

Contextual Bandits w/ policies Π : $\pi \in \Pi$ is a map $\bar{\pi}: \mathcal{C} \rightarrow \mathcal{X}$

Input Π

for $t=1, 2, \dots, T$

Nature reveals context $c_t \in \mathcal{C}$ is iid from \mathcal{D}

Player choose action $x_t \in \mathcal{X}$ (equiv. chooses $\pi_t \in \Pi$ and plays $\pi_t(c_t)$)

Player receives $r_t \in [0, 1]$ where $\mathbb{E}[r_t | c_t, x_t] = v(c_t, x_t)$

Define $V(\pi) = \mathbb{E}_c [v(c, \pi(c))]$

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

$$\max_{\pi \in \Pi} \sum_{t=1}^T v(c_t, \pi) - v(c_t, \pi_t) \quad (2)$$

$$\max_{\pi} \sum_{t=1}^T r_t[\pi(c_t)] - r_t[\pi_t(c_t)] \quad (3)$$

In our stochastic setting these are all within $O(\sqrt{T})$

Idea: play some logging policy μ : in response

to context c_t , policy plays $x_t \in \mathcal{X}$ w.p. $\mu(x | c_t)$

Collect dataset $\{(c_t, x_t, r_t, p_t)\}_{t=1}^T$ $p_t := \mu(x_t | c_t)$

$$\text{Estimate } \hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{p_t} r_t$$

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

$$\mathbb{E}[(\hat{V}(\pi) - V(\pi))^2] \leq \mathbb{E}_{c \sim \mathcal{D}} \left[\frac{1}{\mu(\pi(c) | c)} \right].$$

By Bernstein. w.p. $\geq 1-\delta$, $\forall \pi \in \Pi$ we have

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\mathbb{E}_{c \sim \mathcal{D}} \left[\frac{1}{\mu(\pi(c)|c)} \right]} \cdot \frac{2 \log(2/\delta)}{3} + \frac{2 V_{\max} \log(2/\delta)}{33}$$

where $V_{\max} = \max_{c, x} \frac{1}{\mu(x, c)}$.

\mathcal{J} -Greedy: Algorithm for "model the bias"/general policies

Input Π

for $t=1, 2, \dots, \mathcal{J}$

Nature reveals c_t

Alg plays x_t uniformly at random

Construct $\hat{V}(\pi)$ as above

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

for $t = \mathcal{J}+1, \dots, T$

Nature reveals c_t

$$x_t = \hat{\pi}(c_t)$$

$$\mu(x|c_t) = \frac{1}{|\mathcal{X}|} \text{ for all } x \in \mathcal{X}.$$

$$\Rightarrow \text{w.p. } \geq 1-\delta \quad \forall \pi \in \Pi \quad |\hat{V}(\pi) - V(\pi)| \leq \underbrace{\sqrt{\frac{4|\mathcal{X}| \log(2/\delta)}{3}}}_{:= \epsilon_{\mathcal{J}}}$$

$$\pi_{\star} \in \operatorname{argmax}_{\pi \in \Pi} V(\pi)$$

$$\begin{aligned}
V(\hat{\pi}) &= \underbrace{V(\hat{\pi}) - \hat{V}(\hat{\pi})}_{\geq -\varepsilon_{\mathcal{I}}} + \underbrace{\hat{V}(\hat{\pi}) - \hat{V}(\pi_*)}_{\geq 0} + \underbrace{\hat{V}(\pi_*) - V(\pi_*)}_{\geq -\varepsilon_{\mathcal{I}}} + V(\pi_*) \\
&\geq V(\pi_*) - 2\varepsilon_{\mathcal{I}}
\end{aligned}$$

$$\begin{aligned}
\text{Regret} &= \sum_{t=1}^T V(\pi_*) - V(\pi_t) \\
&= \sum_{t=1}^{\mathcal{I}} V(\pi_*) - V(\pi_t) + \sum_{t=\mathcal{I}+1}^T V(\pi_*) - V(\pi_t) \\
&\leq \mathcal{I} \cdot 1 + (T-\mathcal{I}) \underbrace{(V(\pi_*) - V(\hat{\pi}))}_{\leq 2\varepsilon_{\mathcal{I}}} \\
&= \mathcal{I} + 2(T-\mathcal{I}) \sqrt{\frac{4|\chi| \log(2/\pi/\delta)}{\mathcal{I}}} \\
&= O\left(T^{2/3} (|\chi| \log(2/\pi/\delta))^{1/3}\right)
\end{aligned}$$

$$\mathcal{I} = (|\chi| T^2 \log(2/\pi/\delta))^{1/3}$$

How do we find $\operatorname{argmax}_{\pi \in \Pi} \hat{V}(\pi)$ efficiently?

$$\begin{aligned} \hat{V}(\pi) &= \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} r_t \\ &= \frac{1}{3} \sum_{t=1}^3 \frac{(1 - \mathbb{1}\{x_t \neq \pi(c_t)\})}{P_t} r_t \\ &= \frac{1}{3} \sum_{t=1}^3 \frac{r_t}{P_t} - \frac{1}{3} \sum_{t=1}^3 \mathbb{1}\{x_t \neq \pi(c_t)\} \frac{r_t}{P_t} \end{aligned}$$

$$\operatorname{argmax}_{\pi} \hat{V}(\pi) = \operatorname{argmin}_{\pi} \frac{1}{3} \sum_{t=1}^3 \mathbb{1}\{x_t \neq \pi(c_t)\} \frac{r_t}{P_t}$$

Think of each $\pi: \mathcal{C} \rightarrow \mathcal{X}$ as a classifier over classes $\{1, \dots, |\mathcal{X}|\}$ and examples in \mathcal{C} .

$x_t \in \mathcal{X}$ is a "label", $c_t \in \mathcal{C}$ example, $\frac{r_t}{P_t}$ is weight

Minimize 0/1 loss on $\{(c_t, x_t, \frac{r_t}{P_t})\}_{t=1}^3$
 \uparrow
 weight of t th point

Ex. let $f_{\theta}: \mathcal{C} \rightarrow \mathbb{R}^{\mathcal{X}}$ be a neural network parameterized by parameters $\theta \in \mathbb{R}^d$

$$\text{loss}(\theta) = -\frac{1}{3} \sum_{t=1}^3 \frac{r_t}{P_t} \log \left(\frac{\exp(f_{\theta}(c_t)[x_t])}{\sum_{x \in \mathcal{X}} \exp(f_{\theta}(c_t)[x])} \right)$$

$$\text{Let } \hat{\theta}_3 = \underset{\theta}{\text{argmin}} \text{Loss}(\theta)$$

When c_t arrives we play

$$x_t = \underset{x \in \mathcal{X}}{\text{argmax}} \left[f_{\hat{\theta}_3}(c_t) \right]_x$$

Zooming out: γ -greedy achieves $R_T \leq T^{2/3} (|\mathcal{X}| \log(|\mathcal{T}|/\delta))^{1/3}$.

But $|\mathcal{C}|=1$ and $|\mathcal{T}|=|\mathcal{X}|$ in the standard MAB then we know $\sqrt{|\mathcal{X}|T}$ regret is possible.

→ For contextual Bandits, can we achieve $R_T \leq \sqrt{|\mathcal{X}|T \cdot \log(|\mathcal{T}|)}$?

Answer: yes.

An elimination alg for stochastic Contextual Bandits

Input \mathcal{T}, δ

$$\hat{\mathcal{T}}_1 = \mathcal{T}$$

for $l=1, 2, \dots$

$$\epsilon_l = \frac{1}{2^l}, \quad \begin{array}{ll} \# \text{ measurements} & \mathcal{I}_l \text{ TBD} \\ \text{regularization} & \gamma_l \text{ TBD} \end{array}$$

$$Q_l = \underset{Q \in \Delta_{\hat{\mathcal{T}}_l}}{\text{argmin}} \max_{\pi \in \hat{\mathcal{T}}_l} \mathbb{E}_c \left[\frac{1}{Q^x(\pi(c)|c)} \right]$$

$$\text{where } Q^x(x|c) := \gamma + (1-\gamma)x Q(x|c)$$

$$Q(x|c) := \sum_{\pi: \pi(c)=x} Q(\pi) \leftarrow \text{prob dist over actions}$$

$$T_e = \sum_{i=1}^e \mathcal{I}_i$$

for $t = T_{e-1} + 1, \dots, T_e$ ($T_e - T_{e-1} = \mathcal{I}_e$)

Nature reveals C_t

Player plays $x_t \sim Q_e^\gamma(\cdot | C_t)$ to get reward r_t

$$\hat{V}_e(\pi) = \frac{1}{\mathcal{I}_e} \sum_{t=1}^{\mathcal{I}_e} \frac{\mathbb{1}\{C_t = x_t\}}{P_t} r_t$$

$$\hat{\Pi}_{e+1} = \hat{\Pi}_e \setminus \{ \pi \in \hat{\Pi}_e : \max_{\pi'} \hat{V}(\pi') - \hat{V}(\pi) \geq 2\epsilon_e \}$$

// Q_e is chosen to minimize maximum variance of $\hat{V}(\pi)$:

$$Q_e \approx \underset{Q \in \Delta_{\mathcal{X}_e}}{\text{argmin}} \max_{\pi \in \hat{\Pi}_e} \mathbb{E} [(\hat{V}(\pi) - V(\pi))^2]$$

// Any $\pi \in \hat{\Pi}_e$ we have $V(\pi) - V(\pi^*) \leq 8\epsilon_e$
 \Rightarrow average regret incurred @ stage e is \uparrow

Lemma For any policy Π and dist. over contexts

$$\min_{Q \in \Delta_{\mathcal{X}}} \max_{\pi \in \Pi} \mathbb{E} \left[\frac{1}{Q(\pi(c)|c)} \right] \leq |\mathcal{X}|.$$

Furthermore, if $\gamma < \frac{1}{2K}$ then

$$\min_{Q \in \Delta_{\mathcal{X}}} \max_{\pi \in \Pi} \mathbb{E} \left[\frac{1}{Q^\gamma(\pi(c)|c)} \right] \leq 2|\mathcal{X}|.$$

This \uparrow lemma is a corollary of Kiefer-Wolfowitz which we will show next time.

Consequently we have by Bernstein

$$|\hat{V}_\ell(\bar{\pi}) - V(\bar{\pi})| \leq \sqrt{\frac{(2|\mathcal{X}|) \cdot 2 \log(2|\mathcal{X}|/\delta)}{\gamma_\ell}} + \frac{2 \log(2|\mathcal{X}|/\delta)}{\gamma_\ell^3 \mathfrak{J}_\ell}$$

$$\leq \sqrt{\frac{16|\mathcal{X}| \log(2|\mathcal{X}|/\delta)}{\mathfrak{J}_\ell}} =: \varepsilon_\ell$$

for $\gamma_\ell = \min \left\{ \frac{1}{2|\mathcal{X}|}, \sqrt{\frac{2 \log(2|\mathcal{X}|/\delta)}{9|\mathcal{X}| \mathfrak{J}_\ell}} \right\}$

$$\mathfrak{J}_\ell = 16|\mathcal{X}| \log(2|\mathcal{X}|/\delta) \varepsilon_\ell^{-2}$$

Lemma For all $\ell=1,2,\dots$ we have w.p.

$$\geq 1-\delta, \quad \bar{\pi}^* \in \hat{\mathcal{T}}_\ell \text{ and } \max_{\bar{\pi} \in \hat{\mathcal{T}}_\ell} V(\bar{\pi}^*) - V(\bar{\pi}) \leq 8\varepsilon_\ell.$$

Theorem w.p. $\geq 1-\delta$ we have

$$\sum_{\ell=1}^T V(\bar{\pi}^*) - V(\bar{\pi}_\ell) \leq c \sqrt{|\mathcal{X}| T \log(2|\mathcal{X}|/T/\delta)}.$$

Proof Fix $\nu \in [0,1]$ $\gamma_\ell \approx \frac{\varepsilon_\ell}{|\mathcal{X}|}$

$$\begin{aligned}
\sum_{t=1}^T V(\alpha_t^*) - V(\alpha_t) &\leq \nu T + \sum_{t=1}^{\log_2 \delta \nu^{-1}} (\gamma_t |\mathcal{X}| + (1 - \gamma_t |\mathcal{X}|) \cdot \delta \varepsilon_t) \mathcal{J}_t \\
&\lesssim \nu T + \sum_{t=1}^{\log_2 \nu^{-1}} \varepsilon_t \cdot \mathcal{J}_t \\
&\lesssim \nu T + \sum_{t=1}^{\log_2 \nu^{-1}} \varepsilon_t^{-1} |\mathcal{X}| \log(|\mathcal{X}| T / \delta) \\
&\leq \nu T + \frac{1}{\nu} |\mathcal{X}| \log(|\mathcal{X}| T / \delta)
\end{aligned}$$

minimize over ν yields the result.

Proof of Kiefer-Wolfowitz generalization.

Recall Kiefer-Wolfowitz says for $\mathcal{X} \subset \mathbb{R}^d$

$$\inf_{\lambda \in \Delta_{\mathcal{X}}} \max_{x \in \mathcal{X}} x^T \left(\sum_{x'} \lambda_{x'} x x'^T \right)^{-1} x \leq d.$$

Lemma Let $\eta \in S$ be a.R.V. and

let $\phi(x, z) \in \mathbb{R}^d$ be a feature map.

Then

$$\min_{\lambda \in \Delta_X} \max_{x \in X} \mathbb{E}_z \left[\phi(x, z)^T \left(\sum_{x'} \lambda_{x'} \phi(x', z) \phi(x', z)^T \right)^{-1} \phi(x, z) \right] \leq d$$

(Recovers KW w/ $\phi(x, z) = x$).

Key idea proof: $\frac{\partial}{\partial t} \log \det(A(t)) = \text{Trace} \left(A(t)^{-1} \frac{\partial}{\partial t} A(t) \right)$

$$f(\lambda) = \mathbb{E}_z \left[\log \det \left(\sum_{x'} \lambda_{x'} \phi(x', z) \phi(x', z)^T \right) \right]$$

$$\lambda^* = \text{argmax}_{\lambda} f(\lambda)$$

$$0 \geq \langle \nabla f(\lambda^*), e_x - \lambda^* \rangle$$

$$\geq \mathbb{E}_z \left[\phi(x, z)^T \left(\sum_{x'} \lambda_{x'} \phi(x', z) \phi(x', z)^T \right)^{-1} \phi(x, z) \right] - d$$

$$\langle \nabla f(\lambda), \lambda \rangle \leq d \quad \text{for all } \lambda.$$

To prove $\min_{Q \in \Delta_\pi} \max_{\pi} \mathbb{E}_c \left[\frac{1}{Q(a(c)|c)} \right] \leq |X|$

Consider actions $i=1, 2, \dots, |X|$ and define

$$\pi_c = e_{\pi(c)} \in \{0,1\}^{|\mathcal{X}|}$$

$$(i.e. \phi(\pi, c) = e_{\pi(c)})$$

$$\mathbb{E}_c \left[\pi_c^T \left(\sum_{\pi'} Q(\pi') \pi'_c \pi'_c{}^T \right)^{-1} \pi_c \right]$$

diagonal matrix
w/ entries only summed
over the π

$$\leq |\mathcal{X}|$$

Okay, $\exists Q \in \Delta_{\mathcal{T}}$ s.t. variance bounded by $|\mathcal{X}|$.

In linear bandits we found a sparse Q .

w/ contextual bandits $|\text{support}(Q^*)| \lesssim \min\{\frac{1}{\delta}, |\mathcal{X} \cdot |\mathcal{C}|\}$

Computational Effic. Algorithm for context bandits.

Input Π, δ

Play uniform at random for $\log_2(\Pi)$ steps to get $\hat{\Delta}_0(\pi)$
for $l=1, 2, \dots$

$$\epsilon_l = \frac{1}{2^l}, \quad \begin{array}{l} \# \text{ measurements } \mathcal{I}_l \text{ TBD} \\ \text{regularization } \gamma_l \text{ TBD} \end{array}$$

Q_l is any $Q \in \Delta_\Pi$ such that

$$1) \quad \sum_{\pi} \hat{\Delta}_{l-1}(\pi) Q(\pi) \leq C'' \epsilon_l$$

$$2) \quad \sqrt{\mathbb{E}_c \left[\frac{1}{Q^{\gamma_l}(\pi(c)|c)} \right] \frac{C \log_2(\Pi/\delta)}{\mathcal{I}_l}} \leq \epsilon_l + \hat{\Delta}_{l-1}(\pi)$$

$$\text{where } Q^{\gamma}(x|c) := \gamma + (1-\gamma)Q(x|c)$$

$$Q(x|c) := \sum_{\pi: \pi(c)=x} Q(\pi) \leftarrow \text{prob dist over actions}$$

$$T_l = \sum_{i=1}^l \mathcal{I}_i$$

for $t = T_{l-1} + 1, \dots, T_l$ ($T_l - T_{l-1} = \mathcal{I}_l$)

Nature reveals c_t

Player plays $x_t \sim Q_l^{\gamma}(\cdot | c_t)$ to get reward r_t

$$\hat{V}_l(\pi) = \frac{1}{\mathcal{I}_l} \sum_{t=1}^{\mathcal{I}_l} \frac{\mathbb{1}\{\pi(c_t)=x_t\}}{P_t} r_t, \quad \hat{\Delta}_l(\pi) = \max_{\pi'} \hat{V}_l(\pi') - \hat{V}_l(\pi)$$

By induction one can show

$$\hat{V}_l(\pi) - V(\pi) \leq C' \max\{\epsilon_l, \Delta(\pi)\}$$

$$\hat{\Delta}_l(\pi) \geq \hat{V}_l(\pi^*) - \hat{V}_l(\pi)$$

$$\geq \Delta(\pi) - 2 \max\{\epsilon_l, \Delta(\pi)\}$$

which is $\geq \Delta(\pi) / 2$ when $\epsilon_l < \frac{\Delta(\pi)}{8}$

If $\epsilon_l < \Delta(\pi)$ then $\hat{\Delta}_l(\pi) \geq \Delta(\pi)$

otherwise $\hat{\Delta}_l(\pi) \leq \epsilon_l$

\Rightarrow If Q_l feasible is found then average regret at stage l is just

$$|X| \gamma_l + C'' \epsilon_l \leq (1 + C'') \epsilon_l$$

\Rightarrow Regret analysis is identical to the comp. inf. alg.

So we just need to find
a feasible Q_ℓ efficiently.

Idea: Let $\Pi_j = \{\pi: \Delta(\pi) \leq \epsilon_j\}$. If P_j is

$$P_j = \underset{P \in \Delta_\Pi}{\text{arg min}} \quad \max_{\pi \in \Pi_j} \mathbb{E} \left[\frac{1}{P^\ell(\pi(c)|c)} \right]$$

and $\bar{P}_\ell = \frac{1}{\sum_j \alpha_j} \sum_{j=1}^L \alpha_j P_j.$

Some algebra shows

$$\sum_{\pi} \Delta_\pi \bar{P}_\ell(\pi) \leq c \epsilon_\ell \quad (1)$$

and $\max_{\pi \in \Pi} \mathbb{E} \left[\frac{1}{\bar{P}_\ell(\pi(c)|c)} \right]$ satisfies (2).

Thus \bar{P}_ℓ is feasible for Q_ℓ

Find a \bar{w} that violates (2)

using cost-sensitive / weighted
classification.

iterations per opt problem

$$\leq \frac{1}{\gamma \epsilon}.$$