

1. HW3 almost graded. soln is up
(DDL \rightarrow content \rightarrow HW3-solns)

2. HW4 posted. (course website)

due on May 4, 11:59 pm

no late HW will be accepted.

3. Final presentation (preferences on slot Apr 29 / May 4)
send me an email.

discuss presentation to our OHS.

Thm: UCB satisfies:

$$\textcircled{1} R_T \leq \sum_{a: \Delta(a) > 0} \frac{16 \ln T}{\Delta(a)} + 3K \quad (\text{gap dependent})$$

$$\textcircled{2} R_T \leq \tilde{O}(\sqrt{TK}) \quad (\text{gap independent})$$

$$\text{Lemma 2} \quad \forall a. \quad \mathbb{E} [m_T(a)] \leq \frac{16 \ln T}{\Delta(a)^2} + 3.$$

pf of $\textcircled{2}$:

given a "cutoff" $\Delta > 0$. group the arms

based on $\Delta(a) \leq \Delta$ or $\Delta(a) > \Delta$.

$$R_T = \sum_a \mathbb{E}[m_T(a)] \cdot \Delta(a)$$

$$= \sum_{a: \Delta(a) \in (0, \Delta]} \mathbb{E}[m_T(a)] \Delta(a) + \sum_{a: \Delta(a) > \Delta} \dots$$

$\leq \Delta$.

$$\leq T \Delta + \sum_{a: \Delta(a) > \Delta} \left(\Delta(a) \frac{b \ln T}{\Delta(a)^2} + 3 \Delta(a) \right)$$

$\leq \frac{b \ln T}{\Delta}$

$$\leq T \Delta + \frac{b k \ln T}{\Delta} + 3k$$

$$T \Delta = \frac{k \ln T}{\Delta} \Rightarrow \Delta = \sqrt{\frac{k \ln T}{T}}$$

we focus on $T \geq k$, so

$$\Rightarrow R_T \leq O(\sqrt{T k \ln T})$$

$\sqrt{T k} \geq k$

§ adversarial multi-armed bandits

what if the losses are non-stationary, or even does not come from a distribution?

can we still give algorithms w/ regret guarantees?

random.

Recall MAB:

For $t=1, 2, \dots, T$: ^{not} l_t 's are chosen before round 1. (oblivious adversary)

environment gives $l_t \in [0, 1]^K$.

learner selects action $a_t \in \{1, \dots, K\}$ ^{randomly} $\sim P_t \in \Delta^{K-1}$.

learner suffers loss $l_t(a_t)$. $(l_t(a_t), a_t = a)$ may not be fixed.

Goal: minimize regret.

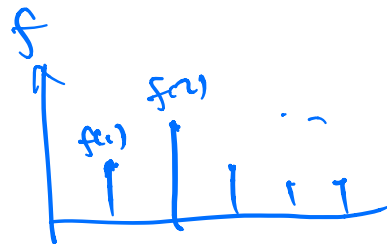
$$R_T = \mathbb{E} \left[\sum_{t=1}^T l_t(a_t) \right] - \min_a \sum_{t=1}^T l_t(a)$$

How to design algorithms w/ sublinear R_T ?

alternative representation of R_T .

$$\mathbb{E}_{a_t \sim P_t} l_t(a_t) = \sum_{a=1}^K P_t(a) \cdot l_t(a) = \langle P_t, l_t \rangle$$

$$\min_a \sum_{t=1}^T l_t(a) = \min_{P \in \Delta^{K-1}} \langle P, \sum_{t=1}^T l_t \rangle$$



$$f_t(P) = \langle P, l_t \rangle$$

f

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \langle P_t, l_t \rangle - \min_{P \in \Delta^{K-1}} \sum_{t=1}^T \langle P, l_t \rangle \right]$$

\sqrt{TK} regret for weather forecasting using OMD (l_t 's are known).

Q: can we re-use the OMD algorithm to develop algorithms w/ low R_T ?

$$P_{\text{opt}} = \underset{P \in \Delta^{K-1}}{\text{argmin}} \eta \cdot \langle \ell_t, P \rangle + D_f(P, P_t)$$

η
 \downarrow
 need to be changed.

$$\frac{a_t \sim P_t \cdot \ell_t(a_t)}{\ell_t}$$

Thompson sampling
 $\ell(a)$

do OMD on unbiased estimators of ℓ_t 's.

$$\tilde{\ell}_t \cdot \mathbb{E}_{a_t \sim P_t} [\tilde{\ell}_t] = \ell_t$$

Low regret wrt $\tilde{\ell}_t$'s \Rightarrow ^{taking expectation} low regret wrt ℓ_t 's.

Constructing $\tilde{\ell}_t$:

$$\tilde{\ell}_t(a) = \begin{cases} 0 & a_t \neq a \\ x = \frac{\ell_t(a)}{P_t(a)} & a_t = a \end{cases}$$

$\ell_t(a_t)$ \rightarrow can be observed.

$$\mathbb{E}_{a_t \sim p_t} [\tilde{Q}_t(a)] = \frac{l_t(a_t)}{p_t(a_t)}$$

$$= (1 - p_t(a_t)) \cdot 0 + p_t(a_t) \cdot x = l_t(a_t)$$

$$\tilde{Q}_t(a) = \frac{l_t(a)}{p_t(a)} \mathbb{I}(a = a_t)$$

7 OMD with \tilde{Q}_t 's w/ $\Omega = \Delta^{k-1}$

$$\psi(p) = \sum_{a=1}^k p(a) \ln p(a) \quad p_1 = \left(\frac{1}{k}, \dots, \frac{1}{k}\right)$$

$$p_{t+1}(a) \stackrel{(*)}{\propto} p_t(a) \exp(-\eta \tilde{Q}_t(a)) \propto \exp\left(-\eta \sum_{s=1}^t \tilde{Q}_s(a)\right)$$

→ EXP3 algorithm.

achieving exploration / exploitation tradeoff?

① $\sum_{s=1}^t \tilde{Q}_s(a) \approx \sum_{s=1}^t l_s(a) \Rightarrow$ exploit $l_t(a) \in [-1, 0]$.
the algorithm will behave greedily

② $\eta < \infty$.

$a_t = a$.

$$\tilde{Q}_t = (0, \dots, 0, \frac{l_t(a)}{p_t(a)}, \dots, 0)$$

↓

$h_t(x)$

P_{t+1} skew toward action other than a_t implies it encourages taking other actions.

Analysis of EXP3:

applying OMD guarantees namely: RHS

$$\sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle - \sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\tilde{a}_t\|_{P_t}^2$$

$\mathbb{E} \left[\sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle \right]$
 $\mathbb{E} \left[\sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle \right] = \sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle$

$\mathbb{E} \left[\sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle \right]$
 $\mathbb{E} \left[\sum_{t=1}^T \langle P_t, \tilde{a}_t \rangle \right]$

$\mathbb{E}_{P_t} \left(\frac{\|\tilde{a}_t\|_{P_t}^2}{P_t(a_t)} \right)$

$\leq \sum \frac{1}{P_t(a_t)}$

cannot be well controlled.

* get a better regret bound for OMD w/

neg-entropy regularizer.

Lemma (Local norm bound)

$\Omega = \Delta^{K-1}$. $\psi(w) = \text{neg-entropy}$

learning rate η . on $\left\{ f_t^*(w) = \langle w, g_t \rangle \right\}_{t=1}^T$ where

$\dots \rightarrow K$

$$g_t \in [0, \infty)$$

$$\sum_{t=1}^T \langle w_t, g_t \rangle - \sum_{t=1}^T \langle w, g_t \rangle \leq \frac{\ln K}{\eta} + \eta \sum_{t=1}^T \|g_t\|_{\text{diag}(w_t)}^2$$

$$\text{diag}(w_t) = \begin{pmatrix} w_t(1) & & \\ & \ddots & \\ & & w_t(K) \end{pmatrix}$$

$$\sum_a w_t(a) g_t(a)^2 \leq \|g_t\|_{\infty}^2$$

Applying Lemma to EXP:

$$\text{RHS} = \frac{\ln K}{\eta} + \eta \sum_{t=1}^T \sum_{a=1}^K p_t(a) \cdot \tilde{l}_t(a)^2$$

$$= \frac{\ln K}{\eta} + \eta \sum_{t=1}^T p_t(a_t) \left(\frac{l_t(a_t)}{p_t(a_t)} \right)^2$$

$$\leq \frac{\ln K}{\eta} + \eta \sum_{t=1}^T \underbrace{\frac{1}{p_t(a_t)}}_K \rightarrow \text{better than } \frac{1}{(p_t(a_t))^2}$$

$$\mathbb{E}[\text{RHS}]$$

$$= \frac{\ln K}{\eta} + \eta \sum_{t=1}^T \mathbb{E} \left[\frac{1}{p_t(a_t)} \right]$$

$$\mathbb{E}_{a_t \sim P_t} \left[\frac{1}{P_t(a_t)} \right] = \sum_{a \in \mathcal{A}} P_t(a) \cdot \frac{1}{P_t(a)} = K$$

$$= \frac{\ln K}{\eta} + \eta T \cdot K$$

$$\eta = \sqrt{\frac{\ln K}{TK}}$$

$$= 2 \sqrt{TK \ln K}$$

In conclusion:

Thm: EXP3 w/ $\eta = \sqrt{\frac{\ln K}{TK}}$ has
 regret $R_T \leq O(\sqrt{TK \ln K})$.

PF of Lemma:

Recall in OMD:

$$\langle \underbrace{w_{t+1}}_w, -w_t, \eta g_t \rangle \leq D_\psi(w, w_t) - D_\psi(w, w_{t+1})$$

$$- D_\psi(w_{t+1}, w_t)$$

$$\langle w_t - w, \eta g_t \rangle$$

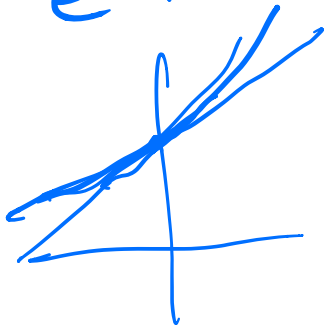
$$\leq \underbrace{\langle w_t - w_{t+1}, \eta g_t \rangle}_{\textcircled{1}} + D_f(w, w_t) - D_f(w, w_{t+1})$$

$$w_{t+1}(a) = \frac{w_t(a) \cdot e^{-\eta g_t(a)}}{\underbrace{\sum_{a'} w_t(a') e^{-\eta g_t(a')}}_{\textcircled{Z_t} \leq 1}}$$

$$\geq w_t(a) \cdot e^{-\eta g_t(a)}$$

$$\textcircled{1} \leq \sum_{a=1}^K w_t(a) \underbrace{\left(1 - e^{-\eta g_t(a)}\right)}_{\geq 1 - \eta g_t(a)} \underbrace{\eta g_t(a)}$$

$$e^x \geq 1+x \leq \eta \sum_{a=1}^K w_t(a) \cdot g_t(a)^2$$



Summing over t , divide by η

~~✗~~

\Rightarrow Lerner

Next class:

stochastic linear bandits.

(contextual)

$$a_1 \rightarrow l_1 = (0 \dots 0, \underbrace{1}_{a_1} \dots 0)$$

$$a_2 \rightarrow l_2 = (0 \dots 0, \underbrace{1}_{a_2} \dots 0)$$

\forall Alg. \exists adversary.

$$R_T \geq \Omega(T).$$

Cover's impossibility result $\Omega(T)$.