Safety-Aware Algorithms for Adversarial Contextual Bandits

Wen Sun

Carnegie Mellon University, USA

Debadeepta Dey, and Ashish Kapoor

Microsoft Research, Redmond, USA

The Robotics Institute Carnegie Mellon University



[Georgia Tech Autorally Project, autorally.github.io]

[Silver, et.al, 2016, Nature]



[Georgia Tech Autorally Project, autorally.github.io]

[Silver, et.al, 2016, Nature]





[Georgia Tech Autorally Project, autorally.github.io]

[Silver, et.al, 2016, Nature] Lee Secol vs. Alpha Go ne 122 (W d12): Blact to pla MNOPQRS ABCDEFG пп















[Coates et.al,08,ICML]





[Coates et.al,08,ICML]





[Georgia Tech Autorally Project, autorally.github.io]



[Coates et.al,08,ICML]











Classic Exploration & Exploitation (e.g., epsilon greedy, upper confidence bound) **Not enough**!



Classic Exploration & Exploitation (e.g., epsilon greedy, upper confidence bound) **Not enough**!





Classic Exploration & Exploitation (e.g., epsilon greedy, upper confidence bound) **Not enough**!



Explore more carefully...

Risk of failures Energy exhaustion Side effect of a treatment in clinical trial

This Work:

We attempt to model this problem in the contextual bandit setting

We introduce extra risk associated with each action

The goal is to maintain small regret for reward while ensuring the cumulative risk is small

Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$





Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$





Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$ and *Risk vector*: $r_t \in [0, 1]^K$



Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$ and *Risk vector*: $r_t \in [0, 1]^K$ *Pre-Defined Risk Threshold*: $\beta \in [0, 1]$

Environment



Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$ and *Risk vector*: $r_t \in [0, 1]^K$ *Pre-Defined Risk Threshold*: $\beta \in [0, 1]$



Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$ and *Risk vector*: $r_t \in [0, 1]^K$ *Pre-Defined Risk Threshold*: $\beta \in [0, 1]$



Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$ and *Risk vector*: $r_t \in [0, 1]^K$ *Pre-Defined Risk Threshold*: $\beta \in [0, 1]$



Context (i.e., features): $s_t \in S$ Actions (finitely many): $a_t \in [K]$ Cost vector: $c_t \in [0, 1]^K$ and *Risk vector*: $r_t \in [0, 1]^K$ *Pre-Defined Risk Threshold*: $\beta \in [0, 1]$



No statistical assumptions on the generation of context, cost or risk....

Goal:

Ideally:

$$\min_{i_1,\ldots,i_T} \sum_{t=1}^T c_t[i_t] - \min_{\pi^*} \sum_{t=1}^T c_t[\pi^*(s_t)], \quad s.t., r_t[i_t] \le \beta, \forall t$$

Goal:

Ideally:

$$\min_{i_1,\ldots,i_T} \sum_{t=1}^T c_t[i_t] - \min_{\pi^*} \sum_{t=1}^T c_t[\pi^*(s_t)], \quad s.t., r_t[i_t] \le \beta, \forall t$$

More realistically:

$$\min_{i_1,\dots,i_T} \sum_{t=1}^T c_t[i_t] - \min_{\pi^*} \sum_{t=1}^T c_t[\pi^*(s_t)], \quad s.t., \sum_{t=1}^T (r_t[i_t] - \beta) \le o(T)$$

Goal:

Ideally:

$$\min_{i_1,\ldots,i_T} \sum_{t=1}^T c_t[i_t] - \min_{\pi^*} \sum_{t=1}^T c_t[\pi^*(s_t)], \quad s.t., r_t[i_t] \le \beta, \forall t$$

More realistically:

$$\min_{i_1,\dots,i_T} \sum_{t=1}^T c_t[i_t] - \min_{\pi^*} \sum_{t=1}^T c_t[\pi^*(s_t)], \quad s.t., \sum_{t=1}^T (r_t[i_t] - \beta) \le o(T)$$

No constraint violation in a long-term perspective!

Decision: $x \in \mathcal{X}$ Convex Loss: $\ell_t(x)$ Convex Constraint: $f_t(x) \leq 0$

Environment

Learner



Convex Loss: $\ell_t(x)$ Convex Constraint: $f_t(x) \leq 0$



Decision: $x \in \mathcal{X}$ Convex Loss: $\ell_t(x)$ Convex Constraint: $f_t(x) \leq 0$



Decision: $x \in \mathcal{X}$ Convex Loss: $\ell_t(x)$ Convex Constraint: $f_t(x) \leq 0$



Decision: $x \in \mathcal{X}$ Convex Loss: $\ell_t(x)$ Convex Constraint: $f_t(x) \leq 0$



No constraint violation in a long-term perspective!

Competing Against...

$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^*} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

Competing Against...


$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^*} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^T f_t(x) \leq 0\}$$

$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^*} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^T f_t(x) \leq 0\}$$
 Decisions satisfy average constraint

$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^*} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^{T} f_t(x) \leq 0\}$$
 Decisions satisfy average constraint
Option 2: $\mathcal{O} - \{x \in \mathcal{X} : f_t(x) \leq 0 \ \forall t\}$

Option 2: $\mathcal{O} = \{x \in \mathcal{X} : f_t(x) \le 0, \forall t\}$

$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^*} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^T f_t(x) \leq 0\}$$
 Decisions satisfy average constraint

Option 2: $\mathcal{O} = \{x \in \mathcal{X} : f_t(x) \leq 0, \forall t\}$ Decisions satisfy every constraint

$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^*} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^{T} f_t(x) \leq 0\}$$
 Decisions satisfy average constraint
Option 2: $\mathcal{O} = \{x \in \mathcal{X} : f_t(x) \leq 0, \forall t\}$ Decisions satisfy every constraint

Option 1 seems quite natural.....

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^T f_t(x) \le 0\}$$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^T f_t(x) \leq 0\}$$

Claim: there exist a sequence of loss and constraints such that for any sequence of decisions that satisfies the average constraint:

$$\limsup_{t \to \infty} \sum_{i=1}^{t} f_i(x_i)/t \leq 0,$$

then, the regret grows linearly when competing against \mathcal{O}' :
$$\limsup_{t \to \infty} (\sum_{i=1}^{t} \ell_i(x_i) - \min_{x^* \in \mathcal{O}'} \sum_{i=1}^{t} \ell_i(x^*)) = \Omega(t)$$

Option 1:
$$\mathcal{O}' = \{x \in \mathcal{X} : \sum_{t=1}^T f_t(x) \leq 0\}$$

Claim: there exist a sequence of loss and constraints such that for any sequence of decisions that satisfies the average constraint:

$$\limsup_{t \to \infty} \sum_{i=1}^{t} f_i(x_i)/t \leq 0,$$

then, the regret grows linearly when competing against \mathcal{O}'
$$\limsup_{t \to \infty} (\sum_{i=1}^{t} \ell_i(x_i) - \min_{x^* \in \mathcal{O}'} \sum_{i=1}^{t} \ell_i(x^*)) = \Omega(t)$$

(Construction adapts a discrete two-player game in [Mannor, et.al, 09])

$$\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^* \in \mathcal{O}} \sum_{t=1}^{T} \ell_t(x^*) \le o(T)$$

s.t., $\sum_{t=1}^{T} f_t(x_t) \le o(T)$

where $\mathcal{O} = \{x \in \mathcal{X} : f_t(x) \le 0, \forall t\}$

$$\begin{split} \sum_{t=1}^{T} \ell_t(x_t) &- \min_{x^* \in \mathcal{O}} \sum_{t=1}^{T} \ell_t(x^*) \le o(T) \\ s.t., \sum_{t=1}^{T} f_t(x_t) \le o(T) \\ \end{split}$$
where $\mathcal{O} = \{x \in \mathcal{X} : f_t(x) \le 0, \forall t\}$

Convex-Concave formulation [Mahdavi et al.,2012]

$$\mathcal{L}_t(x,\lambda) = \ell_t(x) + \lambda f_t(x) - \frac{\delta\mu}{2}\lambda^2, \ \delta \in \mathbb{R}^+$$

$$\begin{split} \sum_{t=1}^{T} \ell_t(x_t) &- \min_{x^* \in \mathcal{O}} \sum_{t=1}^{T} \ell_t(x^*) \le o(T) \\ s.t., \sum_{t=1}^{T} f_t(x_t) \le o(T) \\ \end{split}$$
where $\mathcal{O} = \{x \in \mathcal{X} : f_t(x) \le 0, \forall t\}$

Convex-Concave formulation [Mahdavi et al.,2012]

$$\mathcal{L}_{t}(x,\lambda) = \ell_{t}(x) + \lambda f_{t}(x) - \frac{\delta\mu}{2}\lambda^{2}, \quad \delta \in \mathbb{R}^{+}$$
$$\lambda \in \mathbb{R}^{+} \text{ dual variable}$$
$$\lambda \in [0,\infty)$$



Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent

At iteration t, x_t, λ_t

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent

At iteration t, x_t, λ_t $\mathcal{L}_t(x, \lambda) = \ell_t(x) + \lambda f_t(x) - \frac{\delta \mu}{2} \lambda^2, \quad \delta \in \mathbb{R}^+$

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent

At iteration t, x_t, λ_t $\mathcal{L}_t(x, \lambda) = \ell_t(x) + \lambda f_t(x) - \frac{\delta \mu}{2} \lambda^2, \quad \delta \in \mathbb{R}^+$

$$\mathcal{L}_t(x,\lambda_t)$$

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent

At iteration t, x_t, λ_t $\mathcal{L}_t(x, \lambda) = \ell_t(x) + \lambda f_t(x) - \frac{\delta \mu}{2} \lambda^2, \quad \delta \in \mathbb{R}^+$ $\mathcal{L}_t(x, \lambda_t)$ OMD

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent

At iteration t, x_t, λ_t $\mathcal{L}_t(x, \lambda) = \ell_t(x) + \lambda f_t(x) - \frac{\delta \mu}{2} \lambda^2, \quad \delta \in \mathbb{R}^+$ $\mathcal{L}_t(x, \lambda_t) \qquad \qquad \mathcal{L}_t(x_t, \lambda)$ OMD

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent



Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent



 $x_{t+1} = \operatorname{arg\,min}_{x \in \mathcal{X}} D_R(x, \tilde{x}_{t+1})$

Strongly convex regularizer R(x)

Reduction to Online Mirror Descent + Online Gradient Ascent



 $x_{t+1} = \operatorname{arg\,min}_{x \in \mathcal{X}} D_R(x, \tilde{x}_{t+1})$





$$\sum_{t=1}^{\infty} \mathcal{L}_t(x_t, \lambda) - \sum_{t=1}^{\infty} \mathcal{L}_t(x_t, \lambda_t) \le \frac{1}{\mu} + \frac{1}{2} \sum_{t=1}^{\infty} \left(\frac{-\iota(\tau) - \iota(\tau)}{\partial \lambda_t}\right)$$



Combine them:



$$\sum_{t=1} \mathcal{L}_t(x_t, \lambda) - \sum_{t=1} \mathcal{L}_t(x_t, \lambda_t) \le \frac{1}{\mu} + \frac{1}{2} \sum_{t=1}^{\ell} \left(\frac{1}{2} \frac{1}{2} \frac{1}{2}$$

Combine them:

$$\frac{1}{T} \Big[\sum_{t=1}^{T} \ell_t(x_t) - \min_{x^* \in \mathcal{O}} \sum_{t=1}^{T} \ell_t(x^*) \Big] \le O(1/\sqrt{T})$$
$$\frac{1}{T} \Big[\sum_{t=1}^{T} f_t(x_t) \Big] \le O(T^{-1/4})$$



Special Case

Set regularizer:
$$R(x) = \sum_{i} x[i] \ln x[i]$$

Special Case

Set regularizer:
$$R(x) = \sum_{i} x[i] \ln x[i]$$

$$x_{t+1}[i] = \frac{x_t[i] \exp(-\mu \nabla_x \mathcal{L}_t(x_t, \lambda_t)[i])}{\sum_{j=1}^d x_t[j] \exp(-\mu \nabla_x \mathcal{L}_t(x_t, \lambda_t)[j])}$$

$$\lambda_{t+1} = \max\{0, \lambda_t + \mu \nabla \mathcal{L}_t(x_t, \lambda) |_{\lambda = \lambda_t}\}$$





Expert Setting:



Expert Setting:

Finite Expert Set:
$$\Pi : \pi(s) : \mathcal{S} \to \Delta(K)$$



Expert Setting:

Finite Expert Set: $\Pi : \pi(s) : S \to \Delta(K)$ Decision Set: $\Delta(\Pi)$ (all distributions over expert set)



Expert Setting:

Finite Expert Set: $\Pi : \pi(s) : S \to \Delta(K)$ Decision Set: $\Delta(\Pi)$ (all distributions over expert set) Competing Against: $P = \{w \in \Delta(\Pi) : \mathbb{E}_{i \sim w, j \sim \pi_i(s_t)} r_t[j] \le \beta, \forall t\}$



Expert Setting:

Finite Expert Set: $\Pi : \pi(s) : S \to \Delta(K)$ Decision Set: $\Delta(\Pi)$ (all distributions over expert set) Competing Against: $P = \{w \in \Delta(\Pi) : \mathbb{E}_{i \sim w, j \sim \pi_i(s_t)} r_t[j] \leq \beta, \forall t\}$



Reduction to Full Information Setting (i.e., EXP4 [Auer et al., 02,], EXP4.P [Beygelzimer et al., 11])



Environment


















Expert N
$$p_N = \pi_N(s_t)$$









Player

$$p_t = \sum_{i=1}^{N} w_t[i] p_i \rightarrow a_t \in [K]$$

$$c_t[a_t], r_t[a_t]$$

Player

$$p_t = \sum_{i=1}^{N} w_t[i] p_i \rightarrow a_t \in [K]$$

 $c_t[a_t], r_t[a_t]$

 $\hat{c}_t = [0, 0, ..., c_t[a_t]/p_t[a_t], 0...], \quad \hat{r}_t = [0, 0, ..., r_t[a_t]/p_t[a_t], 0, ...]$

Player

$$p_t = \sum_{i=1}^{N} w_t[i]p_i \rightarrow a_t \in [K]$$

 $c_t[a_t], r_t[a_t]$

 $\hat{c}_t = [0, 0, ..., c_t[a_t]/p_t[a_t], 0...], \quad \hat{r}_t = [0, 0, ..., r_t[a_t]/p_t[a_t], 0, ...]$

$$\hat{y}_t[j] = p_j \cdot \hat{c}_t, \hat{z}_t[j] = p_j \cdot \hat{r}_t$$
 for expert j

Player

$$p_t = \sum_{i=1}^{N} w_t[i]p_i \rightarrow a_t \in [K]$$

 $c_t[a_t], r_t[a_t]$

 $\hat{c}_t = [0, 0, ..., c_t[a_t]/p_t[a_t], 0...], \quad \hat{r}_t = [0, 0, ..., r_t[a_t]/p_t[a_t], 0, ...]$

$$\hat{y}_t[j] = p_j \cdot \hat{c}_t, \hat{z}_t[j] = p_j \cdot \hat{r}_t$$
 for expert j

$$\hat{\ell}_t(w) = w \cdot \hat{y}_t, \hat{f}_t(w) = w \cdot \hat{z}_t$$

Player

$$p_t = \sum_{i=1}^{N} w_t[i] p_i \rightarrow a_t \in [K]$$

 $c_t[a_t], r_t[a_t]$

 $\hat{c}_t = [0, 0, ..., c_t[a_t]/p_t[a_t], 0...], \quad \hat{r}_t = [0, 0, ..., r_t[a_t]/p_t[a_t], 0, ...]$

$$\hat{y}_{t}[j] = p_{j} \cdot \hat{c}_{t}, \hat{z}_{t}[j] = p_{j} \cdot \hat{r}_{t} \text{ for expert j}$$

$$\hat{\ell}_{t}(w) = w \cdot \hat{y}_{t}, \hat{f}_{t}(w) = w \cdot \hat{z}_{t}$$

$$\int \beta$$
Black-Box Learner of OCP
with Constraints
(R(w) set to negative entropy)

Player

$$p_t = \sum_{i=1}^{N} w_t[i] p_i \rightarrow a_t \in [K]$$

 $c_t[a_t], r_t[a_t]$

 $\hat{c}_t = [0, 0, \dots, c_t[a_t]/p_t[a_t], 0\dots], \quad \hat{r}_t = [0, 0, \dots, r_t[a_t]/p_t[a_t], 0, \dots]$

$$\hat{y}_{t}[j] = p_{j} \cdot \hat{c}_{t}, \hat{z}_{t}[j] = p_{j} \cdot \hat{r}_{t} \text{ for expert j}$$

$$\hat{\ell}_{t}(w) = w \cdot \hat{y}_{t}, \hat{f}_{t}(w) = w \cdot \hat{z}_{t}$$

$$\downarrow \beta$$
Black-Box Learner of OCP with Constraints
(R(w) set to negative entropy)

Player

$$p_t = \sum_{i=1}^{N} w_t[i] p_i \rightarrow a_t \in [K]$$

 $c_t[a_t], r_t[a_t]$

 $\hat{c}_t = [0, 0, ..., c_t[a_t]/p_t[a_t], 0...], \quad \hat{r}_t = [0, 0, ..., r_t[a_t]/p_t[a_t], 0, ...]$

$$\hat{y}_{t}[j] = p_{j} \cdot \hat{c}_{t}, \hat{z}_{t}[j] = p_{j} \cdot \hat{r}_{t} \text{ for expert j}$$

$$\hat{\ell}_{t}(w) = w \cdot \hat{y}_{t}, \hat{f}_{t}(w) = w \cdot \hat{z}_{t}$$

$$\downarrow \beta$$
Black-Box Learner of OCP
with Constraints
(R(w) set to negative entropy)
$$\downarrow$$
Update to w_{t+1}

$$EXP4.R (EXP4$$
with Risk
constraints)

Under the assumption that $P \neq \emptyset$, for any sequence of cost and risk vectors, EXP4.R has the following guarantees:

$$\mathbb{E}\left[\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j]\right] \leq O(\sqrt{TK \ln(|\Pi|)})$$
$$\mathbb{E}\left[\sum_{t=1}^{T} r_t[a_t] - \beta\right] \leq O(T^{3/4}(K \ln(|\Pi|))^{1/4})$$

Where $w^* \in \{w \in \Delta(\Pi) : \mathbb{E}_{i \sim w, j \sim \pi_i(s_t)} r_t[j] \le \beta, \forall t\}$

Algorithm

$$\hat{\mathcal{L}}_t(w,\lambda) = w \cdot (\hat{y}_t + \lambda \hat{z}_t) - \delta \mu \lambda^2 / 2$$

$$\hat{\mathcal{L}}_t(w,\lambda) = w \cdot (\hat{y}_t + \lambda \hat{z}_t) - \delta \mu \lambda^2 / 2$$

$$\mathbb{E}_t[\hat{y}_t + \lambda \hat{z}_t] = y_t + \lambda z_t$$

$$\hat{\mathcal{L}}_t(w,\lambda) = w \cdot (\hat{y}_t + \lambda \hat{z}_t) - \delta \mu \lambda^2 / 2$$

$$\mathbb{E}_t[\hat{y}_t + \lambda \hat{z}_t] = y_t + \lambda z_t$$

$$\hat{y}_t + \lambda \hat{z}_t - \kappa \sum_{k=1}^K \frac{\pi_j(s_t)[k]}{p_t[k]}$$

$$\hat{\mathcal{L}}_t(w,\lambda) = w \cdot (\hat{y}_t + \lambda \hat{z}_t) - \delta \mu \lambda^2 / 2$$

$$\mathbb{E}_t[\hat{y}_t + \lambda \hat{z}_t] = y_t + \lambda z_t$$

$$\tilde{\mathcal{L}}_t(w,\lambda) = w \cdot \left(\hat{y}_t + \lambda \hat{z}_t - \kappa \sum_{k=1}^K \frac{\pi_j(s_t)[k]}{p_t[k]} \right) - \delta \mu \lambda^2 / 2$$

We also want high probability statement. Use the trick in EXP3.P (and EXP4.P)

$$\hat{\mathcal{L}}_t(w,\lambda) = w \cdot (\hat{y}_t + \lambda \hat{z}_t) - \delta \mu \lambda^2 / 2$$

$$\mathbb{E}_t[\hat{y}_t + \lambda \hat{z}_t] = y_t + \lambda z_t$$

$$\tilde{\mathcal{L}}_t(w,\lambda) = w \cdot \left(\hat{y}_t + \lambda \hat{z}_t - \kappa \sum_{k=1}^K \frac{\pi_j(s_t)[k]}{p_t[k]} \right) - \delta \mu \lambda^2 / 2$$

EXP4.P.R (EXP4.P with Risk Constraints)

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

When $\epsilon \to 0$: Average Regret-> $O(1/\sqrt{T})$ Avg constraint violation-> O(1)

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

When $\epsilon \to 0$: Average Regret-> $O(1/\sqrt{T})$ Avg constraint violation-> O(1)When $\epsilon \to 0.5$: Average Regret-> O(1) Avg constrain violation-> $O(T^{-1/4})$

Challenge: $\mathbb{E}_t[\hat{y}_t + \lambda_t \hat{z}_t] = y_t + \lambda_t z_t$

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

Challenge:
$$\mathbb{E}_t[\hat{y}_t + \lambda_t \hat{z}_t] = y_t + \lambda_t z_t$$

$$Var(\lambda_t \hat{z}_t) = \frac{\lambda_t^2}{p}$$

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

Challenge:
$$\mathbb{E}_t[\hat{y}_t + \lambda_t \hat{z}_t] = y_t + \lambda_t z_t$$

$$Var(\lambda_t \hat{z}_t) = \frac{\lambda_t^2}{p} \quad \lambda_t = O(\frac{\beta}{\delta\mu})$$

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

Challenge:
$$\mathbb{E}_t[\hat{y}_t + \lambda_t \hat{z}_t] = y_t + \lambda_t z_t$$

$$Var(\lambda_t \hat{z}_t) = \frac{\lambda_t^2}{p} \quad \lambda_t = O(\frac{\beta}{\delta\mu}) \sqrt{T}$$

Theorem: Under the assumption that $P \neq \emptyset$, for any sequence of cost, risk vectors and any $\epsilon \in (0, 0.5)$, we have with high probability 1- v:

$$\sum_{t=1}^{T} c_t[a_t] - \sum_{t=1}^{T} \mathbb{E}_{i \sim w^*, j \sim \pi_i(s_t)} c_t[j] \le O(\sqrt{T^{\epsilon+1/2} K \ln(\Pi/v)})$$
$$\sum_{t=1}^{T} (r_t[a_t] - \beta) \le O(T^{1-\epsilon/2} \sqrt{K \ln(\Pi)})$$

Challenge:
$$\mathbb{E}_t[\hat{y}_t + \lambda_t \hat{z}_t] = y_t + \lambda_t z_t$$

$$Var(\lambda_t \hat{z}_t) = \frac{\lambda_t^2}{p} \quad \lambda_t = O(\frac{\beta}{\delta\mu}) \longrightarrow Var(\lambda_t \hat{z}_t) = \frac{T}{p}$$







Context is the RBF feature with respect to the nine way points.

We have 4^9 experts. Namely each expert suggests one action at each waypoint

We ran the EXP4.R with different risk thresholds






Conclusion and Future Work

- 1. We consider sequential decision making problem with additional adversarial constraints.
- 2. In our applications these constraints are used to model safety related issues in decision making process.
- 3. Is there any algorithm that can achieve \sqrt{T} total regret and \sqrt{T} total constraint violation simultaneously?
- 4. Is there better heuristic we can leverage to achieve tighter regret and constrain violation in high probability?

Thanks

Wen Sun [wensun@cs.cmu.edu]

The Robotics Institute Carnegie Mellon University

