# AN EFFICIENT PESSIMISTIC-OPTMISITIC ALGORITHM FOR CONSTRAINED STOCHASTIC LINEAR BANDITS

Xin Liu, Assistant Professor @ ShanghaiTech

Bin Li
Pennsylvania State University

Pengyi Shi
Purdue University

Lei Ying
University of Michigan, Ann Arbor

# CONSTRAINED BANDITS: ONLINE DISPATCHING



Crowdsourcing

❑ Jobs arrive in a dynamic way
❑ Dispatch jobs to servers
❑ Observe reward, cost, and budget

# CONSTRAINED BANDITS: ONLINE DISPATCHING

Crowdsourcing

Healthcare

❑ Jobs arrive in a dynamic way
❑ Dispatch jobs to servers
❑ Observe reward, cost, and budget

# CONSTRAINED BANDITS: ONLINE DISPATCHING

# CONSTRAINED BANDITS: ONLINE DISPATCHING

# CONSTRAINED BANDITS: ONLINE DISPATCHING

# CONSTRAINED BANDITS: ONLINE DISPATCHING

$$\max_{\pi} E\left[\sum_{t=1}^{T} R(t)\right]$$

$$s.t. \quad E\left[\sum_{t=1}^{\tau} W(t)\right] \leq E\left[\sum_{t=1}^{\tau} U(t)\right]$$

# CONSTRAINED BANDITS: ONLINE DISPATCHING

$$\max_{\pi} \mathrm{E}\left[\sum_{t=1}^{T} \mathrm{R}(c_t, \mathrm{A}^{\pi}(t))\right]$$

$$\mathrm{s.t.} \quad \mathrm{E}\left[\sum_{t=1}^{\tau} \mathrm{W}(c_t, \mathrm{A}^{\pi}(t))\right] \leq \mathrm{E}\left[\sum_{t=1}^{\tau} \mathrm{U}(t)\right]$$

$$\mathrm{R}(\, \text{⚙}, \text{👷} \,) = 0.1$$

$$\mathrm{W}(\, \text{⚙}, \text{👷} \,) = 90$$

amazon mechanical turk

# CONSTRAINED BANDITS: ONLINE DISPATCHING

$$\max_{\pi} E\left[\sum_{t=1}^{T} R(c_t, A^{\pi}(t))\right]$$

$$\text{s.t. } E\left[\sum_{t=1}^{\tau} W(c_t, A^{\pi}(t))\right] \leq E\left[\sum_{t=1}^{\tau} U(t)\right]$$

# CONSTRAINED BANDITS: ONLINE DISPATCHING



V.S.

**Goals:**

❑ Achieve optimal performance:

$$\text{Regret} = \text{TotalReward}(\pi^*) - \text{TotalReward}(\pi).$$

❑ Guarantee zero violation:

$$\text{Violation} = \text{TotalCost}(\pi) - \text{TotalBudget}.$$

# CONSTRAINED STOCHASTIC LINEAR BANDITS

## Model:

1. N servers and T time slots

2. A task arrives with feature $c_t$ at time slot $t$.

3. Rewards are unknown:
$$R(c_t, j) = <\textcolor{red}{\theta_*}, \phi(c_t, j)> + \eta_t$$

4. Costs are known:
$$W(c_t, j)$$

5. Stochastic (anytime) accumulative constraints:
$$E[\textstyle\sum_{t=1}^{\tau} W(c_t, A(t))] \leq E[\textstyle\sum_{t=1}^{\tau} U(t)]$$

$c_t$

amazon
mechanical turk

# CONSTRAINED STOCHASTIC LINEAR BANDITS

## Model:

1. N servers and T time slots
2. A task arrives with feature $c_t$ at time slot $t$.
3. Rewards are unknown:

$$R(c_t, j) = <\textcolor{red}{\theta_*}, \phi(c_t, j)> + \eta_t$$

4. Costs are known:

$$W(c_t, j)$$

5. Stochastic (anytime) accumulative constraints:

$$E[\sum_{t=1}^{\tau} W(c_t, A(t))] \leq E[\sum_{t=1}^{\tau} U(t)]$$

$$R(\text{⚙}, \text{👷}) = 0.6$$

amazon
mechanical turk

$c_t$

$$R(\text{⚙}, \text{👷}) = 0.8$$

# CONSTRAINED STOCHASTIC LINEAR BANDITS

## Model:

$W(\text{⚙}, \text{👷}) = 30$

1. N servers and T time slots

2. A task arrives with feature $c_t$ at time slot $t$.

3. Rewards are unknown:
$$R(c_t, j) = <\ \theta_*, \phi(c_t, j)\ > +\ \eta_t$$

4. Costs are known:
$$W(c_t, j)$$

$c_t$

5. Stochastic (anytime) accumulative constraints:
$$E\left[\sum_{t=1}^{\tau} W(c_t, A(t))\right] \leq E\left[\sum_{t=1}^{\tau} U(t)\right]$$

$W(\text{⚙}, \text{👷}) = 10$

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation



$B_t$

$\cdot\, \theta_*$

$\cdot\, \hat{\theta}_t$

$\hat{r}(c_t, j)$

$\widehat{\text{reward}}(\text{⚙}, \text{👷}) = 0.9$

$\text{true reward}(\text{⚙}, \text{👷}) = 0.8$

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation

$B_t$ $\quad \bullet \theta_*$

$\bullet \hat{\theta}_t$

$\hat{r}(c_t, j)$

2. Pessimistic action

$$A(t) = \underset{j}{\mathrm{argmax}}\ \widehat{\mathrm{reward}}(c_t, j) - \mathrm{violation}(c_t, j)$$

$$\hat{r}(\{\!\!\{\!\circ\!\}\!\!\}, \text{👷}) > \hat{r}(\{\!\!\{\!\circ\!\}\!\!\}, \text{👷})$$

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation

$$B_t \qquad \overset{\bullet}{\theta_*}$$

$$\cdot \hat{\theta}_t$$

$$\to \hat{r}(c_t, j)$$

2. Pessimistic action

$$A(t) = \underset{j}{\text{argmax}} \ \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation



$$\hat{r}(c_t, j)$$

2. Pessimistic action

$$A(t) = \underset{j}{\text{argmax}} \; \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

3. Calibration

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation



$\hat{r}(c_t, j)$

2. Pessimistic action

$$A(t) = \underset{j}{\text{argmax}} \ \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

3. Calibration on reward

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation

$B_t$   $\cdot\, \theta_*$   $\cdot\, \hat{\theta}_t$   $\hat{r}(c_t, j)$

2. Pessimistic action

$$A(t) = \underset{j}{\arg\max}\ \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

3. Calibration on reward

reward

(job, action)

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation



2. Pessimistic action

$$A(t) = \underset{j}{\text{argmax}}\ \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

3. Calibration on reward

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation



$\hat{r}(c_t, j)$

2. Pessimistic action

$$A(t) = \underset{j}{\text{argmax}} \ \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

3. Calibration on reward

# BANDIT LEARNING-BASED ONLINE ALGORITHM

1. Optimistic reward estimation



2. Pessimistic action

$$A(t) = \underset{j}{\text{argmax}} \; \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

3. Calibration on violation

$$\text{violation}(t + 1) = \underbrace{\text{violation}(t)}_{5} + \underbrace{\text{cost}(t) - \text{budget}(t)}_{3}$$

# PESSIMISTIC-OPTIMISTIC ONLINE ALGORITHM

1. Optimistic reward estimation

$\hat{r}(c_t, j)$

$B_t$ $\cdot\theta_*$

$\cdot\hat{\theta}_t$

2. Pessimistic action

$$A(t) = \underset{j}{\arg\max}\ \widehat{reward}(c_t, j) - violation(c_t, j)$$

3. Calibration on violation

$$violation(t+1) = \underbrace{violation(t)}_{5} + \underbrace{cost(t) - budget(t)}_{3}$$

8

5

# PESSIMISTIC-OPTIMISTIC ONLINE ALGORITHM

> **Theorem (Informal):**
>
> Pessimistic-optimistic algorithm achieves Regret$(\tau) = O(\sqrt{\tau})$ and Violation$(\tau) = 0$ after some constant rounds.

[LiuLiShiYing21] First efficient online algorithm to achieve optimal regret & violation (anytime).

# PESSIMISTIC-OPTIMISTIC ONLINE ALGORITHM

> **Theorem (Informal):**
>
> Pessimistic-optimistic algorithm achieves Regret$(\tau)$ = O$(\sqrt{\tau})$ and Violation$(\tau)$ = 0 after some constant rounds.

[LiuLiShiYing21] First efficient online algorithm to achieve optimal regret & violation (anytime).

| Related Work | Constriant Type |
|---|---|
| AD14, AD16, BKS18, CER20 | Constraints imposed at the end of time horizon |
| AAT19 | Anytime action constraints |
| PGBJ20 | Anytime policy constraints |



cumulative constraints

anytime policy constraints

anytime action constraints

anytime cumulative constraints

PGBJ20    AAT19    Ours

# PESSIMISTIC-OPTIMISTIC ONLINE ALGORITHM

> **Theorem (Informal):**
>
> Pessimistic-optimistic algorithm achieves Regret$(\tau) = O(\sqrt{\tau})$ and Violation$(\tau) = 0$ after some constant rounds.

[LiuLiShiYing21] First efficient online algorithm to achieve
optimal regret & violation (anytime).

| Related Work | Constriant Type |
|:---:|:---:|
| AD14, AD16, BKS18, CER20 | Constraints imposed at the end of time horizon |
| AAT19 | Anytime action constraints |
| PGBJ20 | Anytime policy constraints |



Primal-dual approach with adaptive optimism in primal and pessimism in dual.

27

# ADAPTIVE OPTIMISM-PESSIMISM IN PRIMAL-DUAL

**Primal (action):**

$$A(t) = \underset{j}{\text{argmax}} \quad \widehat{\text{reward}}(c_t, j) - \text{violation}(c_t, j)$$

$$= \underset{j}{\text{argmax}} \quad V_t \hat{r}(c_t, j) - \text{violation}(c_t, j)$$

**Dual (calibration):**

$$\text{violation}(t + 1) = \text{violation}(t) + \text{cost}(t) - \text{budget}(t) + \epsilon_t$$

time $t$

pessimism

$\theta_*$

$V_t = \sqrt{t}$
increases

$\epsilon_t = 1/\sqrt{t}$
decreases

$\theta_*$

$\vdots$

$\theta_*$

optimism

$\vdots$

# CONCLUSION

Pessimistic-optimistic online algorithm:
- achieve optimal regret & violation (anytime).
- a novel drift analysis framework to bridge regret and violation.

# THANK YOU!