

# Risk-averse Distributional Reinforcement Learning

## A CVaR optimization approach

Silvestr Stanko<sup>1</sup>

<sup>1</sup>Department of Computer Science  
Czech Technical University

Thursday 31<sup>st</sup> May, 2018

# Outline

- 1 Introduction
  - Motivation
  - Reinforcement Learning
  - Risk
  - Risk-averse Reinforcement Learning
- 2 CVaR Value Iteration
  - Previous results
  - Linear-time improvement
- 3 CVaR Q-learning
  - Var-based policy improvement
- 4 Deep CVaR Q-learning

# Outline

- 1 Introduction
  - Motivation
  - Reinforcement Learning
  - Risk
  - Risk-averse Reinforcement Learning
- 2 CVaR Value Iteration
  - Previous results
  - Linear-time improvement
- 3 CVaR Q-learning
  - Var-based policy improvement
- 4 Deep CVaR Q-learning

# Ultimate goals of AI

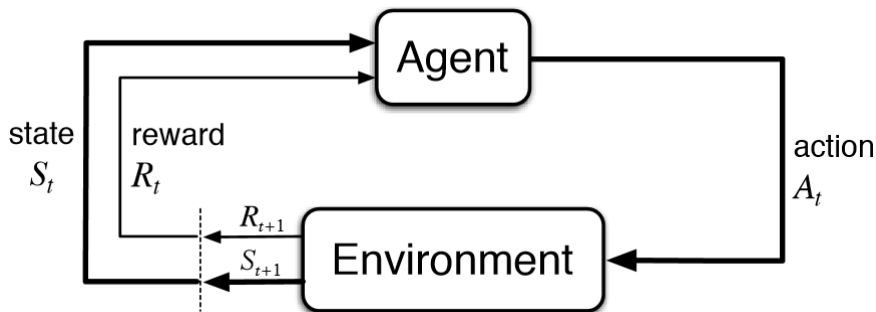
## General AI

- Learning from experience
- Learning *tabula rasa*
- Beyond purpose-specific AI
- Beyond human-level performance

## Safe AI

- Avoiding catastrophic events
- Robust to environment changes or adversaries

# Reinforcement Learning



# Recent successes

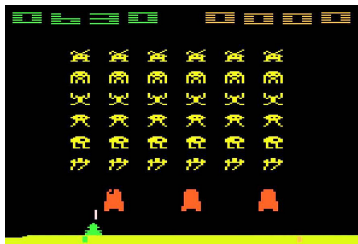


Figure: Atari games

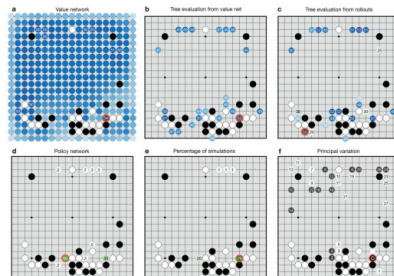


Figure: AlphaGo

# Recent successes



Figure: LOXM - online trader

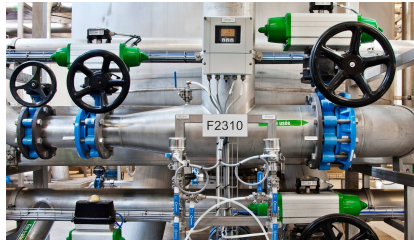


Figure: Google data centers: cooling

- Advertising recommendations (Microsoft)
- Dialog systems (Siri)

# Markov Decision Process - example



# Reinforcement Learning - Goal

## Definition

$Z^\pi(x_t)$  Is a random variable representing the discounted reward along a trajectory generated by the MDP by following the policy  $\pi$ , starting at state  $x_t$ .

$$Z^\pi(x_t) = \sum_{t=0}^{\infty} \gamma^t R(x_t, \pi(x_t))$$

## Reinforcement Learning goals

Our goal is to find a globally optimal policy  $\pi^*$

$$\pi^* = \arg \max_{\pi} \mathbb{E} Z^\pi(x_0)$$

# Potential problems

## Problems

- RL is sample inefficient  $\rightarrow$  expensive training
- Solutions must be **robust** to small model changes
- Solutions must avoid catastrophic events and be **safe**

## Solution

Instead of maximizing the expected reward, focus on other criteria that take into account the **risk** of the potential reward.

# Motivation



Figure: Simulation vs Real world



Figure: Critical Applications

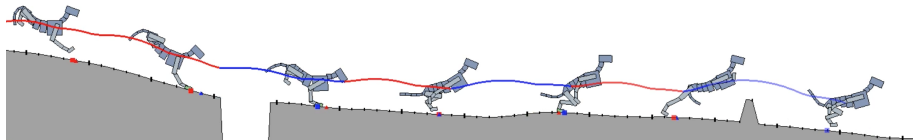


Figure: AI safety

# Risk

## Definition

Risk is the potential of gaining or losing something of value.

**Risk-averse:** disinclined or reluctant to take risks

Risk-neutral: indifferent to or balanced with respect to risk.

Risk-seeking: inclined or eager to take risks

## Example

Choose between paying:

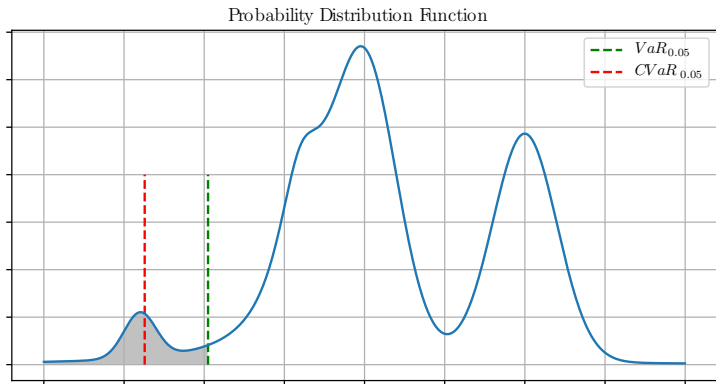
- ① \$100 in 100% cases
- ② \$1,000 in 10% cases and \$0 in 90% cases

# Risk-averse Reinforcement Learning - example

Figure: Greedy agent

Figure: Risk-averse agent

# Value-at-Risk, Conditional Value-at-Risk



# Risk-averse Reinforcement Learning - goals

## Definition

$Z^\pi(x_t)$  is a random variable representing the discounted reward along a trajectory generated by the MDP by following the policy  $\pi$ , starting at state  $x_t$ .

$$Z^\pi(x_t) = \sum_{t=0}^{\infty} \gamma^t R(x_t, \pi(x_t))$$

## Reinforcement Learning with CVaR

For a given  $\alpha$ , our goal is to find a globally optimal policy  $\pi^*$

$$\pi^* = \arg \max_{\pi} \text{CVaR}_{\alpha}(Z^\pi(x_0))$$

# Original Contributions

## 1 Faster CVaR Value Iteration

- Polynomial  $\rightarrow$  linear time.
- Formally proved for increasing, unbounded distributions.
- Experimentally verified for general distributions.

## 2 CVaR Q-learning

- Sampling version of CVaR Value Iteration.
- Based on the distributional approach.
- Experimentally verified.

## 3 Distributional Policy improvement

- Proved monotonic improvement for distributional RL.
- Used as a heuristic for extracting  $\pi^*$  from CVaR Q-learning.

## 4 Deep CVaR Q-learning

- TD update  $\rightarrow$  loss function.
- Experimentally verified in a deep learning context.



# Outline

- 1 Introduction
  - Motivation
  - Reinforcement Learning
  - Risk
  - Risk-averse Reinforcement Learning
- 2 CVaR Value Iteration
  - Previous results
  - Linear-time improvement
- 3 CVaR Q-learning
  - Var-based policy improvement
- 4 Deep CVaR Q-learning

# Value Iteration

## Definition

Value function  $V(x)$  represents the expected return when starting in state  $x$  and following the optimal policy  $\pi^*$  thereafter.

## Value Iteration

Initialize  $V_0(x)$  for each state (arbitrary value, e.g. 0).

Update each state:

$$V_{k+1}(x) = \max_a \left[ R(x, a) + \gamma \sum_{x'} p(x'|x, a) V_k(x') \right]$$

Repeat.

The algorithm converges to the optimal policy  $\pi^*$ :  $\lim_{k \rightarrow \infty} V_k(x) = V(x)$

# CVaR Value Iteration

## Theorem (CVaR Value Iteration)

*The following Bellman operator is a contraction:*

$$\mathbf{T}C(x, y) = \max_a \left[ R(x, a) + \gamma \min_{\xi} \sum_{x'} p(x'|x, a) \xi(x') C(x', y\xi(x')) \right]$$

The operator  $\mathbf{T}$  describes the following relationship:

$$\mathbf{T}CVaR_y(Z(x)) = \max_a \left[ R(x, a) + \gamma CVaR_y(Z(x')) \right]$$

$$x' \sim p(\cdot|x, a)$$

# CVaR VI computational complexity

## Problem

- CVaR Value Iteration requires computing a Linear Program for each state and atom.
- LP computation is slow

## Solution

CVaR Value Iteration with quantile representation.

# $\alpha\text{CVaR}_\alpha$ describes a quantile function

## Lemma

*Any discrete distribution has a piece-wise linear  $\alpha\text{CVaR}_\alpha$  function. Similarly, any a piece-wise linear  $\alpha\text{CVaR}_\alpha$  function can be seen as representing a certain discrete distribution.*

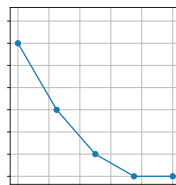
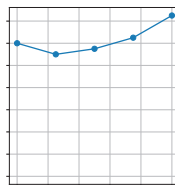
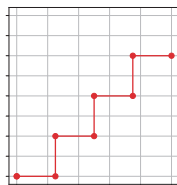
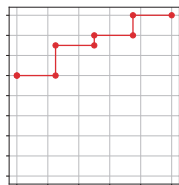
$$\alpha\text{CVaR}_\alpha \Rightarrow \text{VaR}$$

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

$$\alpha\text{CVaR}_\alpha \Leftarrow \text{VaR}$$

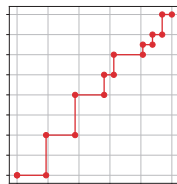
$$\alpha\text{CVaR}_\alpha(Z) = \int_0^\alpha \text{VaR}_\beta(Z) d\beta$$

# Next state CVaR computation



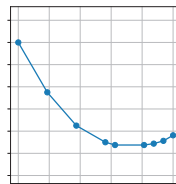
0.25

0.75

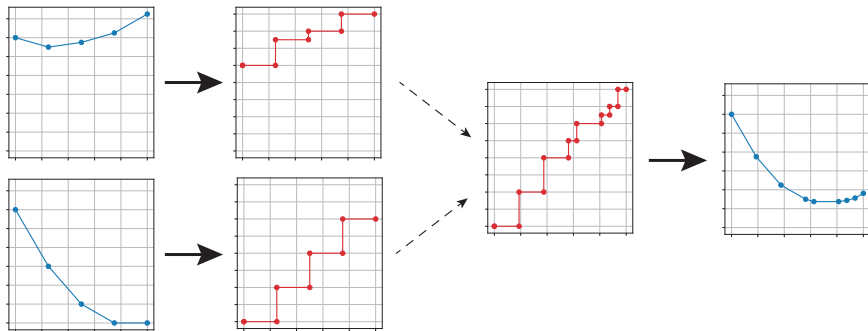


0.25

0.75



# Next state CVaR computation



# Linear-time Computation

## Theorem

*Solution to minimization problem*

$$\min_{\xi \in \mathcal{U}_{CVaR}(\alpha, p(\cdot|x, a))} \sum_{x'} p(x'|x, a) \xi(x') CVaR_{\xi(x')\alpha} (Z^\pi(x'))$$

*can be computed by setting*

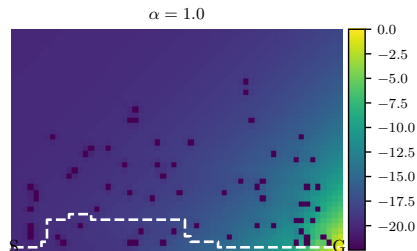
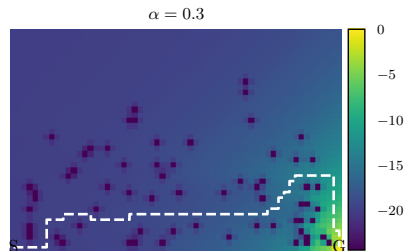
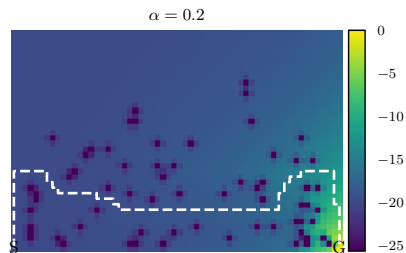
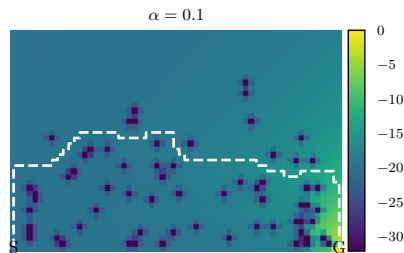
$$\xi(x') = \frac{F_{Z(x')}(F_{Z(x,a)}^{-1}(\alpha))}{\alpha}$$

*The computational complexity is  $O(n \cdot m)$  where  $n$  is the number of transition states and  $m$  is the number of atoms.*

- Proved for increasing unbounded distributions



# CVaR Value Iteration - Experiments



# Outline

- 1 Introduction
  - Motivation
  - Reinforcement Learning
  - Risk
  - Risk-averse Reinforcement Learning
- 2 CVaR Value Iteration
  - Previous results
  - Linear-time improvement
- 3 CVaR Q-learning
  - Var-based policy improvement
- 4 Deep CVaR Q-learning

# Q-learning

## Problem

- In practice, we often don't have access to the transition probabilities  $p(x'|x, a)$
- We need to learn through direct interaction with the environment.

## Solution

Q-learning: Sampling version of Value Iteration.

# Action Value function

## Definition

Value function  $V(x)$  represents the expected return when starting in state  $x$  and following the optimal policy  $\pi^*$  thereafter.

## Definition

Action-Value function  $Q(x, a)$  represents the expected return when starting in state  $x$ , performing action  $a$  and following the optimal policy  $\pi^*$  thereafter.

# Q-learning

## Value Iteration

$$Q(x, a) \leftarrow \mathcal{T}Q(x, a)$$

$$Q(x, a) \leftarrow R(x, a) + \gamma \sum_{x'} p(x'|x, a) \max_{a'} Q(x', a')$$

## Q-learning

$$Q(x, a) \leftarrow (1 - \beta)Q(x, a) + \beta \mathcal{T}Q(x, a)$$

$$Q(x, a) \leftarrow (1 - \beta)Q(x, a) + \beta \left[ R(x, a) + \gamma \max_{a'} Q(x', a') \right]$$

# CVaR estimation

Requirements:

- Store a single value
- Expected value of updates is CVaR

## Recursive CVaR Estimation

$$V_{t+1} = V_t + \beta_t \left[ 1 - \frac{1}{\alpha} \mathbb{1}_{(V_t \geq r)} \right]$$
$$C_{t+1} = (1 - \beta_t) C_t + \beta_t \left[ V_t + \frac{1}{\alpha} (r - V_t)^- \right]$$

# CVaR Q-learning

## Pseudocode

- 1 Sample a transition  $x, a, x', r$
- 2 Create a target distribution  $\mathbf{d}$
- 3 Sample from the target distribution
- 4 Update current estimates of VaR and CVaR towards the sample

## Improvement

- ~~Sample from target distribution~~
- Update proportionally to the target distribution

# CVaR Q-learning

## CVaR Q-learning: General case

```

1: input:  $x, a, x', r$ 
2: for each  $i$  do
3:    $C(x', y_i) = \max_{a'} C(x', a', y_i)$ 
4: end for
5:  $\mathbf{d} = \text{extractDistribution}(C(x', \cdot), \mathbf{y})$ 
6: for each  $i$  do
7:    $V(x, a, y_i) = V(x, a, y_i) + \beta \mathbb{E}_j \left[ 1 - \frac{1}{y_i} \mathbb{1}_{(V(x, a, y_i) \geq r + \gamma d_j)} \right]$ 
8:    $C(x, a, y_i) =$ 
        $(1 - \beta)C(x, a, y_i) + \beta \mathbb{E}_j \left[ V(x, a, y_i) + \frac{1}{y_i} (r + \gamma d_j - V(x, a, y_i))^- \right]$ 
9: end for

```



# Optimal policy extraction

## Standard RL: Optimal policy

$$\pi^*(x) = \arg \max_a Q(x, a)$$

## CVaR VI: Optimal policy

$$\pi^*(x_0) = \arg \max_a C(x, a, \alpha)$$

$$\pi^*(x_1) = \arg \max_a C(x_1, a, \alpha \xi^*(x_0))$$

$$\vdots$$

$$\pi^*(x_t) = \arg \max_a C(x_t, a, y_{t-1} \xi^*(x_{t-1}))$$

# VaR-based Policy Improvement

## Theorem

Let  $\pi$  be a fixed policy,  $\alpha \in (0, 1]$ . By following policy  $\pi'$  from the following algorithm, we will improve  $CVaR_\alpha(Z)$  in expectation:

$$CVaR_\alpha(Z^\pi) \leq CVaR_\alpha(Z^{\pi'})$$

## VaR-based Policy Improvement

$a = \arg \max_a CVaR_\alpha(Z(x_0, a))$

$s = VaR_\alpha(Z(x_0, a))$

Take action  $a$ , observe  $x, r$

**while**  $x$  is not terminal **do**

$$s = \frac{s - r}{\gamma}$$

$a = \arg \max_a \mathbb{E}[(Z(x, a) - s)^-]$

Take action  $a$ , observe  $x, r$

**end while**

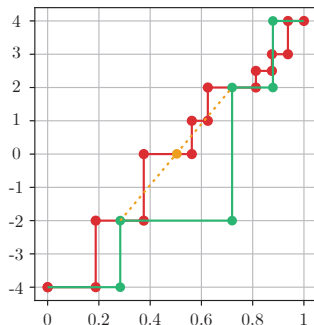
# VaR-based heuristic

## Problem

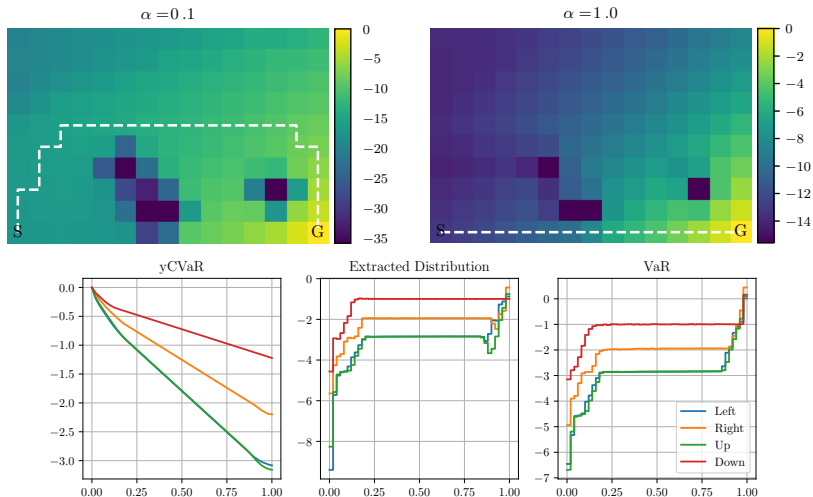
When using quantile discretization, we don't have access to the exact  $VaR$ .

## Solution

Use linear interpolation as a heuristic.



# CVaR Q-learning - Experiments



# Outline

- 1 Introduction
  - Motivation
  - Reinforcement Learning
  - Risk
  - Risk-averse Reinforcement Learning
- 2 CVaR Value Iteration
  - Previous results
  - Linear-time improvement
- 3 CVaR Q-learning
  - Var-based policy improvement
- 4 Deep CVaR Q-learning

# Approximate Q-learning

## Problem

Q-learning is intractable for large state spaces.

## Solution

Use approximate Q-learning.

- Formulate CVaR Q-learning update as a minimizing argument
- Use methods of convex optimization to find the optimal point

# TD update $\rightarrow$ loss function

## Standard RL

$$Q(x, a) = (1 - \beta)Q(x, a) + \beta \mathcal{T}Q(x, a)$$

$$\downarrow$$

$$\min_{\theta} \mathbb{E} \left[ (Q_{\theta}(x, a) - \mathcal{T}Q(x, a))^2 \right]$$

## CVaR RL

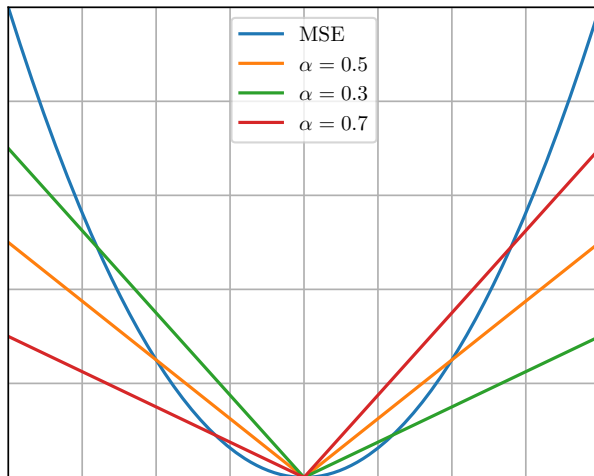
$$V(x, a, y_i) = V(x, a, y_i) + \beta \mathbb{E}_j \left[ 1 - \frac{1}{y_i} \mathbb{1}_{(V(x, a, y_i) \geq \mathcal{T}d_j)} \right]$$

$$C(x, a, y_i) = (1 - \beta)C(x, a, y_i) + \beta \mathbb{E}_j \left[ V(x, a, y_i) + \frac{1}{y_i} (r + \gamma d_j - V(x, a, y_i))^- \right]$$

$$\downarrow$$

$$???$$

# Quantile loss





# TD update $\rightarrow$ loss function

## VaR loss

$$\mathcal{L}_{\text{VaR}} = \sum_{i=1}^N \mathbb{E}_j \left[ (r + \gamma d_j - V_i(x, a))(y_j - \mathbb{1}_{(V_i(x, a) \geq r + \gamma d_j)}) \right]$$

## CVaR loss

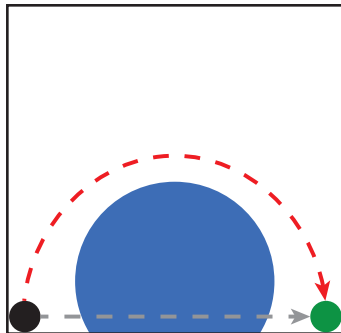
$$\mathcal{L}_{\text{CVaR}} = \sum_{i=1}^N \mathbb{E}_j \left[ \left( V_i(x, a) + \frac{1}{y_i} (r + \gamma d_j - V_i(x, a))^- - C_i(x, a) \right)^2 \right]$$

$$\mathcal{L} = \mathbb{E} [\mathcal{L}_{\text{VaR}} + \mathcal{L}_{\text{CVaR}}]$$

# Deep CVaR Q-learning

- Model: Convolutional Neural Network
  - 1 Input:  $84 \times 84 \times 4$
  - 2 Convolution:  $8 \times 8 \times 32$  (stride 4)
  - 3 Convolution:  $4 \times 4 \times 64$  (stride 2)
  - 4 Convolution:  $3 \times 3 \times 32$  (stride 1)
  - 5 Fully connected: 256 hidden units
  - 6 Output:  $|\mathcal{A}| \times 100$
- Replay Memory
- Target network  $C'$
- Optimizer: Adam (Stochastic Gradient Descent)

# Deep CVaR Q-learning - Experiments



- ① Video:  $\alpha = 1$
- ② Video:  $\alpha = 0.3$

# Summary

## ① Faster CVaR Value Iteration

- Polynomial  $\rightarrow$  linear time.
- Formally proved for increasing, unbounded distributions.
- Experimentally verified for general distributions.

## ② CVaR Q-learning

- Sampling version of CVaR Value Iteration.
- Based on the distributional approach.
- Experimentally verified.

## ③ Distributional Policy improvement

- Proved monotonic improvement for distributional RL.
- Used as a heuristic for extracting  $\pi^*$  from CVaR Q-learning.

## ④ Deep CVaR Q-learning

- TD update  $\rightarrow$  loss function.
- Experimentally verified in a deep learning context.