

## Contextual Bandits

Input actions/arms  $\mathcal{X}$

for  $t=1, 2, \dots$

Nature reveals context  $c_t \sim \mathcal{V}$

Player chooses action  $x_t \in \mathcal{X}$

Nature reveals reward  $V_t(c_t, x_t) \in [0, 1]$

where  $\mathbb{E}[V_t(c_t, x_t) | c_t, x_t] = V(c_t, x_t)$

for some unknown  $V: \mathcal{C} \times \mathcal{X} \rightarrow [0, 1]$ .

Suppose context space was finite  $|\mathcal{C}| < \infty$ .

Suppose we run  $|\mathcal{C}|$  separate bandit algorithms in parallel,  
one for each  $c \in \mathcal{C}$ .

$$\max_x \sum_{t=1}^T \mathbb{1}\{c_t = c\} (V(c_t, x) - V(c_t, x_t))$$

$$\leq \sqrt{|\mathcal{X}| T_c \log(T_c / \delta)}$$

$$\leq \sqrt{|\mathcal{X}| T_c \log(T / \delta)}$$

← w.p.  $\geq 1 - \delta$  using  
elim. or UCB

$$T_c = \sum_{t=1}^T \mathbb{1}\{c_t = c\}$$

$$\sum_{t=1}^T \max_x (V(c_t, x) - V_t(c_t, x_t))$$

$$= \sum_{t=1}^T \sum_{c \in \mathcal{C}} \max_x \mathbb{1}\{c_t = c\} (V(c_t, x) - V(c_t, x_t))$$

$$= \sum_{c \in \mathcal{C}} \max_x \sum_{t=1}^T \mathbb{1}\{c_t = c\} (V(c_t, x) - V(c_t, x_t))$$

$$\leq \sum_{c \in \mathcal{C}} \sqrt{|\mathcal{X}| \cdot T_c \cdot \log(T / \delta)}$$

$$\begin{aligned} \text{Cauchy-Schwarz} &\leq \sqrt{\left(\sum_{c \in \mathcal{C}} T_c\right) \left(\sum_{c \in \mathcal{C}} |x| L_c(T/d)\right)} \\ &= \sqrt{|\mathcal{C}| |x| T \log(T/d)} \end{aligned}$$

What if we ignored context and played just a single band it?

$$\max_x \sum_{t=1}^T V(c_t, x) - V(c_t, x) \leq \sqrt{|x| T \log(T/d)}$$

Total Reward acts like

$$\sum_{t=1}^T V(c_t, x_t) \geq \max_x \sum_{t=1}^T V(c_t, x) - \sqrt{|x| T \log(T/d)}.$$

But compare this to  $|\mathcal{C}|$ -bandit approach that achieves

$$\sum_{t=1}^T V(c_t, x_t) \geq \sum_{t=1}^T \max_x V(c_t, x) - \sqrt{|x| |\mathcal{C}| T \log(T/d)}$$

If  $T \ll |\mathcal{C}|$  then may be far better to play just a single bandit than a bandit-per-context b/c latter will never learn.

Now suppose we have a set of

policies  $\Pi = \{\pi: \mathcal{D} \rightarrow \mathcal{X}\}$ .

At each time  $t$ , choose  $\pi_t \in \Pi$

and play  $a_t = \pi(C_t)$ .

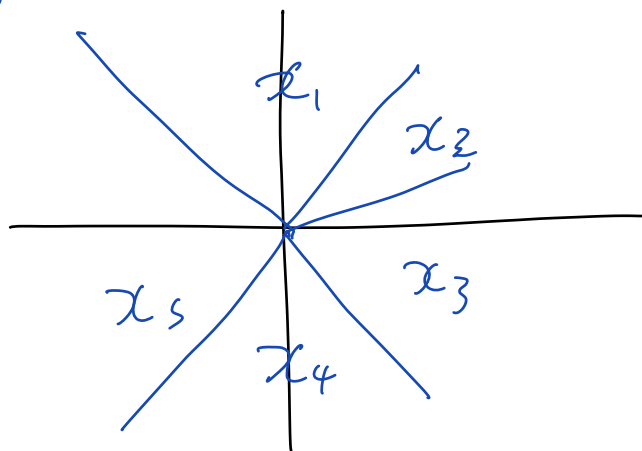
Ex. suppose  $C_t \in \mathbb{R}^d$ , then take

$\Pi$  to be linear multiclass classifiers

$$\pi \in \Pi \Leftrightarrow \pi(C_t) = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \langle w_x, C_t \rangle$$

for some  $\{w_x\}_{x \in \mathcal{X}}$

$$|\mathcal{X}| = 5$$



---

Goal is to minimize policy regret!

$$\max_{\pi \in \Pi} \mathbb{E} \left[ \sum_{t=1}^T V(c_t, \pi(c_t)) - V(c_t, x_t) \right]$$

$\uparrow = \pi_t(c_t)$

Ex. The bandit-per-context strategy is encoded w/  $|\Pi| = |\mathcal{X}|^{|\mathcal{C}|}$

Ex. The single-bandit (ignore context)  $|\Pi| = |\mathcal{X}|$

$$V(\pi) = \mathbb{E}_{c \sim \nu} [V(c, \pi(c))]$$

$$\text{Regret} \quad \max_{\pi} T V(\pi) - \sum_{t=1}^T V(\pi_t)$$

Could think of the policies themselves as "arms"

so that I have  $|\Pi|$  arms, and then play

any MAB alg (UCB) for  $T$  times where

at each time I choose a policy and receive reward.

for  $t=1, 2, \dots$

Player chooses  $\pi_t \in \Pi$

Nature reveals reward  $r_t = V_t(c_t, \pi(c_t))$

$$\text{Regret} \quad \max_{\pi} \sum_{t=1}^T V_t(c_t, \pi(c_t)) - r_t \leq \sqrt{|\Pi| T \log(T)}$$



Regret bound is independent of #actions, #context,

But horrible if  $|T|$  is large.

---

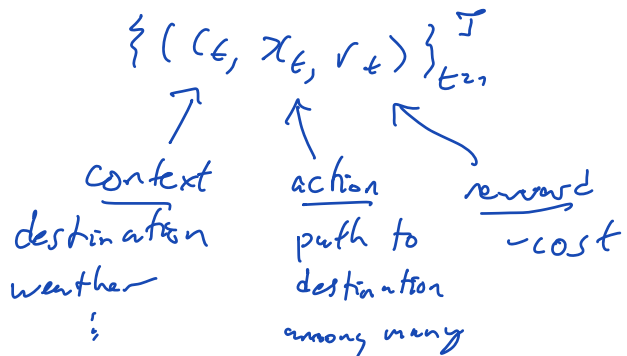
If  $\Pi$  is a class of multi-class classifiers

(e.g. linear, neural nets) then we know  
for classification we can learn a good classifier  
from historical data. What if I have  
historical data for a deployed policy?

Ex. USPS observes <sup>destination</sup> location of a package (context)

and then they choose a route to deliver it  
(action). Reward signal is time or cost.

Mail has been delivered for long time by some  
(implicit) policy  $\hat{\pi}$  resulting in data



Want to learn new policy from old data.

Idea: learn a function  $f: \mathcal{C} \times \mathcal{X} \rightarrow \mathbb{R}$  in  $\mathcal{F}$

$$\text{s.t. } \hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{t=1}^T (f(c_t, a_t) - r_t)^2$$

↑  
neural net  
decision tree  
whatever

To learn a new policy:

$$\pi_{\text{new}}(c) = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \hat{f}(c, x).$$

"Model the world"

What could go wrong?

USPS collected data under policy  $\pi_{\text{old}}$ , meaning  
in response to context  $c_t$  they took action  $a_t = \pi_{\text{old}}(c_t)$ .

By observing  $v_t(c_t, a_t)$  that means we did  
not observe  $v(c_t, a)$  for  $a \neq a_t$ .

$\Rightarrow$  We may not be able to learn  $\underset{a}{\operatorname{argmax}} v(c_t, a)$ !

Solution? Randomization will give you coverage  
of all actions and you can get  
consistent estimation of  $v(c, a) \forall c, a$ .

In general is necessary.

But sometimes you can get lucky:  
 Suppose you defined a feature map

$$\phi: \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d \quad \text{and}$$

$$\mathcal{F} = \{f(c, a) = \langle \phi(c, a), \theta \rangle, \forall \theta \in \mathbb{R}^d\}$$

Assume  $\exists \theta_* \in \mathbb{R}^d : v(c, a) = \langle \phi(c, a), \theta_* \rangle$

Assume  $r_t = V_t(c_t, a_t) = V(c_t, a_t) + \zeta_t$   
 $\zeta_t \sim \mathcal{N}(0, 1)$

Then  $\operatorname{argmin}_{f \in \mathcal{F}} \sum_{t=1}^T (f(c_t, a_t) - r_t)^2$

is equiv. to

$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{t=1}^T (\langle \phi(c_t, a_t), \theta \rangle - r_t)^2$$

$$= \left( \sum_t \phi(c_t, a_t) \phi(c_t, a_t)^T \right)^{\vee} \sum_t \phi(c_t, a_t) r_t$$

$$= \theta_* + \left( \sum_t \phi(c_t, a_t) \phi(c_t, a_t)^T \right)^{\vee} \sum_t \phi(c_t, a_t) \zeta_t$$

$$\|\hat{\theta} - \theta_*\| \left( \sum_t \phi(c_t, a_t) \phi(c_t, a_t)^T \right) \leq \text{const.}$$

vector martingale  
self-normalized bound.

This linear example demonstrates that  
naively fitting a function to data can work.  
But not always.

Suppose the data was collected  
randomly s.t.  $P(a_t = a | c_t) = \mu(a | c_t)$   
for some distribution  $\mu$ .

Moreover the probability  $P_t = \mu(a_t | c_t)$   
was logged as well to give you

$$\left\{ (c_t, a_t, P_t, r_t) \right\}_{t=1}^T$$



Consider the IPS (inverse propensity <sup>score</sup> estimator)

$$\hat{V}_t(c_t, a) = \frac{\mathbb{I}\{a_t = a\} r_t}{P_t}$$

$$\mathbb{E}_{r_t, a_t} \left[ \hat{V}_t(c_t, a) \mid c_t \right]$$

$$= \sum_{a' \in \mathcal{X}} \mu(a' | c_t) \mathbb{E} \left[ \frac{\mathbb{I}\{a_t = a\} r_t}{P_t} \mid c_t, a_t = a' \right]$$

$$= \sum_{a' \in \mathcal{X}} \mu(a' | c_t) \mathbb{E} \left[ \frac{\mathbb{I}\{a_t = a\} r_t}{\mu(a_t | c_t)} \mid c_t, a_t = a' \right]$$

$$= \sum_{a' \in \mathcal{X}} \cancel{\mu(a' | c_t)} \frac{\mathbb{I}\{a' = a\} v(c_t, a')}{\cancel{\mu(a' | c_t)}}$$

$$= \sum_{a' \in \mathcal{X}} \mathbb{I}\{a' = a\} v(c_t, a')$$

$$= v(c_t, a)$$

$$\hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \hat{V}_t(c_t, \pi(c_t))$$

$$\mathbb{E}[\hat{V}(\pi)] = \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\hat{V}_t(c_t, \pi(c_t))]$$

$$= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{c_t \sim \mu} [V(c_t, \pi(c_t))]$$

$$= V(\pi)$$

Great! Output  $\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \hat{V}(\pi)$ .

What is  $V(\pi_*) - V(\hat{\pi})$ ?

$$\mathbb{E}[(\hat{V}(\pi) - V(\pi))^2]$$

$$= \frac{1}{T^2} \sum_{t=1}^T \mathbb{E}_{c_t} \left[ \underbrace{\mathbb{E}[(\hat{V}_t(c_t, \pi(c_t)) - V(c_t, \pi(c_t)))^2 | c_t]}_{\leq \frac{1}{\mu(\pi(c_t) | c_t)}} \right]$$

Fact! for any R.V.  $X$ ,  $\operatorname{argmin}_a \mathbb{E}[(X-a)^2] = \mathbb{E}[X]$

$$\mathbb{E} \left[ \left( \hat{v}_t(c_t, a) - v(c_t, a) \right)^2 \mid c_t \right]$$

$$\leq \mathbb{E} \left[ \hat{v}_t(c_t, a)^2 \mid c_t \right]$$

$$= \mathbb{E} \left[ \frac{\mathbb{1}\{a_t = a\} r_t^2}{p_t^2} \mid c_t \right] \quad \text{Assume } r_t \in [0, 1]$$

$$= \sum_{a' \in \mathcal{X}} \mu(a' \mid c_t) \mathbb{E} \left[ \frac{\mathbb{1}\{a_t = a\} r_t^2}{p_t^2} \mid c_t, a_t = a' \right]$$

$$= \sum_{a' \in \mathcal{X}} \mu(a' \mid c_t) \mathbb{E} \left[ \frac{\mathbb{1}\{a_t = a\} r_t^2}{\mu(a_t \mid c_t)^2} \mid c_t, a_t = a' \right]$$

$$\leq \sum_{a'} \frac{\mathbb{1}\{a' = a\}}{\mu(a' \mid c_t)} \quad (r_t \in [0, 1])$$

$$= \frac{1}{\mu(a \mid c_t)}$$

$$\text{Var} \left[ v_t(c_t, a) \mid c_t \right] \leq \frac{1}{\mu(a \mid c_t)}$$

$$\begin{aligned}\hat{V}(\pi) &= \frac{1}{J} \sum_{t=1}^J \hat{V}_t(c_t, \pi(c_t)) \\ &= \frac{1}{J} \sum_{t=1}^J \frac{\mathbb{1}\{a_t = \pi(c_t)\}}{P_t} r_t\end{aligned}$$

$$\mathbb{E}[\hat{V}(\pi)] = V(\pi)$$

$$\begin{aligned}\mathbb{E}[(\hat{V}(\pi) - V(\pi))^2] &\leq \mathbb{E}\left[\frac{1}{\mu(\pi(c)|c)}\right] \quad \forall \pi \in \Pi \\ &= \sum_c \gamma_c \frac{1}{\mu(\pi(c)|c)}\end{aligned}$$

$\mu(\pi(c)|c)$  can be arbitrarily small, which means

$P_t$  can be arbitrarily small, which means

$\hat{V}(\pi)$  is heavy-tailed (not sub-Gaussian).

Assume  $\mu(a|c) > \gamma \quad \forall c, a$ . (said another way  $\gamma := \min_{c,a} \mu(a|c)$ )

Now  $\underbrace{\frac{\mathbb{1}\{a_t = \pi(c_t)\}}{P_t}}_{= X_t} r_t$  are unbiased, have variance  $\leq \mathbb{E}\left[\frac{1}{\mu(\pi(c)|c)}\right]$  (which could be small)

$r_t \in [0, 1]$  and have support in  $[0, \frac{1}{\gamma}]$ .

Bernstein's Inequality, Let  $X_1, \dots, X_n$  be

independent R.V.s w/  $\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] \leq \sigma^2$

and  $|X_i| \leq b$ . Then

$$\left| \frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right| \leq \sqrt{\frac{2\sigma^2 \log(2/\delta)}{n}} + \frac{2b \log(1/\delta)}{3n}$$

w.p.  $\geq 1 - \delta$ .

Compare to Hoeffding. All we have is that

$|X_i| \leq b$  so Hoeffding says

$$\left| \frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right| \leq b \sqrt{\frac{2 \log(1/\delta)}{n}}$$

w.p.  $\geq 1 - \delta$ .

$$\hat{V}(\pi) = \frac{1}{3} \sum_{t=1}^3 \underbrace{\frac{\mathbb{1}\{a_t = \pi(c_t)\}}{P_t}}_{= X_t} r_t \quad \leftarrow \text{"model the bias"}$$

$$\mathbb{E}[X_t] = \mathbb{E}_{c_t} \left[ \mathbb{E} \left[ \frac{\mathbb{1}\{a_t = \pi(c_t)\}}{P_t} r_t \mid c_t \right] \right]$$

$$= \mathbb{E}_{c_t} [V(c_t, \pi(c_t))] = V(\pi)$$

$$\mathbb{E}[(X_t - \mathbb{E}[X_t])^2]$$

$$= \mathbb{E}_{c_t} \left[ \mathbb{E} \left[ \left( \frac{\mathbb{1}\{a_t = \pi(c_t)\}}{p_t} r_t - V(\pi) \right)^2 \mid c_t \right] \right]$$

$$\leq \mathbb{E} \left[ \frac{1}{\mu(\pi(c) | c)} \right] = \sigma^2$$

$$X_t \in [0, \frac{1}{\delta}) \quad b = \frac{1}{\delta}$$

$\Rightarrow$  Bernstein applied: For any  $\pi \in \Pi$

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\mathbb{E} \left[ \frac{1}{\mu(\pi(c) | c)} \right] \frac{2 \log(2/\delta)}{3}} + \frac{2 \log(2/\delta)}{33\gamma}$$

w.p.  $\geq 1 - \delta$ .

To hold for all  $\pi \in \Pi$  take union bound  $\delta \mapsto \frac{\delta}{|\Pi|}$ .

We said uniform was a good idea

for exploration:  $\mu(a | c) = \frac{1}{|\mathcal{X}|}$  ( $|\mathcal{X}|$  actions)

That implies that for all  $\pi \in \Pi$  simultaneously

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\frac{2|\mathcal{X}| \log(2|\Pi|/\delta)}{3}} + \frac{2|\mathcal{X}| \log(2|\Pi|/\delta)}{33\gamma}$$

w.p.  $\geq 1 - \delta$ .

$$\leq \sqrt{\frac{4|\mathcal{X}| \log(2|\mathcal{T}|/\delta)}{3}}$$

$$\Rightarrow \text{w.p.} \geq 1-\delta$$

$$\hat{\pi} = \arg \max_{\pi \in \mathcal{T}} \hat{V}(\pi)$$

$$\begin{aligned} V(\pi_*) - V(\hat{\pi}) &= V(\pi_*) - \hat{V}(\pi_*) + \underbrace{\hat{V}(\pi_*) - \hat{V}(\hat{\pi})}_{\leq 0} \\ &\quad + \hat{V}(\hat{\pi}) - V(\hat{\pi}) \end{aligned}$$

$$\leq 2 \max_{\pi \in \mathcal{T}} |\hat{V}(\pi) - V(\pi)|$$

$$\leq 2 \sqrt{\frac{4|\mathcal{X}| \log(2|\mathcal{T}|/\delta)}{3}}.$$

Great: you will show<sup>or hw</sup> an explore-then-commit  
 stat that exploits this to achieve  $(T)^{2/3}$  regret.  
 Can we do better? yes.

But first how do you solve

$$\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \hat{V}(\pi)$$

$$= \operatorname{argmax}_{\pi \in \Pi} \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{a_t = \pi(c_t)\} r_t}{P_t}$$

$$= \operatorname{argmax}_{\theta \in \Theta} \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{a_t = \pi_{\theta}(c_t)\} r_t}{P_t}$$

$$= \operatorname{argmin}_{\theta \in \Theta} \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{a_t \neq \pi_{\theta}(c_t)\} r_t}{P_t}$$

$$\hat{f} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{3} \sum_{t=1}^3 \left( \frac{r_t}{P_t} \right) \operatorname{softmax} \left( f_{\theta}(c_t), a_t \right)$$

$$\hat{\pi}(c) = \hat{f}(c)$$

—  
Doubly robust estimator

Given  $\{(c_t, a_t, P_t, r_t)\}_{t=1}^3$  and  $\hat{f}: \mathcal{C} \times \mathcal{X} \rightarrow [0, 1]$

$$\hat{V}_t^{(DR)}(c_t, a) = \hat{f}(c_t, a) + (r_t - \hat{f}(c_t, a)) \frac{\mathbb{1}\{a_t = a\}}{P_t}$$

Unbiased  $\mathbb{E}[\hat{V}_t^{(DR)}(c_t, a) | c_t] = V(c_t, a) \neq \hat{f}$



(IPS estimator takes  $\hat{f}(c, a) = 0$ )