Stochastic multi-armed bandits

Protocol Repeat for $t = 1, 2, \ldots, T$: LEARNER plays action $I_t \in \{1, \ldots, K\}$. Environment generates rewards $X_{t,i} \sim v_i$. LEARNER gains and observes reward $X_{t,I_{+}}$. Notation: Rean rewards: $\mu_i = \mathbb{E}[X_{t,i}]$ st best arm: $i^* = \arg \max_{i \in [K]} \mu_i$ Rean reward of best arm: $\mu^* = \max_{i \in [K]} \mu_i$ suboptimality gaps: $\Delta_i = \mu^* - \mu_i$ \aleph number of draws of arm i until round t: $N_{t,i}$ GOAL: minimize regret $R_{T} = \mu^{*}T - \sum_{t=1}^{l} \mathbb{E}\left[X_{t,I_{t}}\right] = \sum_{i=1}^{K} \Delta_{i}\mathbb{E}\left[N_{T,i}\right]$ Assumption: σ^2 -subgaussian rewards $\mathbb{E}\left[e^{y(X_{t,i}-\mathbb{E}[X_{t,i}])}\right] \leq e^{\sigma^2 y^2/2}$ Lower bound: For Gaussian rewards: $R_T \gtrsim \sum_{i=/3*} \frac{\sigma^2 \log T}{\Delta_i} \qquad \qquad R_T \gtrsim \sigma \sqrt{KT}$

The elephant in the room

Boltzmann exploration Initialize: $\hat{\mu}_{1,i} = 0$ for all $i \in [K]$. Repeat for $t = 1, 2, \ldots, T$: Compute distribution $p_{t,i} \propto e^{\eta_t \hat{\mu}_{t,i}}$. Play action $I_t \sim p_t$ and observe r_{t,I_t} . Update empirical means $\hat{\mu}_{t,i} = \frac{\sum_{s=1}^{t} X_{s,i} \mathbb{I}\{I_s = i\}}{N_{t,i}}.$

Broadly used exploration strategy in RL, but very little theory to support it!

Boltzmann Exploration Done Right

Claudio Gentile INRIA Lille – Nord Europe Villeneuve d'Ascq, France

Gábor Lugosi ICREA & Universitat Pompeu Fabra Barcelona, Spain

Boltzmann exploration done wrong

Main result

For any monotone sequence of learning rates η_t , Boltzmann exploration will suffer suboptimal regret.

Two regimes with no middle ground:

- η_t grows too slowly \Rightarrow too much time on exploration / too slow to zoom in on i^*
- $\lambda \eta_t$ grows too quickly \Rightarrow high probability of missing i^*

Proposition 1: over-exploration Regret on any 2-armed bandit problem with known means: $\lambda \eta_t = \frac{\log(t\Delta^2)}{\Delta} \Rightarrow R_T \approx \frac{\log T}{\Delta}.$ $\$ \eta_{t} = \frac{\log(t\Delta^{2})}{(1+\alpha)\Delta} \Rightarrow R_{T} \approx T^{\frac{\alpha}{(1+\alpha)}} \cdot \left(\frac{1}{\Delta}\right)^{\frac{1-\alpha}{1+\alpha}}.$

Proof idea:

$$\begin{split} \mathsf{R}_\mathsf{T} &= \sum_{t=1}^\mathsf{T} \mathbb{P}\left[\mathsf{I}_t = 2\right] = \sum_{t=1}^\mathsf{T} \frac{1}{1 + e^{\eta_t \Delta}} \geq \sum_{t=1}^\mathsf{T} \frac{e^{-\eta_t \Delta}}{2} \\ &= \frac{1}{2} \sum_{t=1}^\mathsf{T} \left(t \Delta^2 \right)^{1+\alpha} \\ &\approx \begin{cases} \frac{\log \mathsf{T}}{\Delta}, & \text{if } \alpha = 0, \\ \mathsf{T}^{\frac{\alpha}{(1+\alpha)}} \cdot \left(\frac{1}{\Delta} \right)^{\frac{1-\alpha}{1+\alpha}}, & \text{if } \alpha > 0. \end{cases} \end{split}$$

Proposition 2: under-exploration There exists a 2-armed stochastic bandit problem where BE using any $\eta_t > 2 \log t$ has regret $R_T = \Omega(T)$.

Proof idea:

- \bigstar Two arms: $X_{t,2} = \frac{1}{2}$ and $X_{t,1} \sim \text{Bernoulli}\left(\frac{1}{2} + \Delta\right)$
- Bad event:

 $E_0 = \{arm \ 1 \text{ gives } 0 \text{ reward in first } t_0 \text{ rounds} \}$

- \aleph Under E₀, BE will not draw arm 1 after round t_0 due to η_t growing too fast
- $\mathbb{P}\left[\mathsf{E}_{0}\right] \geq \left(\frac{1}{2} \Delta\right)^{\mathsf{t}_{0}} = \mathsf{const.}$

Gergely Neu Universitat Pompeu Fabra Barcelona, Spain

A quick fix

$$\begin{split} \text{Theorem} \\ \text{Let } \tau &= \frac{16e\mathsf{K}\log\mathsf{T}}{\Delta^2}. \ \text{Then the regret of BE with} \\ \text{the learning rate sequence} \\ \eta_t &= \mathbb{I}\{t < \tau\} + \frac{\log(t\Delta^2)}{\Delta}\mathbb{I}\{t \geq \tau\} \text{ satisfies} \\ \mathsf{R}_\mathsf{T} &\leq \frac{16e\mathsf{K}\log\mathsf{T}}{\Delta^2} + \frac{9\mathsf{K}}{\Delta^2}. \end{split}$$

* near-optimal performance guarantees ③ $\texttt{\texttt{K}}$ requires prior knowledge of Δ and T $\odot \odot \odot$

Boltzmann exploration done right

What's wrong with Boltzmann exploration? It doesn't reason about uncertainty of

reward estimates!

Our solution: arm-dependent learning rates!

Key tool: "Gumbel-softmax trick"

$$I_t = \arg \max_{i \in [K]} \{ \hat{\mu}_{t,i} + Z_{t,i} \},$$

follows Boltzmann distribution if $Z_{t,i}$ are i.i.d. standard Gumbel random variables. Idea: account for uncertainty by scaling $Z_{t,i}$ differently for each arm!

Boltzmann–Gumbel exploration Initialize: $\hat{\mu}_{1,i} = 0$ for all $i \in [K]$. Repeat for $t = 1, 2, \ldots, T$:

- Draw $Z_{t,i}$ i.i.d. from standard Gumbel distribution.
- Play action

$$I_{t} = \operatorname{arg\,max}_{i \in [K]} \left\{ \widehat{\mu}_{t,i} + \sqrt{\frac{C^{2}}{N_{t,i}}} \cdot Z_{t,i} \right\}$$

Observe $r_{t,I_{+}}$ and update empirical means

$$u_{t,i} = \frac{\sum_{s=1}^{t} X_{s,i} \mathbb{I}\{I_s = i\}}{N_{t,i}}$$

For
BGI
$$R_T \lesssim$$

Proof sketch:

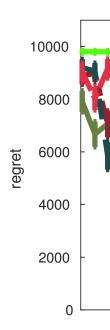
Let
Set
Key
$$E_{t,i}^{\hat{\mu}}$$

 $E_{t,i}^{\hat{\mu}}$
 $\mathbb{E}[\mathbb{N}]$

per
$$\stackrel{1}{\times} \sum_{t=1}^{T}$$

 $\beta_{t,i} = \sqrt{C^2/N_{t,i}}$ and $\tilde{\mu}_{t,i} = \hat{\mu}_{t,i} + \beta_{t,i}Z_{t,i}$ thresholds $x_i = \mu_i + \frac{\Delta}{3}$ and $y_i = \mu_1 - \frac{\Delta}{3}$ events: $= \{ \hat{\mu}_{t,i} \leq x_i \} \sim \text{arm i well-estimated}$ $= \{ \tilde{\mu}_{t,i} \leq y_i \} \sim \text{small perturbation on arm } i$ $J_{t,i}$] decomposed into 3 terms: $\sum_{t=1}^{\Gamma} \mathbb{P}\left[I_t = i, E_{t,i}^{\tilde{\mu}}, E_{t,i}^{\hat{\mu}}\right] \sim \text{interaction between}$ erturbations $Z_{t,1}$ and fluctuations of $\hat{\mu}_{t,1}$ $\sum_{t=1}^{I} \mathbb{P}\left[I_t = i, E_{t,i}^{\hat{\mu}}, E_{t,i}^{\hat{\mu}}\right] \sim \text{large perturbations}$ $\sum_{t=1}^{T} \mathbb{P}\left[I_t = i, \overline{E_{t,i}^{\hat{\mu}}}\right] \sim \text{large deviations}$

First and last terms bounded by
$$\begin{split} \sum_{k=0}^{T-1} \mathbb{E} \left[\exp \left(\frac{\mu_{i} - \hat{\mu}_{\tau_{k},i}}{\beta_{\tau_{k},i}} \right) \right] e^{-\frac{\Delta_{i}\sqrt{k}}{3C}} &\leq e^{\sigma^{2}/2C^{2}} \cdot \sum_{k=1}^{T-1} e^{-\frac{\Delta_{i}\sqrt{k}}{3C}} \\ &\leq \frac{18C^{2}e^{\sigma^{2}/2C^{2}}}{\Delta_{i}^{2}} \end{split}$$
Middle term bounded by



Analysis

Theorem	-
σ^2 -subgaussian rewards	s, the regret of
E with $C = \sigma$ satisfies	
$\lesssim \sum_{i \neq i^*} \frac{\sigma^2 log^2 \left(T \Delta_i^2 / \sigma^2 \right)}{\Delta_i}$	$R_T \lesssim \sigma \sqrt{KT} log K$

$$\frac{2C^2\log^2_+(T\Delta_i^2/c^2)+c^2e^{\gamma}}{\Delta_i^2}$$

for any c > 0

Technique extends to other subgaussian mean estimators for heavy-tailed rewards

Empirical illustration

Sensitivity of BGE and BE with learning rates $\eta_t = C, \, \eta_t = C/\log t, \, \eta_t = C/\sqrt{t}$ to various settings of C:

