdot, \cdot | s_t = s)} \left[\left(c_t^{\theta} + V_{t+1}(s_{t+1}^{\theta}; \theta) - \lambda_m^* \right)_+ \left(\nabla_{\theta} \log \pi^{\theta}(a | s_t) \Big|_{a=a_t^{\theta}} \right) \right] \\ & + \mathbb{E}_{\mathbb{P}^{\theta}(\cdot, \cdot | s_t = s)} \left[\left(\nabla_{\theta} V_{t+1}(s'; \theta) \Big|_{s'=s_{t+1}^{\theta}} \right) \xi_m^*(a_t^{\theta}, s_{t+1}^{\theta}) \right], \end{aligned}$$

Actor-critic style algorithm composed of two interleaved procedures:

- Value function estimation given a policy
- Policy update given a value function

Policy Gradient

- We wish to optimise the value function over policies θ via a policy gradient method:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} V(\cdot; \theta)$$

[Gradient of V , Coache et al., 2022]

Under some regularity assumptions, the gradient of the value function at any period $t \in \mathcal{T}$ and any state $s \in \mathcal{S}$ for dynamic spectral risk measures with finite support spectrum is

$$\begin{aligned} \nabla_{\theta} V_t(s; \theta) = & \sum_{m=1}^{k-1} \frac{p_m}{1 - \alpha_m} \mathbb{E}_{\mathbb{P}^{\theta}(\cdot, \cdot | s_t = s)} \left[\left(c_t^{\theta} + V_{t+1}(s_{t+1}^{\theta}; \theta) - \lambda_m^* \right)_+ \left(\nabla_{\theta} \log \pi^{\theta}(a | s_t) \Big|_{a=a_t^{\theta}} \right) \right] \\ & + \mathbb{E}_{\mathbb{P}^{\theta}(\cdot, \cdot | s_t = s)} \left[\left(\nabla_{\theta} V_{t+1}(s'; \theta) \Big|_{s'=s_{t+1}^{\theta}} \right) \xi_m^*(a_t^{\theta}, s_{t+1}^{\theta}) \right], \end{aligned}$$

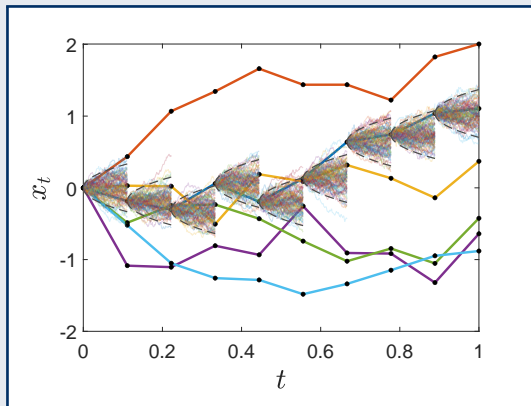
Actor-critic style algorithm composed of two interleaved procedures:

- Value function estimation given a policy
- Policy update given a value function

Estimation of V

Previous approaches: nested simulations [Tamar et al., 2016; Coache and Jaimungal, 2021]

- Generate (outer) episodes and (inner) transitions for every visited state
- Computationally expensive...



Elicitability

We leverage the elicibility to efficiently estimate dynamic risk measures

ρ is *elicitable* iff there exists a scoring function $S : \mathbb{R} \times \mathbb{Y} \rightarrow \mathbb{R}$ s.t.

$$\rho(Y) = \arg \min_{a \in \mathbb{R}} \mathbb{E}_{Y \sim F_Y} [S(a, Y)].$$

$\rho(Y)$	Mean	Median	VaR_α	CVaR_α
$S(a, y)$	$(a - y)^2$	$ a - y $	$\mathbb{1}_{a \leq y} - \alpha$	\emptyset

Non-elicitable mappings can be components of an elicitable vector-valued mapping:

- Elicitability of (static) spectral risk measures [Fissler and Ziegel, 2016]
- Characterization of their scoring function S

Elicitability

We leverage the elicibility to efficiently estimate dynamic risk measures

ρ is *elicitable* iff there exists a scoring function $S : \mathbb{R} \times \mathbb{Y} \rightarrow \mathbb{R}$ s.t.

$$\rho(Y) = \arg \min_{a \in \mathbb{R}} \mathbb{E}_{Y \sim F_Y} [S(a, Y)].$$

$\rho(Y)$	Mean	Median	VaR_α	CVaR_α
$S(a, y)$	$(a - y)^2$	$ a - y $	$\mathbb{1}_{a \leq y} - \alpha$	\emptyset

Non-elicitable mappings can be components of an elicitable vector-valued mapping:

- Elicitability of (static) spectral risk measures [Fissler and Ziegel, 2016]
- Characterization of their scoring function S

Conditional Elicitability

Example: $(\text{VaR}_\alpha(Y), \text{CVaR}_\alpha(Y))$ is elicitable, i.e.

$$(\text{VaR}_\alpha(Y), \text{CVaR}_\alpha(Y)) = \arg \min_{(\mathbf{a}_1, \mathbf{a}_2) \in \mathbb{R}^2} \mathbb{E}_{Y \sim F_Y} [S(\mathbf{a}_1, \mathbf{a}_2, Y)]$$

In our RL problem, the costs are supported by observed features, i.e. the states $s \in \mathcal{S}$

$$(\text{VaR}_\alpha(Y|s_t = s), \text{CVaR}_\alpha(Y|s_t = s)) = \arg \min_{h_1, h_2: \mathcal{S} \rightarrow \mathbb{R}} \mathbb{E}_{Y \sim F_Y} [S(h_1(s), h_2(s), Y)]$$

- Model $V_t(s; \theta)$ with NNs $H_t^\psi(s), V_t^\phi(s)$
- Use empirical estimates based on observed data

$$\arg \min_{\psi, \phi} \sum_{t \in \mathcal{T}} \sum_{i=1}^n S\left(\underbrace{H_t^\psi(s^{(i)})}_{\text{VaR}_\alpha}, \underbrace{V_t^\phi(s^{(i)})}_{\text{CVaR}_\alpha}, \underbrace{c_t^{(i)} + V_{t+1}^\phi(s_{t+1}^{(i)})}_{\text{random costs}}\right)$$

Similar results for a class of spectral risk measures

Conditional Elicitability

Example: $(\text{VaR}_\alpha(Y), \text{CVaR}_\alpha(Y))$ is elicitable, i.e.

$$(\text{VaR}_\alpha(Y), \text{CVaR}_\alpha(Y)) = \arg \min_{(a_1, a_2) \in \mathbb{R}^2} \mathbb{E}_{Y \sim F_Y} [S(a_1, a_2, Y)]$$

In our RL problem, the costs are supported by observed features, i.e. the states $s \in \mathcal{S}$

$$(\text{VaR}_\alpha(Y|s_t = s), \text{CVaR}_\alpha(Y|s_t = s)) = \arg \min_{h_1, h_2: \mathcal{S} \rightarrow \mathbb{R}} \mathbb{E}_{Y \sim F_Y} [S(h_1(s), h_2(s), Y)]$$

- Model $V_t(s; \theta)$ with NNs $H_t^\psi(s), V_t^\phi(s)$
- Use empirical estimates based on observed data

$$\arg \min_{\psi, \phi} \sum_{t \in \mathcal{T}} \sum_{i=1}^n S\left(\underbrace{H_t^\psi(s^{(i)})}_{\text{VaR}_\alpha}, \underbrace{V_t^\phi(s^{(i)})}_{\text{CVaR}_\alpha}, \underbrace{c_t^{(i)} + V_{t+1}^\phi(s_{t+1}^{(i)})}_{\text{random costs}}\right)$$

Similar results for a class of spectral risk measures

Conditional Elicitability

Example: $(\text{VaR}_\alpha(Y), \text{CVaR}_\alpha(Y))$ is elicitable, i.e.

$$(\text{VaR}_\alpha(Y), \text{CVaR}_\alpha(Y)) = \arg \min_{(a_1, a_2) \in \mathbb{R}^2} \mathbb{E}_{Y \sim F_Y} [S(a_1, a_2, Y)]$$

In our RL problem, the costs are supported by observed features, i.e. the states $s \in \mathcal{S}$

$$(\text{VaR}_\alpha(Y|s_t = s), \text{CVaR}_\alpha(Y|s_t = s)) = \arg \min_{h_1, h_2: \mathcal{S} \rightarrow \mathbb{R}} \mathbb{E}_{Y \sim F_Y} [S(h_1(s), h_2(s), Y)]$$

- Model $V_t(s; \theta)$ with NNs $H_t^\psi(s), V_t^\phi(s)$
- Use empirical estimates based on observed data

$$\arg \min_{\psi, \phi} \sum_{t \in \mathcal{T}} \sum_{i=1}^n S \left(\underbrace{H_t^\psi(s^{(i)})}_{\text{VaR}_\alpha}, \underbrace{V_t^\phi(s^{(i)})}_{\text{CVaR}_\alpha}, \underbrace{c_t^{(i)} + V_{t+1}^\phi(s_{t+1}^{(i)})}_{\text{random costs}} \right)$$

Similar results for a class of spectral risk measures

Accuracy of the Elicitable Approach

- We can approximate the value function to an arbitrary accuracy using this framework

[Approximation of V , Coache et al., 2022]

Suppose π^θ is a fixed policy, with its corresponding value function $V_t(s; \theta)$. Then for any $\varepsilon_1^*, \dots, \varepsilon_k^* > 0$, there exist NNs denoted $H_{1,t}^{\psi_1}, \dots, H_{k,t}^{\psi_k}$ such that for any $t \in \mathcal{T}$, we have

$$\operatorname{ess\,sup}_{s \in \mathcal{S}} \left\| V_t(s; \theta) - \left(H_{k,t}^{\psi_k}(s; \theta) + \sum_{m=1}^{k-1} p_m \sum_{l=1}^m H_{l,t}^{\psi_l}(s; \theta) \right) \right\| < \varepsilon^*.$$

Portfolio Allocation

Consider a market with d assets. An agent

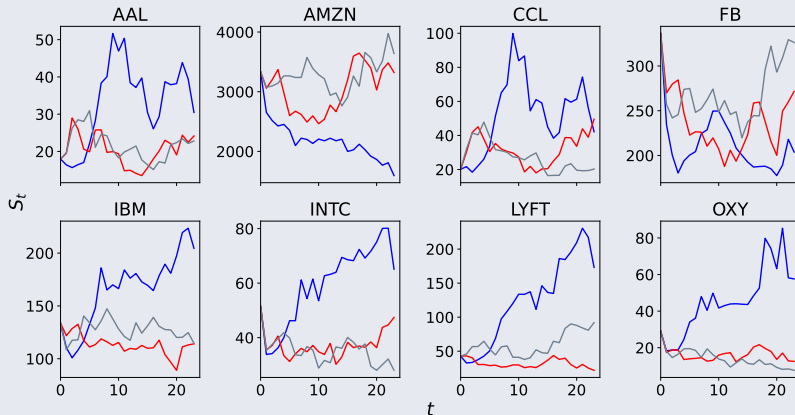
- observes the time t and asset prices $\{S_t^{(i)}\}_{i=1,\dots,d}$
- decides on the proportion of its wealth $\pi_t^{(i)}$ to invest in asset i
- receives feedback from P&L differences $y_t - y_{t+1}$, where its wealth y_t varies according to

$$dy_t = y_t \left(\sum_{i=1}^d \pi_t^{(i)} \frac{dS_t^{(i)}}{S_t^{(i)}} \right), \quad y_0 = 1.$$

We assume a null interest rate, no leveraging nor short-selling.

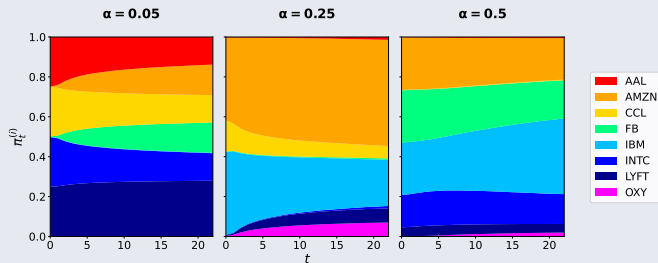
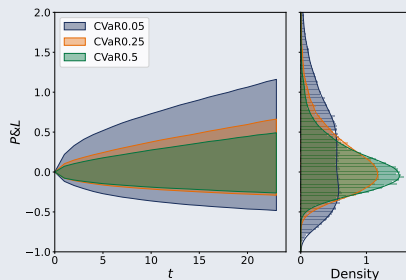
Portfolio Allocation

Co-integration model using daily data from eight different stocks listed on the NASDAQ exchange between September 31, 2020 and December 31, 2021.



Portfolio Allocation

Co-integration model using daily data from eight different stocks listed on the NASDAQ exchange between September 31, 2020 and December 31, 2021.



Contributions

A **practical, flexible framework** for risk-aware **RL with dynamic risk measures**

- Novel setting utilising *elicitability* for efficient & accurate estimation
- Performance validation on several benchmark optimisation problems

Future directions

- Robustification to protect against model uncertainty
- DDPG approach for dynamic risk measures
- Risk-aware dynamic RL for multi-agent systems

Thank you!

Paper, code and slides available at: anthonycoache.ca

References

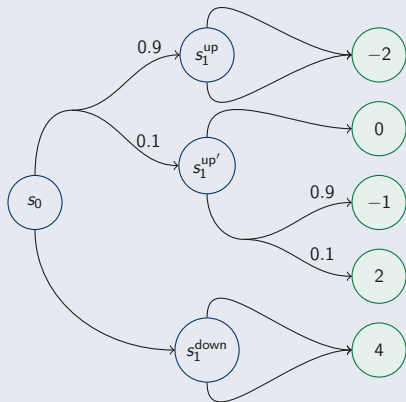
- Bielecki, T. R., Cialenco, I., and Ruszczyński, A. (2022). Risk filtering and risk-averse control of Markovian systems subject to model uncertainty. *arXiv preprint arXiv:2206.09235*.
- Cheng, Z. and Jaimungal, S. (2022). Markov decision processes with Kusuoka-type conditional risk mappings. *arXiv preprint arXiv:2203.09612*.
- Coache, A. and Jaimungal, S. (2021). Reinforcement learning with dynamic convex risk measures. *arXiv preprint arXiv:2112.13414*.
- Coache, A., Jaimungal, S., and Cartea, Á. (2022). Conditionally elicitable dynamic risk measures for deep reinforcement learning. *arXiv preprint arXiv:2206.14666*.
- Di Castro, D., Oren, J., and Mannor, S. (2019). Practical risk measures in reinforcement learning. *arXiv preprint arXiv:1908.08379*.
- Fissler, T. and Ziegel, J. F. (2016). Higher order elicitability and Osband's principle. *The Annals of Statistics*, 44(4):1680–1707.
- Marzban, S., Delage, E., and Li, J. Y. (2021). Deep reinforcement learning for equal risk pricing and hedging under dynamic expectile risk measures. *arXiv preprint arXiv:2109.04001*.
- Nass, D., Belousov, B., and Peters, J. (2019). Entropic risk measure in policy search. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1101–1106. IEEE.
- Ruszczyński, A. (2010). Risk-averse dynamic programming for Markov decision processes. *Mathematical Programming*, 125(2):235–261.
- Tamar, A., Chow, Y., Ghavamzadeh, M., and Mannor, S. (2016). Sequential decision making with coherent risk. *IEEE Transactions on Automatic Control*, 62(7):3323–3338.

Time-Consistency Issue...

Let us minimize $\text{CVaR}_{0.9}$ of the terminal cost.

- *Optimal actions at s_0 : Move up, then down*
- *Optimal actions at $s_1^{up'}$: Move up*

Contradiction with the initial optimal strategy...



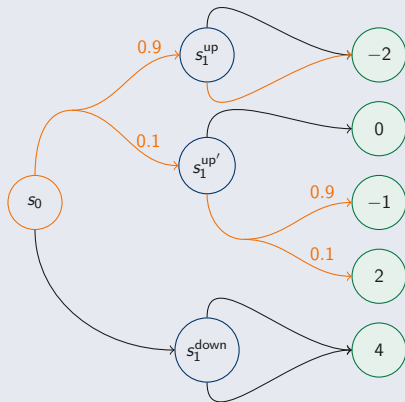
Optimizing static risk measures leads to optimal precommitment policies!

Time-Consistency Issue...

Let us minimize $\text{CVaR}_{0.9}$ of the terminal cost.

- *Optimal actions at s_0* : Move up, then down
- *Optimal actions at $s_1^{up'}$* : Move up

Contradiction with the initial optimal strategy...



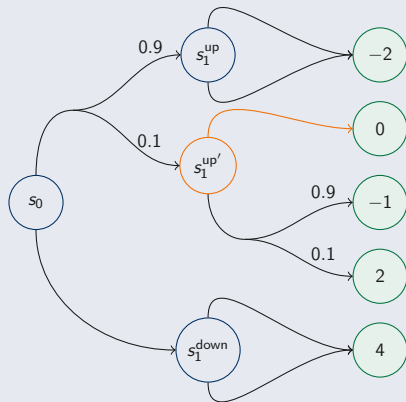
Optimizing static risk measures leads to optimal precommitment policies!

Time-Consistency Issue...

Let us minimize $\text{CVaR}_{0.9}$ of the terminal cost.

- *Optimal actions at s_0* : Move up, then down
- *Optimal actions at $s_1^{up'}$* : Move up

Contradiction with the initial optimal strategy...



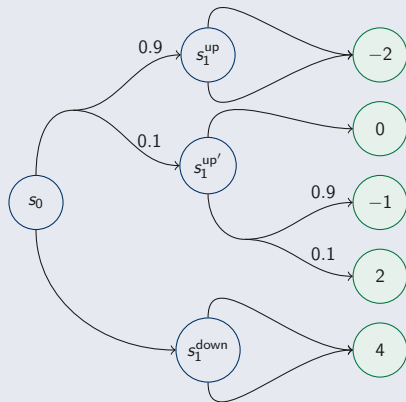
Optimizing static risk measures leads to optimal precommitment policies!

Time-Consistency Issue...

Let us minimize $\text{CVaR}_{0.9}$ of the terminal cost.

- *Optimal actions at s_0* : Move up, then down
- *Optimal actions at $s_1^{up'}$* : Move up

Contradiction with the initial optimal strategy...



Optimizing **static risk** measures leads to optimal **precommitment policies**!

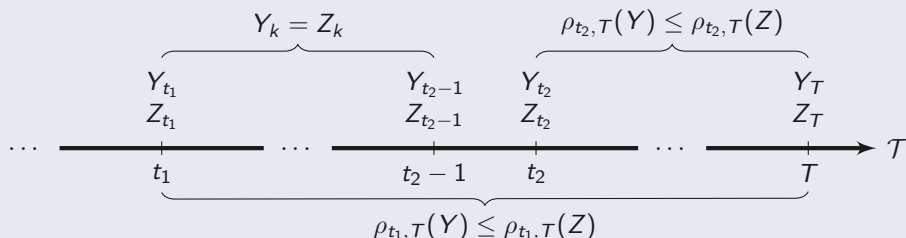
Time-Consistency

Strong time-consistency

$\{\rho_{t,T}\}_t$ is *strongly time-consistent* iff for any $Y, Z \in \mathcal{Y}_{t_1,T}$ and $0 \leq t_1 < t_2 \leq T$, we have

$$Y_k = Z_k, \forall k = t_1, \dots, t_2 - 1 \text{ and } \rho_{t_2,T}(Y_{t_2}, \dots, Y_T) \leq \rho_{t_2,T}(Z_{t_2}, \dots, Z_T)$$

implies that $\rho_{t_1,T}(Y_{t_1}, \dots, Y_T) \leq \rho_{t_1,T}(Z_{t_1}, \dots, Z_T)$.



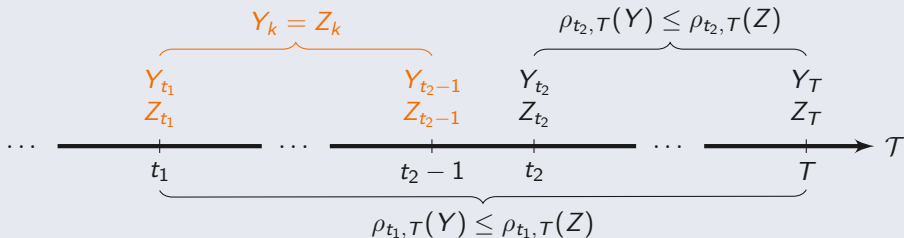
Time-Consistency

Strong time-consistency

$\{\rho_{t,T}\}_t$ is *strongly time-consistent* iff for any $Y, Z \in \mathcal{Y}_{t_1,T}$ and $0 \leq t_1 < t_2 \leq T$, we have

$$Y_k = Z_k, \forall k = t_1, \dots, t_2 - 1 \text{ and } \rho_{t_2,T}(Y_{t_2}, \dots, Y_T) \leq \rho_{t_2,T}(Z_{t_2}, \dots, Z_T)$$

implies that $\rho_{t_1,T}(Y_{t_1}, \dots, Y_T) \leq \rho_{t_1,T}(Z_{t_1}, \dots, Z_T)$.



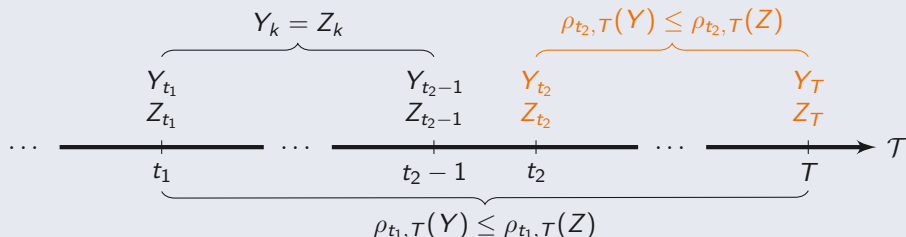
Time-Consistency

Strong time-consistency

$\{\rho_{t,T}\}_t$ is *strongly time-consistent* iff for any $Y, Z \in \mathcal{Y}_{t_1,T}$ and $0 \leq t_1 < t_2 \leq T$, we have

$$Y_k = Z_k, \forall k = t_1, \dots, t_2 - 1 \text{ and } \rho_{t_2,T}(Y_{t_2}, \dots, Y_T) \leq \rho_{t_2,T}(Z_{t_2}, \dots, Z_T)$$

implies that $\rho_{t_1,T}(Y_{t_1}, \dots, Y_T) \leq \rho_{t_1,T}(Z_{t_1}, \dots, Z_T)$.



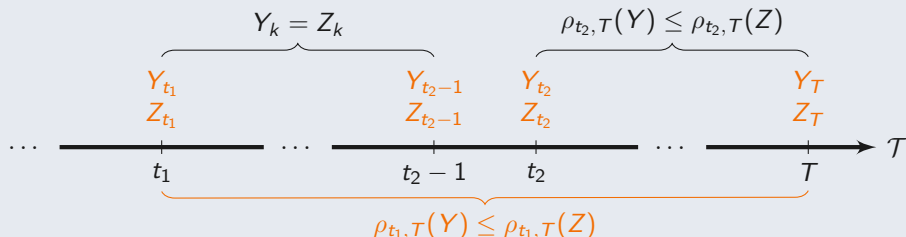
Time-Consistency

Strong time-consistency

$\{\rho_{t,T}\}_t$ is *strongly time-consistent* iff for any $Y, Z \in \mathcal{Y}_{t_1,T}$ and $0 \leq t_1 < t_2 \leq T$, we have

$$Y_k = Z_k, \forall k = t_1, \dots, t_2 - 1 \text{ and } \rho_{t_2,T}(Y_{t_2}, \dots, Y_T) \leq \rho_{t_2,T}(Z_{t_2}, \dots, Z_T)$$

implies that $\rho_{t_1,T}(Y_{t_1}, \dots, Y_T) \leq \rho_{t_1,T}(Z_{t_1}, \dots, Z_T)$.



Algorithms

Algorithm 1: Actor-critic algorithm – Elicitable approach

Input: NNs π^θ , V^ϕ , numbers of epochs K 's, mini-batch sizes B 's

```

1 Set initial learning rates for  $\phi, \theta$ ;
2 for each iteration  $k = 1, \dots, K$  do
3   for each epoch  $k^\phi = 1, \dots, K^\phi$  do
4     Simulate a mini-batch of  $B^\phi$  episodes induced by  $\pi^\theta$ ;
5     Compute the loss  $\mathcal{L}(\phi)$ : minimization of the expected consistent score;
6     Update  $\phi$  by performing an Adam optimisation step, tune the learning rate for  $\phi$ ;
7     if  $k^\phi \bmod K^* = 0$  then
8       | Update the target networks  $\tilde{\phi}$ ;
9   for each epoch  $k^\theta = 1, \dots, K^\theta$  do
10    Simulate a mini-batch of  $\lceil B^\theta / (1 - \alpha) \rceil$  episodes induced by  $\pi^\theta$ ;
11    Compute the loss  $\mathcal{L}(\theta)$ : policy gradient;
12    Update  $\theta$  by performing an Adam optimisation step, tune the learning rate for  $\theta$ ;

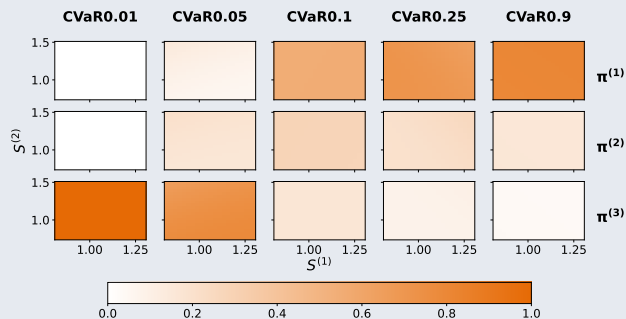
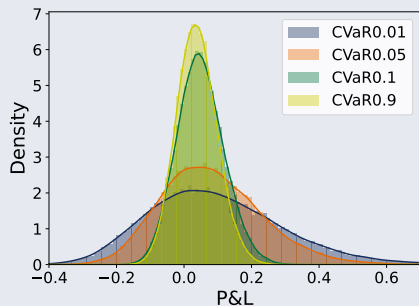
```

Output: Optimal policy π^θ and its value function V^ϕ

Portfolio Allocation

$$dS_t^{(i)} = \mu^{(i)} S_t^{(i)} dt + \sigma^{(i)} S_t^{(i)} dW_t^{(i)}$$

Drifts and volatilities are $\mu = [0.03; 0.06; 0.09]$ and $\sigma = [0.06; 0.12; 0.18]$



Portfolio Allocation

$$dX_t^{(i)} = -\kappa X_t^{(i)} dt + \sigma^{(i)} dW_t^{(i)} \quad \text{with} \quad S_t^{(i)} = e^{X_t^{(i)} + \mu^{(i)} t - (\sigma^{(i)})^2 \frac{1 - e^{-2\kappa t}}{4\kappa}}$$

Drifts and volatilities are $\mu = [0.03; 0.06; 0.09]$ and $\sigma = [0.06; 0.12; 0.18]$

