

Risk-Averse Distributional Reinforcement Learning: Bonus Materials

May 2019

1 Proofs

Proof of Lemma 1

Proof. The fact that discrete distributions have a piecewise linear $y\text{CVaR}_y$ function has already been shown by [Rockafellar and Uryasev, 2000].

According to definition (2) we have

$$y\text{CVaR}_y(Z) = y \frac{1}{y} \int_0^y \text{VaR}_\beta(Z) d\beta = \int_0^y \text{VaR}_\beta(Z) d\beta$$

by taking the y derivative, we have

$$\frac{\partial}{\partial y} y\text{CVaR}_y(Z) = \frac{\partial}{\partial y} \int_0^y \text{VaR}_\beta(Z) d\beta = \text{VaR}_y(Z)$$

□

Proof of Theorem 1

Proof. Since we are interested in the minimal argument, we can ease the computation by focusing on the αCVaR_α function instead of CVaR_α . When working with two states, the equation of interest simplifies to

$$\begin{aligned} \alpha\text{CVaR}_\alpha(Z(x, a)) &= \min_{\xi} p_1 \xi_1 \alpha\text{CVaR}_{\xi_1 \alpha}(Z(x'_1)) + p_2 \xi_2 \alpha\text{CVaR}_{\xi_2 \alpha}(Z(x'_2)) \\ \text{s.t. } & p_1 \xi_1 + p_2 \xi_2 = 1 \\ & 0 \leq \xi_1 \leq \frac{1}{\alpha} \\ & 0 \leq \xi_2 \leq \frac{1}{\alpha} \end{aligned}$$

therefore

$$\begin{aligned}\alpha\text{CVaR}_\alpha(Z(x, a)) &= \min_{\xi} p_1 \xi_1 \alpha \text{CVaR}_{\xi_1 \alpha}(Z(x'_1)) + (1 - p_1) \frac{1 - p_1 \xi_1}{1 - p_1} \alpha \text{CVaR}_{\frac{1 - p_1 \xi_1}{1 - p_1} \alpha}(Z(x'_2)) \\ &= \min_{\xi} p_1 \int_0^{\xi_1 \alpha} \text{VaR}_{\beta}(Z(x'_1)) d\beta + (1 - p_1) \int_0^{\frac{1 - p_1 \xi_1}{1 - p_1} \alpha} \text{VaR}_{\beta}(Z(x'_2)) d\beta\end{aligned}$$

To find the minimal argument, we find the first derivative w.r.t. ξ_1

$$\begin{aligned}\frac{\partial \alpha \text{CVaR}_\alpha}{\partial \xi_1} &= p_1 \alpha \text{VaR}_{\xi_1 \alpha}(Z(x'_1)) + (1 - p_1) \alpha \frac{-p_1}{1 - p_1} \text{VaR}_{\frac{1 - p_1 \xi_1}{1 - p_1} \alpha}(Z(x'_2)) \\ &= p_1 \text{VaR}_{\xi_1 \alpha}(Z(x'_1)) - p_1 \text{VaR}_{\frac{1 - p_1 \xi_1}{1 - p_1} \alpha}(Z(x'_2))\end{aligned}$$

By setting the derivative to 0, we get

$$\text{VaR}_{\xi_1 \alpha}(Z(x'_1)) \stackrel{!}{=} \text{VaR}_{\frac{1 - p_1 \xi_1}{1 - p_1} \alpha}(Z(x'_2)) = \text{VaR}_{\xi_2 \alpha}(Z(x'_2))$$

[Bernard and Vanduffel, 2015] have shown that in the case of strictly increasing c.d.f. with unbounded support, it holds that

$$\begin{aligned}\text{VaR}_{\xi_1 \alpha}(Z(x'_1)) &= \text{VaR}_{\xi_2 \alpha}(Z(x'_2)) \\ &= \text{VaR}_\alpha(Z(x, a)) \\ F_{Z(x'_1)}^{-1}(\xi_1 \alpha) &= F_{Z(x'_2)}^{-1}(\xi_2 \alpha) \\ &= F_{Z(x, a)}^{-1}(\alpha)\end{aligned}$$

and we can extract the values of $\xi_1 \alpha, \xi_2 \alpha$ using the

$$\begin{aligned}F_{Z(x'_1)}^{-1}(\xi_1 \alpha) &= F_{Z(x, a)}^{-1}(\alpha) & / F_{Z(x'_1)} \\ F_{Z(x'_1)} \left(F_{Z(x'_1)}^{-1}(\xi_1 \alpha) \right) &= F_{Z(x'_1)} \left(F_{Z(x, a)}^{-1}(\alpha) \right) \\ \xi_1 \alpha &= F_{Z(x'_1)} \left(F_{Z(x, a)}^{-1}(\alpha) \right)\end{aligned}$$

And similarly for ξ_2 .

Since the problem is convex, we have found the optimal point. \square

Proof of Theorem 2

Proof. Let s^* be a solution to $\max_s \frac{1}{\alpha} \mathbb{E}[(Z^\pi(x_0) - s)^-] + s$. Then by optimizing $\frac{1}{\alpha} \mathbb{E}[(Z^\pi - s^*)^-]$ over π , we monotonously improve the optimization criterion $\text{CVaR}_\alpha(Z(x_0))$.

$$\begin{aligned}\text{CVaR}_\alpha(Z^\pi) &= \max_s \frac{1}{\alpha} \mathbb{E}[(Z^\pi - s)^-] + s &= \frac{1}{\alpha} \mathbb{E}[(Z^\pi - s^*)^-] + s^* \\ &\leq \max_{\pi'} \frac{1}{\alpha} \mathbb{E}[(Z^{\pi'} - s^*)^-] + s^* &= \frac{1}{\alpha} \mathbb{E}[(Z^{\pi^*} - s^*)^-] + s^* \\ &\leq \max_{s'} \frac{1}{\alpha} \mathbb{E}[(Z^{\pi^*} - s')^-] + s' &= \text{CVaR}_\alpha(Z^{\pi^*})\end{aligned}$$

When optimizing w.r.t. π we can ignore the scaling term $\frac{1}{\alpha}$ and a constant term s^* without affecting the optimal argument. We can therefore focus on optimization of $\mathbb{E}[(Z^\pi(x_0) - s^*)^-]$.

$$\begin{aligned}
\mathbb{E}[(Z_t - s)^-] &= \mathbb{E}[(Z_t - s)\mathbb{1}(Z_t \leq s)] = \mathbb{E}\left[(r_t + \gamma Z_{t+1} - s)\mathbb{1}(Z_{t+1} \leq \frac{s - r_t}{\gamma})\right] \\
&= \sum_{x_{t+1}, r_t} P(x_{t+1}, r_t | x_t, a) \mathbb{E}\left[(r_t + \gamma Z(x_{t+1}) - s)\mathbb{1}(Z(x_{t+1}) \leq \frac{s - r_t}{\gamma})\right] \\
&= \sum_{x_{t+1}, r_t} P(x_{t+1}, r_t | x_t, a) \mathbb{E}\left[\gamma \left(Z(x_{t+1}) - \frac{s - r_t}{\gamma}\right) \mathbb{1}(Z(x_{t+1}) \leq \frac{s - r_t}{\gamma})\right] \\
&= \gamma \sum_{x_{t+1}, r_t} P(x_{t+1}, r_t | x_t, a) \mathbb{E}\left[\left(Z(x_{t+1}) - \frac{s - r_t}{\gamma}\right) \mathbb{1}(Z(x_{t+1}) \leq \frac{s - r_t}{\gamma})\right] \\
&= \gamma \sum_{x_{t+1}, r_t} P(x_{t+1}, r_t | x_t, a) \mathbb{E}\left[\left(Z(x_{t+1}) - \frac{s - r_t}{\gamma}\right)^-\right]
\end{aligned} \tag{1}$$

where we used the definition of return $Z_t = R_t + \gamma Z_{t+1}$ and the fact that probability mixture expectations can be computed as $\mathbb{E}[f(Z)] = \sum_i p_i \mathbb{E}[f(Z_i)]$ for any function f .

Now let's say we sampled reward r_t and state x_{t+1} , we are still trying to find a policy π^* that maximizes

$$\begin{aligned}
\pi^* &= \arg \max_{\pi} \mathbb{E}[(Z(x_t) - s)^- | x_{t+1}, r_t] \\
&= \arg \max_{\pi} \mathbb{E}\left[\left(Z(x_{t+1}) - \frac{s - r_t}{\gamma}\right)^-\right]
\end{aligned} \tag{2}$$

Where we ignored the unsampled states, since these are not a function of x_{t+1} , and the multiplicative constant γ that will not affect the maximum argument.

At the starting state, we set $s = s^*$. At each following state we select an action according to equation (2). By induction we maximize the criterion (??) in each step. \square

2 Algorithms

Algorithm 1 CVaR Computation via Quantile Representation

| | |
|--|--|
| <pre> function extractDistribution input: vectors C, y # Note: $y_0 = C(x', y_0) = 0$ for $i \in \{1, \dots, \mathbf{y} \}$ do $d_i = \frac{C(x', y_i) - C(x', y_{i-1})}{y_i - y_{i-1}}$ end for output vector d </pre> | <pre> function extractC input: vectors d, p $C_0 = 0$ for $i \in \{1, \dots, \mathbf{p} \}$ do $C_i = C_{i-1} + d_i \cdot p_i$ end for output vector C </pre> |
|--|--|

```

function mixDistributions
  input: tuples  $(\mathbf{d}^{(1)}, p^{(1)}), \dots, (\mathbf{d}^{(K)}, p^{(K)})$  and vector y
  #  $\sum_{k=1}^K p_k = 1$ 
  for  $i, k \in \{1, \dots, K\} \times \{1, \dots, |\mathbf{y}|\}$  do
    # Weigh atom probabilities by transitions
     $p_i^{(k)} = p^{(k)} \cdot (y_i - y_{i-1})$ 
  end for
  # Join all tuples together:
   $atoms = \{(d_1^{(1)}, p_1^{(1)}), \dots, (d_N^{(1)}, p_N^{(1)}), (d_1^{(2)}, p_1^{(2)}), \dots, (d_N^{(K)}, p_N^{(K)})\}$ 
  Sort atoms by d
  Unwrap vectors d, p from sorted tuples
  output d, p

```

```

# Main
input: tuples  $(C(x'_i, \cdot), p^{(1)}), \dots, (C(x'_i, \cdot), p^{(K)})$  and vector y
for  $i \in \{1, \dots, K\}$  do
   $\mathbf{d}^{(i)} = \text{extractDistribution}(C(x'_i, \cdot), \mathbf{y})$ 
end for
 $\mathbf{d}_{\text{mix}}, \mathbf{y}_{\text{mix}} = \text{mixDistributions}((\mathbf{d}^{(1)}, p^{(1)}), \dots, (\mathbf{d}^{(K)}, p^{(K)}), \mathbf{y})$ 
 $\mathbf{C}_{\text{out}} = \text{extractC}(\mathbf{d}_{\text{mix}}, \mathbf{y}_{\text{mix}})$ 
output:  $\mathbf{C}_{\text{out}}$ 

```

Algorithm 2 CVaR Q-learning policy

input: α , converged V, C
 $x = x_0$
 $a = \arg \max_a C(x, a, \alpha)$
 $s = V(x, a, y)$
while x is not terminal **do**
 $\mathbf{d}_a = \text{extractDistribution}(C(x', a, \cdot), \mathbf{y})$ for each a
 $a = \arg \max_a \text{expMinInterp}(s, \mathbf{d}, V(x', a, \cdot), \mathbf{y})$
 Take action a , observe r, x'
 $s = \frac{s - r}{\gamma}$
 $x = x'$
end while

Compute $\mathbb{E}[(\mathbf{d}_a - s)^-]$ with linear interpolation

function expMinInterp
 input: s , vectors $\mathbf{d}, \mathbf{V}, \mathbf{y}$
 $z = 0$
 for $i \in \{1, \dots, |\mathbf{y}|\}$ **do**
 if $s < V_i$ **then**
 break
 end if
 $z = z + d_i \cdot (y_i - y_{i-1})$
 end for
 $p_{\text{last}} = \frac{s - V_{i-1}}{V_i - V_{i-1}}(y_i - y_{i-1})$
 $z = z + d_i \cdot p_{\text{last}}$
 output z

Algorithm 3 Deep CVaR Q-learning with experience replay

```
Initialize replay memory  $M$ 
Initialize the VaR function  $V$  with random weights  $\theta_v$ 
Initialize the CVaR function  $C$  with random weights  $\theta_c$ 
Initialize target CVaR function  $C'$  with weights  $\theta'_c = \theta_c$ 
for each episode do
   $x = x_0$ 
  while  $x$  is not terminal do
    Choose  $a$  using a policy derived from  $C$  ( $\epsilon$ -greedy)
    Take action  $a$ , observe  $r, x'$ 
    Store transition  $(x, a, r, x')$  in  $M$ 
     $x = x'$ 
    Sample random transitions  $(x_j, a_j, r_j, x'_j)$  from  $M$ 
    Build the loss function  $\mathcal{L}_{\text{VaR}} + \mathcal{L}_{\text{CVaR}}$  (Algorithm ??)
    Perform a gradient step on  $\mathcal{L}_{\text{VaR}} + \mathcal{L}_{\text{CVaR}}$  w.r.t.  $\theta_v, \theta_c$ 
    Every  $N_{\text{target}}$  steps set  $\theta'_c = \theta_c$ 
  end while
end for
```

3 Figures

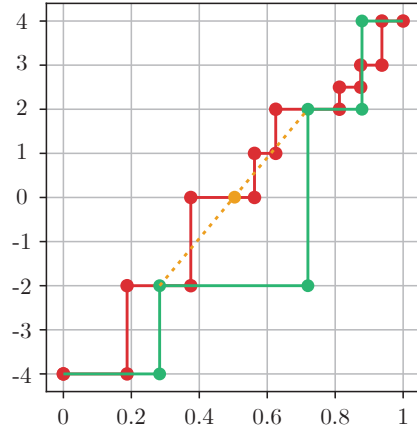


Figure 1: Visualization of the VaR-based heuristic. Quantile function of the exact distribution (unknown to the model) is shown in red and the VaR estimates at selected α -levels are shown in green. Let's say we now want to know y where $\text{VaR}_y = 0$. We use linear interpolation between the nearest known VaRs, shown in orange. In this case the interpolation estimate is $y = 0.5$.

References

- Bernard, C. and Vanduffel, S. (2015). Quantile of a mixture with application to model risk assessment. *Dependence Modeling*, 3(1).
- Rockafellar, R. T. and Uryasev, S. (2000). Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42.