

Machine Learning Theory (CS 6783)

Lecture 16: Learning With Bandit Feedback

1 Online Linear Optimization With Bandit Feedback

The bandit linear online learning problem is as follows:

For $t = 1$ to n

1. Learner picks $\hat{\mathbf{y}}_t \in \mathcal{F}$
2. Adversary picks ∇_t simultaneously
3. Learner observes and suffers loss $\hat{\mathbf{y}}_t^\top \nabla_t$

End For

Goal: Minimize expected regret

$$n\text{Reg}_n = \sum_{t=1}^n \hat{\mathbf{y}}_t^\top \nabla_t - \inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \nabla_t$$

Main Idea:

1. Obtain $\hat{\mathbf{y}}_t$ from a full information algorithm
2. Randomize move such that $\mathbb{E}[\hat{\mathbf{y}}_t] = \hat{\mathbf{y}}_t$
3. Play $\hat{\mathbf{y}}_t$ and receive feedback $\hat{\mathbf{y}}_t^\top \nabla_t$
4. Build unbiased estimate of ∇_t based on feedback as $\mathbb{E}[\tilde{\nabla}_t] = \nabla_t$
5. Feed $\tilde{\nabla}_t$ to a full information online linear algorithms

For bandit algorithms, adaptive vs oblivious adversaries make a difference. Adaptive adversaries are ones that know the internal randomization of the learner up to given point while oblivious adversaries only know the learning algorithms but not the random bits produced by the adversary. That is, we can think of adversary (knowing the learning algorithm) first prefixing $\nabla_1, \dots, \nabla_n$ and producing them one by one. In this case, note that:

$$\begin{aligned}
n\mathbb{E}[\text{Reg}_n] &= \mathbb{E} \left[\sum_{t=1}^n \hat{y}_t^\top \nabla_t \right] - \inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \nabla_t \\
&= \mathbb{E} \left[\sum_{t=1}^n \hat{y}_t^\top \nabla_t \right] - \inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \mathbb{E}[\tilde{\nabla}_t] \\
&\leq \mathbb{E} \left[\sum_{t=1}^n \hat{y}_t^\top \nabla_t \right] - \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \tilde{\nabla}_t \right] \\
&= \mathbb{E} \left[\sum_{t=1}^n \mathbb{E}[\hat{y}_t]^\top \nabla_t \right] - \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \tilde{\nabla}_t \right] \\
&= \mathbb{E} \left[\sum_{t=1}^n \hat{y}_t^\top \nabla_t \right] - \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \tilde{\nabla}_t \right] \\
&= \mathbb{E} \left[\sum_{t=1}^n \hat{y}_t^\top \mathbb{E}[\tilde{\nabla}_t] \right] - \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \tilde{\nabla}_t \right] \\
&= \mathbb{E} \left[\sum_{t=1}^n \hat{y}_t^\top \tilde{\nabla}_t \right] - \inf_{f \in \mathcal{F}} \sum_{t=1}^n f^\top \tilde{\nabla}_t \\
&= \mathbb{E} \left[n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \right]
\end{aligned}$$

Hence overall we conclude that for this procedure:

$$n\mathbb{E}[\text{Reg}_n] \leq \mathbb{E} \left[n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \right] \quad (1)$$

The above basically tells us that we have a reduction from bandit algorithm to full information algorithm. Hence, using this unbiased gradient estimate trick, we can apply a full information algorithm on the estimates and the expected regret of this full information algorithm will be the bound for our bandit algorithm. A word of caution though: typically our full information algorithms depend on norms of gradient being bounded under some appropriate norm. Eg. exponential weights algorithm on ℓ_∞ norm and gradient descent on ℓ_2 norm. However, while ∇_t 's might have this norm bounded, our estimates can have very large norms and this can cause our bounds to blow up. To this end, we have two options: either we modify our full information algorithm (albeit at the cost of worse bounds) so that the estimates have smaller norms or alternatively we are more careful to get an adaptive bound for our full information algorithm to get tighter bounds in expectation. We will now pursue the latter approach and tighten our mirror descent algorithm to have bounds in terms of so called local norms.

2 Mirror Descent with Local Norms (full information case)

We have shown that if we are able to find a function R that is strongly convex w.r.t. some norm $\|\cdot\|$ then mirror descent algorithm with step size η using this function R has the following bound on regret:

$$n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \leq \frac{\eta}{2} \sum_{t=1}^n \|\tilde{\nabla}_t\|_*^2 + \frac{1}{\eta} \sup_{f \in \mathcal{F}} \Delta_R(f|\hat{y}_1)$$

where $\|\cdot\|_*$ is the dual to the norm $\|\cdot\|$. We will now modify this result to replace the dual norm with a local norm. This will turn out to be useful to obtain bandit algorithms.

Recall the mirror descent algorithm:

$$\nabla R(\hat{\mathbf{y}}'_{t+1}) = \nabla R(\hat{\mathbf{y}}_t) - \eta \nabla_t \quad \& \quad \hat{\mathbf{y}}_{t+1} = \operatorname{argmin}_{f \in \mathcal{F}} \Delta_R(f | \hat{\mathbf{y}}'_{t+1})$$

Assume that the function R is twice differentiable and let $\nabla^2 R(f)$ denote the Hessian of the function at a point R . Recall the following claim.

Lemma 1. *For any twice differentiable convex R , if we run mirror descent using step size η , then*

$$n \operatorname{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \leq \frac{\eta}{2} \sum_{t=1}^n \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)}^2 + \frac{1}{\eta} \sup_{f \in \mathcal{F}} \Delta_R(f | \hat{\mathbf{y}}_1)$$

where z_t is some convex combination of $\hat{\mathbf{y}}_t$ and $\hat{\mathbf{y}}'_{t+1}$ (here matrix M , $\|x\|_M^2 = x^\top M x$)

Proof. We will recall the upper bound from the mirror descent proof of the form:

$$\langle \tilde{\nabla}_t, \hat{\mathbf{y}}_t - f^* \rangle \leq \langle \tilde{\nabla}_t, \hat{\mathbf{y}}_t - \hat{\mathbf{y}}'_{t+1} \rangle + \frac{1}{\eta} (\Delta_R(f^* | \hat{\mathbf{y}}_t) - \Delta_R(f^* | \hat{\mathbf{y}}_{t+1}) - \Delta_R(\hat{\mathbf{y}}'_{t+1} | \hat{\mathbf{y}}_t))$$

Now the key trick is that we start with the definition of Bregman divergence and use Taylor's theorem. Note that:

$$\Delta_R(\hat{\mathbf{y}}'_{t+1} | \hat{\mathbf{y}}_t) = R(\hat{\mathbf{y}}'_{t+1}) - R(\hat{\mathbf{y}}_t) - \langle R(\hat{\mathbf{y}}_t), \hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t \rangle$$

Now using Taylor's theorem (+ intermediate value theorem) there exists a point z_t that is some convex combination of $\hat{\mathbf{y}}'_{t+1}$ and $\hat{\mathbf{y}}_t$ such that

$$R(\hat{\mathbf{y}}'_{t+1}) - R(\hat{\mathbf{y}}_t) - \langle R(\hat{\mathbf{y}}_t), \hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t \rangle = \frac{1}{2} (\hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t)^\top \nabla^2 R(z_t) (\hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t) = \frac{1}{2} \|\hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t\|_{\nabla^2 R(z_t)}^2$$

Hence using this we can conclude that

$$\langle \tilde{\nabla}_t, \hat{\mathbf{y}}_t - f^* \rangle \leq \langle \tilde{\nabla}_t, \hat{\mathbf{y}}_t - \hat{\mathbf{y}}'_{t+1} \rangle + \frac{1}{\eta} (\Delta_R(f^* | \hat{\mathbf{y}}_t) - \Delta_R(f^* | \hat{\mathbf{y}}_{t+1})) - \frac{1}{2\eta} \|\hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t\|_{\nabla^2 R(z_t)}^2$$

Now note that for any invertible matrix M , $\|\cdot\|_{M^{-1}}$ is the dual norm to the norm $\|\cdot\|_M$ and hence using the fact (as we did in the earlier mirror descent proof) that

$$\langle \tilde{\nabla}_t, \hat{\mathbf{y}}_t - \hat{\mathbf{y}}'_{t+1} \rangle \leq \frac{\eta}{2} \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)}^2 + \frac{1}{2\eta} \|\hat{\mathbf{y}}'_{t+1} - \hat{\mathbf{y}}_t\|_{\nabla^2 R(z_t)}^2$$

we conclude that

$$\langle \tilde{\nabla}_t, \hat{\mathbf{y}}_t - f^* \rangle \leq \frac{\eta}{2} \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)}^2 + \frac{1}{\eta} (\Delta_R(f^* | \hat{\mathbf{y}}_t) - \Delta_R(f^* | \hat{\mathbf{y}}_{t+1}))$$

Summing over t and simplifying the telescoping sum over the Bregman divergences we we obtain that

$$\begin{aligned} n \operatorname{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) &\leq \frac{\eta}{2} \sum_{t=1}^n \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)}^2 + \frac{1}{\eta} (\Delta_R(f^* | \hat{\mathbf{y}}_1) - \Delta_R(f^* | \hat{\mathbf{y}}_{n+1})) \\ &\leq \frac{\eta}{2} \sum_{t=1}^n \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)}^2 + \frac{1}{\eta} \sup_{f \in \mathcal{F}} \Delta_R(f | \hat{\mathbf{y}}_1) \end{aligned}$$

□

3 Example: Multi-armed Bandit With Exponential Weights

For this example we will use the exponential weights algorithm as the full information algorithm. Recall that

$$R(f) = \sum_{i=1}^N f[i] \log(f[i])$$

and note that

$$\mathbf{y}'_{t+1}[j] = \hat{\mathbf{y}}_t[j] \times \exp(-\eta \nabla_t[j])$$

For a given $\hat{\mathbf{y}}_t$ we will simply draw $i_t \sim \hat{\mathbf{y}}_t$ and play $\hat{y}_t = e_{i_t}$ (that is expert i_t) and note that clearly $\mathbb{E}[\hat{y}_t] = \hat{\mathbf{y}}_t$. Further, for the unbiased estimate of loss, we use:

$$\tilde{\nabla}_t = \frac{e_{i_t}^\top \nabla_t}{q_t(i_t)} e_{i_t}$$

So that: $\mathbb{E}[\tilde{\nabla}_t] = \sum_{i=1}^d \hat{\mathbf{y}}_t(i) \frac{\nabla_t[i]}{\hat{\mathbf{y}}_t(i)} e_i = \nabla_t$

Now note that:

$$\nabla^2 R(f) = \begin{bmatrix} \frac{1}{f[1]} & 0 & 0 & 0 \\ 0 & \frac{1}{f[2]} & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \frac{1}{f[N]} \end{bmatrix}$$

Hence note that

$$\nabla^2 R(f)^{-1} = \begin{bmatrix} f[1] & 0 & 0 & 0 \\ 0 & f[2] & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & f[N] \end{bmatrix}$$

Hence using this in the local norm lemma we have the bound:

$$\begin{aligned} n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) &\leq \frac{\eta}{2} \sum_{t=1}^n \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)^{-1}}^2 + \frac{1}{\eta} \sup_{f \in \mathcal{F}} \Delta_R(f | \hat{y}_1) \\ &\leq \frac{\eta}{2} \sum_{t=1}^n \sum_{i=1}^N \tilde{\nabla}_t^2[i] z_t[i] + \frac{1}{\eta} \log(N) \end{aligned}$$

Now recall that z_t is of the form $z_t = \alpha_t \hat{\mathbf{y}}_t + (1 - \alpha_t) \hat{\mathbf{y}}'_{t+1}$ for some $\alpha_t \in [0, 1]$ and so

$$z_t[j] = \alpha_t \hat{\mathbf{y}}_t[j] + (1 - \alpha_t) \hat{\mathbf{y}}_t[j] \times \exp(-\eta \nabla_t[j]) \leq \hat{\mathbf{y}}_t[j] (1 + \exp(\eta))$$

Hence we conclude that

$$n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \leq \eta \sum_{t=1}^n \sum_{i=1}^N \tilde{\nabla}_t^2[i] \hat{\mathbf{y}}_t[i] + \frac{1}{\eta} \log(N)$$

Now plugging in the form of $\tilde{\nabla}_t = \frac{e_{i_t}^\top \nabla_t}{\hat{\mathbf{y}}_t(i_t)} e_{i_t}$

$$n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \leq \eta \sum_{t=1}^n \frac{\nabla_t^2[i_t]}{\hat{\mathbf{y}}_t[i_t]} + \frac{1}{\eta} \log(N)$$

Now using this regret bound for full information algorithm within the bandit bounds as before from Eq. 1 we conclude that:

$$\begin{aligned}
n\mathbb{E}[\text{Reg}_n] &\leq \mathbb{E}\left[n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n)\right] \\
&\leq \mathbb{E}\left[\eta \sum_{t=1}^n \frac{\nabla_t^2[i_t]}{\hat{\mathbf{y}}_t[i_t]}\right] + \frac{1}{\eta} \log(N) \\
&= \mathbb{E}\left[\eta \sum_{t=1}^n \mathbb{E}_{i_t \sim \hat{\mathbf{y}}_t} \left[\frac{\nabla_t^2[i_t]}{\hat{\mathbf{y}}_t[i_t]}\right]\right] + \frac{1}{\eta} \log(N) \\
&= \mathbb{E}\left[\eta \sum_{t=1}^n \sum_{i=1}^N \hat{\mathbf{y}}_t[i] \frac{\nabla_t^2[i]}{\hat{\mathbf{y}}_t[i]}\right] + \frac{1}{\eta} \log(N) \\
&= \eta \sum_{t=1}^n \|\nabla_t\|_2^2 + \frac{1}{\eta} \log(N) \\
&\leq \eta nN + \frac{1}{\eta} \log(N)
\end{aligned}$$

Setting $\eta = \sqrt{\log(N)/nN}$ we get,

$$\mathbb{E}[\text{Reg}_n] \leq 2\sqrt{\frac{N \log(N)}{n}}$$

4 Example: Multi-armed Bandit With Log Barrier

For this example we will use R to be the log barrier function,

$$R(f) = -\sum_{i=1}^N \log(f[i])$$

and note that

$$\mathbf{y}_{t+1}^{\hat{}}[j] = \frac{\hat{\mathbf{y}}_t[j]}{1 + \eta \hat{\mathbf{y}}_t[j] \tilde{\nabla}_t[j]}$$

This is indeed a strictly convex function. If we perform Mirror Descent using this function R . Just as before, for a given $\hat{\mathbf{y}}_t$ we will simply draw $i_t \sim \hat{\mathbf{y}}_t$ and play $\hat{y}_t = e_{i_t}$ (that is expert i_t) and note that clearly $\mathbb{E}[\hat{y}_t] = \hat{\mathbf{y}}_t$. Further, for the unbiased estimate of loss, we use:

$$\tilde{\nabla}_t = \frac{e_{i_t}^\top \nabla_t}{q_t(i_t)} e_{i_t}$$

So that: $\mathbb{E}[\tilde{\nabla}_t] = \sum_{i=1}^d \hat{\mathbf{y}}_t(i) \frac{\nabla_t[i]}{\hat{\mathbf{y}}_t(i)} e_i = \nabla_t$ Now note that for this R ,

$$\nabla^2 R(f) = \begin{bmatrix} \frac{1}{f^2[1]} & 0 & 0 & 0 \\ 0 & \frac{1}{f^2[2]} & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \frac{1}{f^2[N]} \end{bmatrix}$$

Hence note that

$$\nabla^2 R(f)^{-1} = \begin{bmatrix} f^2[1] & 0 & 0 & 0 \\ 0 & f^2[2] & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & f^2[N] \end{bmatrix}$$

Hence using this in the local norm lemma we have the bound:

$$\begin{aligned} n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) &\leq \frac{\eta}{2} \sum_{t=1}^n \|\tilde{\nabla}_t\|_{\nabla^2 R(z_t)^{-1}}^2 + \frac{1}{\eta} \sup_{f \in \mathcal{F}} \Delta_R(f|\hat{y}_1) \\ &\leq \frac{\eta}{2} \sum_{t=1}^n \sum_{i=1}^N \tilde{\nabla}_t^2[i] z_t^2[i] + \frac{1}{\eta} N \log(N) \end{aligned}$$

Now recall that z_t is of the form $z_t = \alpha_t \hat{y}_t + (1 - \alpha_t) \hat{y}'_{t+1}$ for some $\alpha_t \in [0, 1]$ and so

$$z_t[j]^2 \leq 2\hat{y}_t[j]^2$$

Hence we conclude that

$$n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \leq \eta \sum_{t=1}^n \sum_{i=1}^N \tilde{\nabla}_t^2[i] \hat{y}_t^2[i] + \frac{1}{\eta} N \log(N)$$

Now plugging in the form of $\tilde{\nabla}_t = \frac{e_{i_t}^\top \nabla_t}{\hat{y}_t(i_t)} e_{i_t}$

$$n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \leq \eta \sum_{t=1}^n \nabla_t^2[i_t] + \frac{1}{\eta} N \log(N)$$

Now using this regret bound for full information algorithm within the bandit bounds as before from Eq. 1 we conclude that:

$$\begin{aligned} n\mathbb{E}[\text{Reg}_n] &\leq \mathbb{E} \left[n\text{Reg}_n(\tilde{\nabla}_1, \dots, \tilde{\nabla}_n) \right] \\ &\leq \mathbb{E} \left[\eta \sum_{t=1}^n \nabla_t^2[i_t] \right] + \frac{1}{\eta} N \log(N) \\ &= \mathbb{E} \left[\eta \sum_{t=1}^n \mathbb{E}_{i_t \sim \hat{y}_t} [\nabla_t^2[i_t]] \right] + \frac{1}{\eta} N \log(N) \\ &= \mathbb{E} \left[\eta \sum_{t=1}^n \sum_{i=1}^N \hat{y}_t[i] \nabla_t^2[i] \right] + \frac{1}{\eta} N \log(N) \\ &= \eta \sum_{t=1}^n \mathbb{E} [\nabla_t^2[i_t]] + \frac{1}{\eta} N \log(N) \end{aligned}$$

If $\nabla_t[i] \in [0, 1]$ then note that,

$$\sum_{t=1}^n \mathbb{E} [\nabla_t[i_t]] - \min_{j \in [N]} \sum_{t=1}^n \nabla_t[j] \leq \eta \sum_{t=1}^n \mathbb{E} [\nabla_t[i_t]] + \frac{1}{\eta} N \log(N)$$

Hence we conclude that,

$$\sum_{t=1}^n \mathbb{E} [\nabla_t[i_t]] - \min_{j \in [N]} \sum_{t=1}^n \nabla_t[j] \leq \frac{\eta}{1-\eta} \min_{j \in [N]} \sum_{t=1}^n \nabla_t[j] + \frac{1}{\eta(1-\eta)} N \log(N)$$

Note that if we knew $\min_{j \in [N]} \sum_{t=1}^n \nabla_t[j]$ we could set

$$\eta = \min \left\{ 1/2, \sqrt{N \log(N)} \sqrt{\min_{j \in [N]} \sum_{t=1}^n \nabla_t[j]} \right\}$$

to get that

$$\mathbb{E} [\text{Reg}_n] \leq O \left(\frac{\sqrt{N \log(N)} \min_{j \in [N]} \sum_{t=1}^n \nabla_t[j]}{n} + \frac{N \log N}{n} \right)$$

This bound is nice since it says that if the best expert has low cumulative loss, say of constant order, then we in fact get a $N \log(N)/n$ rate.