

# EXPBYOPT FOR CONVEX BANDITS

June 3, 2025

## 1 Setup and Notation

Reuse the notation from [Lattimore \[2024\]](#).

**Assumption 1.** The following hold:

- The losses are in  $\mathcal{F}_b$ .
- There is no noise so that  $Y_t = f_t(X_t)$ .

**Assumption 2.**  $\mathcal{C}$  is finite subset of  $K$  such that:

- $\log(\mathcal{C}) \leq \tilde{\mathcal{O}}(d)$ .
- For all  $f \in \mathcal{F}_b$  there exists  $x \in \mathcal{C}$  such that  $f(x) \leq \inf_{x' \in K} f(x') + \frac{1}{n}$ .

## 2 Exponential Weights with Importance Sampling

Let  $(\hat{s}_t)_{t=1}^n : \mathcal{C} \rightarrow \mathbb{R}$  be a sequence of functions and

$$q_t(x) = \frac{\exp(-\eta \sum_{u=1}^{t-1} \hat{s}_u(x))}{\sum_{y \in \mathcal{C}} \exp(-\eta \sum_{u=1}^{t-1} \hat{s}_u(y))}, \quad x \in \mathcal{C}. \quad (1)$$

The following theorem gives a bound on the regret of the exponential weights algorithm in the original setting. *Here, it's better to think of  $\hat{s}_t$  as the losses observed at time  $t$  and not the estimates of the losses.*

**Theorem 3** ([Lattimore \[2024\]](#), Theorem 8.11). *For any  $y \in \mathcal{C}$  we have*

$$\sum_{t=1}^n \langle q_t, \hat{s}_t \rangle - \hat{s}_t(y) \leq \frac{\log(|\mathcal{C}|)}{\eta} + \frac{1}{\eta} \sum_{t=1}^n \mathcal{S}_t(\eta \hat{s}_t),$$

where  $\mathcal{S}_t(u) = D_{R^*}(R'(q_t) - u, R'(q_t))$ .

The following theorem applies to the setting when  $\langle f_t, X_t \rangle$  is observed and not  $f_t$ . The idea is to use importance sampling and estimates  $\hat{s}_t$  to compute the losses.

```

1  args: learning rate  $\eta > 0$ 
2  let  $\mathcal{C} \in K$  be finite
3  for  $t = 1$  to  $n$ :
4      compute  $q_t(x) = \frac{\exp(-\eta \sum_{u=1}^{t-1} \hat{s}_u(x))}{\sum_{y \in \mathcal{C}} \exp(-\eta \sum_{u=1}^{t-1} \hat{s}_u(x))}$  for all  $x \in \mathcal{C}$ 
5      find distribution  $p_t$  as a function of  $q_t$ 
6      sample  $X_t \sim p_t$ , and observe  $Y_t = f_t(X_t)$ 
7      compute  $\hat{s}_t(x) \forall x \in \mathcal{C}$  using  $p_t, q_t, X_t, Y_t$ 

```

**Algorithm 1:** Exponential Weights with Importance Sampling

**Theorem 4** (Lattimore [2024], Theorem 8.14). *Let  $x_\star = \arg \min_{x \in \mathcal{C}} \sum_{t=1}^n f_t(x)$  and  $p_\star \in \Delta(\mathcal{C})$  be a Dirac on  $x_\star$ . The expected regret of Algorithm 1 is bounded by*

$$\mathbb{E}[\mathfrak{R}_n(x_\star)] = \frac{\log(|\mathcal{C}|)}{\eta} + \sum_{t=1}^n \mathbb{E} \left[ \langle p_t - p_\star, f_t \rangle + \langle p_\star - q_t, \hat{s}_t \rangle + \frac{1}{\eta} \mathcal{S}_t(\eta \hat{s}_t) \right].$$

*Proof.* Immediate from Theorem 3. □

### 3 Exploration By Optimisation

Let

- $\mathcal{G}$  be the set of all functions  $g : \mathcal{C} \rightarrow \mathbb{R}$ , i.e.,  $\mathcal{G} = \mathbb{R}^{|\mathcal{C}|}$ .
- $\mathcal{E}$  be the set of all functions  $E : \mathcal{C} \times \mathbb{R} \rightarrow \mathcal{G}$ .

The idea is to bound the term inside the expectation in Theorem 4 uniformly over all  $t \in [n]$ . To this end, define

$$\Lambda_\eta(q, p, E, r, f) = \frac{1}{\eta} \mathbb{E} \left[ \langle p - r, f \rangle + \langle r - q, E(X, Y) \rangle + \frac{1}{\eta} \mathcal{S}_q(\eta E(X, Y)) \right].$$

Note that the randomness is only through  $X$  and  $Y$ , where  $X \sim p$  and  $Y = f(X)$ . For the learner, the degree of freedom is in the choice of the estimator  $\hat{s}_t$  and the exploration distribution  $p_t$ . Therefore, fix  $t \in [n]$  and define

$$\Lambda_\eta(q) = \inf_{p \in \Delta(\mathcal{C}), E \in \mathcal{E}} \sup_{r \in \Delta(\mathcal{C}), f \in \mathcal{F}_b} \frac{1}{\eta} \Lambda_\eta(q, p, E, r, f).$$

**Remark 5.** There is a bit of waste here since the supremum is taken over all  $r \in \Delta(\mathcal{C})$ , but it could've been restricted to dirac distributions.

It's worth noting that Algorithm 2 is a instantiation of Algorithm 1 where the estimator  $\hat{s}_t$  and the exploration distribution  $p_t$  are chosen by solving an optimisation problem.

```

1 args: learning rate  $\eta > 0$ , precision  $\epsilon > 0$ 
2 let  $\mathcal{C} \in K$  be finite
3 for  $t = 1$  to  $n$ :
4     compute  $q_t(x) = \frac{\exp(-\eta \sum_{u=1}^{t-1} \hat{s}_u(x))}{\sum_{y \in \mathcal{C}} \exp(-\eta \sum_{u=1}^{t-1} \hat{s}_u(x))}$  for all  $x \in \mathcal{C}$ 
5     find distribution  $p_t \in \Delta(\mathcal{C})$  and  $E_t \in \mathcal{E}$  such that
6          $\Lambda_\eta(q_t, p_t, E_t) \leq \inf_{p, E} \Lambda_\eta(q_t) + \epsilon$ 
7     sample  $X_t \sim p_t$ , and observe  $Y_t = f_t(X_t)$ 
8     compute  $\hat{s}_t(x) = E_t(X_t, Y_t)$ 

```

**Algorithm 2:** Exploration by Optimisation

**Theorem 6** (Lattimore [2024], Theorem 8.15). *The expected regret of Algorithm 2 is bounded by*

$$\mathbb{E}[\mathfrak{R}_n(x_\star)] \leq \frac{\log(|\mathcal{C}|)}{\eta} + n\eta \sup_{q \in \Delta(\mathcal{C})} \Lambda_\eta(q) + n\eta\epsilon.$$

Therefore, from Theorem 6, to bound the regret of Algorithm 2 it suffices to bound

$$\Lambda_{\mathcal{C}}^\star = \sup_{\eta > 0, q \in \Delta(\mathcal{C})} \inf_{p \in \Delta(\mathcal{C}), E \in \mathcal{E}} \sup_{y \in \mathcal{C}, f \in \mathcal{F}_b} \Lambda_\eta(q, p, E, y, f).$$

The following lemma shows that the order of the infimum and supremum can be interchanged.

**Lemma 7.** *For all  $\eta > 0$  and  $q \in \Delta(\mathcal{C})$ , we have*

$$\inf_{p \in \Delta(\mathcal{C}), E \in \mathcal{E}} \sup_{r \in \Delta(\mathcal{C}), f \in \mathcal{F}_b} \Lambda_\eta(q, p, E, r, f) = \sup_{r \in \Delta(\mathcal{C}), f \in \mathcal{F}_b} \inf_{p \in \Delta(\mathcal{C}), E \in \mathcal{E}} \Lambda_\eta(q, p, E, r, f). \quad (2)$$

### 3.1 Unrestricted estimator class $\mathcal{E}$

The minimizer  $E$  in the general case has the following close form.

**Lemma 8.** *Given  $r \in \Delta(\mathcal{C})$ ,  $f \in \mathcal{F}_b$ , and  $p \in \Delta(\mathcal{C})$  define  $G_{r,f,p} \in \mathcal{E}$  as*

$$G_{r,f,p}(x, y) = \frac{1}{\eta} (R'(q) - R'(\mathbb{E}[p_\star | f(x) = y])), \quad (3)$$

*then we have*

$$\inf_{E \in \mathcal{E}} \Lambda_\eta(q, p, E, r, f) = \Lambda_\eta(q, p, G_{r,f,p}, r, f).$$

### 3.2 Restricted estimator class $\mathcal{E}$

Define  $E$  through a probability kernel  $T$  such that

$$E(x, y) = \frac{T(x|\cdot)y}{p(x)} \quad (= \hat{s}_x(\cdot)),$$

$$\text{and} \quad \sum_{x \in \mathcal{C}} T(x|y) = 1 \quad \text{for all } y \in \mathcal{C},$$

Again, fix  $r \in \Delta(\mathcal{C})$ ,  $f \in \mathcal{F}_b$ , and  $p \in \Delta(\mathcal{C})$ . We are interested in the  $T$  that minimizes

$$\begin{aligned} \Lambda(T) &:= \mathbb{E} \left[ \langle p - r, f \rangle + \langle r - q, E(X, Y) \rangle + \frac{1}{\eta} \mathcal{S}_q(\eta E(X, Y)) \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \langle p - r, f \rangle + \langle r - q, \frac{T(X|\cdot)Y}{p(X)} \rangle + \frac{1}{\eta} \mathcal{S}_q \left( \eta \frac{T(X|\cdot)Y}{p(X)} \right) \middle| X \right] \right]. \end{aligned}$$

Define the shorthand  $T_x = T(X|\cdot) \in \mathbb{R}^{|\mathcal{C}|}$ . Then through differentiation we have

$$\begin{aligned} \nabla_{T_x} \Lambda(T) &= \nabla_{T_x} \mathbb{E} \left[ \langle p - r, f \rangle + \langle r - q, \frac{T(X|\cdot)Y}{p(X)} \rangle + \frac{1}{\eta} \mathcal{S}_q \left( \eta \frac{T(X|\cdot)Y}{p(X)} \right) \middle| X = x \right] \\ &= \nabla_{T_x} \mathbb{E} \left[ \langle r - q, \frac{T_x Y}{p(X)} \rangle + \frac{1}{\eta} \left( R^* \left( R'(q) - \eta \frac{T_x Y}{p(X)} \right) - R^*(R'(q)) - \nabla R^*(R'(q))^\top (-\eta \frac{T_x Y}{p(X)}) \right) \middle| X = x \right] \\ &= \nabla_{T_x} \mathbb{E} \left[ \langle r, \frac{T_x Y}{p(X)} \rangle - \frac{Y}{p(X)} q^\top T_x + \frac{1}{\eta} \left( R^* \left( R'(q) - \eta \frac{T_x Y}{p(X)} \right) + \frac{\eta Y}{p(X)} q^\top T_x \right) \middle| X = x \right] \\ &= \nabla_{T_x} \mathbb{E} \left[ \langle r, \frac{T_x Y}{p(X)} \rangle + \frac{1}{\eta} R^* \left( R'(q) - \eta \frac{T_x Y}{p(X)} \right) \middle| X = x \right] \\ &= \mathbb{E} \left[ \frac{Y}{p(x)} \left( r - \nabla R^* \left( R'(q) - \eta \frac{T_x Y}{p(x)} \right) \right) \middle| X = x \right] \\ &= \frac{1}{p(x)} \mathbb{E} \left[ Y \left( r - \exp \left( R'(q) - \eta \frac{T_x Y}{p(x)} \right) \right) \middle| X = x \right] \\ &= \frac{1}{p(x)} \left( f(x)r - \mathbb{E} \left[ Y \exp \left( R'(q) - \frac{\eta}{p(x)} Y T_x \right) \middle| X = x \right] \right) \end{aligned}$$

we want to solve for  $T_x$  such that the gradient is zero. Therefore, for coordinate  $z \in \mathcal{C}$  we have

$$\mathbb{E} \left[ Y \exp \left( R'(q)_z - \frac{\eta}{p(x)} Y T_{x,z} \right) \middle| X = x \right] = f(x)r_z,$$

which is equivalent to

$$\begin{aligned} \sum_{x' \in \mathcal{C}} p(x') f(x') \exp\left(R'(q)_z - \frac{\eta}{p(x)} f(x') T_{x,z}\right) &= f(x) r_z \\ \Leftrightarrow \sum_{x' \in \mathcal{C}} p(x') f(x') \exp\left(-\frac{\eta}{p(x)} f(x') T_{x,z}\right) &= \exp(-R'(q)_z) f(x) r_z \end{aligned}$$

## References

Tor Lattimore. Bandit convex optimisation. *arXiv preprint arXiv:2402.06535*, 2024.