

Online Linear Optimization

— Full feedback and Bandit feedback

Source: Bandit Algorithms, L. & S. Chap 28.

Setting: Evolves over time steps $t=1, 2, \dots, n$.

At each time:

Player: Chooses action $a_t \in A \subseteq \mathbb{R}^d$ (A a convex set).

Adversary: Chooses loss $y_t \in \mathcal{L} \subseteq \mathbb{R}^d$. y_t can depend on the player policy (same as the usual adversarial setting).

Observation: (Full feedback Setting, for now).

Loss incurred at time $t = \langle a_t, y_t \rangle$

Player also gets to observe y_t

Regret (w.r.t. best fixed action in hindsight):

$$R_n = \max_{a \in A} \sum_{t=1}^n \langle a_t - a, y_t \rangle.$$

Note: Action is not randomized in the full feedback setting.

Two Related Algorithms: Mirror Descent and Follow the Regularized Leader (FTRL).

Mirror Descent: Parameters: n (time horizon), \mathcal{A} (set of actions, $\mathcal{A} \subseteq \mathbb{R}^d$), η (learning rate), $F: \mathbb{R}^d \rightarrow \mathbb{R}$, with $\text{domain}(F) = \mathcal{D}$ (Legendre function).

Potential function

Initial: $a_1 = \arg \min_{a \in \mathcal{A}} F(a)$.

$t \geq 1$: $a_{t+1} = \arg \min_{a \in \mathcal{A}} \left(\eta \langle a, y_t \rangle + D_F(a, a_t) \right)$

typically, \mathcal{A} is a convex set

Bregman Divergence

FTRL (Follow the Regularized Leader): (parameters as above).

Initial: $a_1 = \arg \min_{a \in \mathcal{A}} F(a)$.

$$t \geq 1: \quad a_{t+1} = \operatorname{argmin}_{a \in A} \left(\eta \sum_{s=1}^t \langle a, y_s \rangle + F(a) \right)$$

Note 1: $g(x) = b^T x$ (linear function)
 $\Rightarrow D_g(x, y) = 0$

FTRL
implementation

$$\therefore F \text{ Legendre}, g \text{ linear} \Rightarrow D_{F+g}(x, y) = D_F(x, y).$$

Thus, for A convex, $D = \mathbb{R}^d$, FTRL
is equivalent to:

more carefully, $\tilde{a}_{t+1} = \operatorname{argmin}_{a \in A} \left(\eta \sum_{s=1}^t \langle a, y_s \rangle + F(a) \right)$.

if $D = \mathbb{R}^d$, $a_{t+1} = \operatorname{argmin}_{a \in A} D_F(a, \tilde{a}_{t+1})$
 this reduces to the unconstrained problem.

Simplifying: (taking grad. and setting to zero).

$$\tilde{a}_{t+1} = \operatorname{argmin}_a \left(\langle a, \eta \sum_{s=1}^t y_s \rangle + F(a) \right).$$

$$\therefore \eta \sum_{s=1}^t y_s + \nabla F(\tilde{a}_{t+1}) = 0$$

$$\Rightarrow \tilde{a}_{t+1} = \nabla F^{-1} \left(-\eta \sum_{s=1}^t y_s \right)$$

$$= \nabla F^* \left(-\eta \sum_{s=1}^t y_s \right).$$

$a_{t+1} = \Pi_{F, D} \left(\nabla F^* \left(-\eta \sum_{s=1}^t y_s \right) \right)$

FTRL

Bregman Projection

Note 2: $\tilde{a}_{t+1} = \underset{a}{\operatorname{argmin}} \left(\eta \langle a, y_t \rangle + D_F(a, a_t) \right)$

more carefully, $a \in D$

Mirror Descent
Implementation

$$a_{t+1} = \underset{a \in A}{\operatorname{argmin}} D_F(a, \tilde{a}_{t+1})$$

Simplify: (Take gradient, and set to zero)

$$\eta y_t + \nabla F(\tilde{a}_{t+1}) - \nabla F(a_t) = 0$$

$$\Rightarrow \tilde{a}_{t+1} = \nabla F^* \left(\nabla F(a_t) - \eta y_t \right).$$

$a_{t+1} = \Pi_{F, D} \left(\nabla F^* \left(\nabla F(a_t) - \eta y_t \right) \right)$

Note 3: Follow the Leader (FTL):

FTL is
bad.

$$a_{t+1} = \underset{a \in A}{\operatorname{argmin}} \sum_{s=1}^t \langle a, y_s \rangle$$

This is a bad algorithm, because a_t can vary too much over time steps chasing the best action-to-date in hindsight, and can result in large regret. To see this,

$$A = [-1, 1], \quad y_1 = \frac{1}{2}, \quad y_2 = y_4 = y_6 = \dots = -1$$

$$a_1 = 0 \text{ (arbitrary)} \quad y_3 = y_5 = y_7 = \dots = +1$$

$$a_2 = \underset{a \in [-1, 1]}{\operatorname{argmin}} \langle a, y_1 \rangle = -1. \quad \begin{array}{l} \langle a, y_1 \rangle = \frac{a}{2} \\ a = -1 \end{array}$$

$$\langle a, y_1 \rangle + \langle a, y_2 \rangle = \left(\frac{a}{2} - a \right) = -\frac{a}{2} \quad a = (+1)$$

$$a_3 = \underset{a \in [-1, 1]}{\operatorname{argmin}} \left(\frac{a}{2} - a \right) = +1$$

$$\langle a, y_1 \rangle + \langle a, y_2 \rangle + \langle a, y_3 \rangle = \frac{a}{2} - a + a = \frac{a}{2} \quad a_4 = -1$$

$$a_4 = \underset{a \in [-1, 1]}{\operatorname{argmin}} \left(\frac{a}{2} - a + a \right) = -1$$

$$\text{i.e., } \{a_2, a_3, a_4, \dots\} = \{-1, +1, -1, +1, \dots\}$$

$$\begin{aligned} \text{Loss} &= \langle a_2, y_2 \rangle + \langle a_3, y_3 \rangle + \dots = \underbrace{1 + 1 + \dots + 1}_{(n-1)} \\ &= (n-1) \quad (\text{at time } n) \end{aligned}$$

for the best fixed action! Choose $a = 1 \forall t \geq 1$.

$$\text{Loss} = \frac{1}{2} - 1 + 1 - 1 + 1 - \dots = \begin{cases} \frac{1}{2} & \text{if } n \text{ odd} \\ -\frac{1}{2} & \text{if } n \text{ even} \end{cases}$$

\Rightarrow Linear regret w.r.t best arm in hindsight!

Adding a regularizer prevents actions from changing too often, and thus, FTRL has much better regret scaling.

See also: Online Learning and Online Convex Optimization, by Shai Shalev-Shwartz, Found. and Trends in ML, Vol. 4, No. 2, 2011

Note 4: FTRL and Mirror Descent are related, and equivalent when $D \subseteq A$.

(Note this is not a typical setting; typically it is the other way around where $D \supseteq A$)

In this case, we have $\tilde{a}_{t+1} = a_{t+1}$, because the second step in the implementation (the projection onto Δ using Bregman divergence) is unnecessary.

In the calculation below, we are further assuming $\Delta = \mathbb{R}^d$.

$$\text{FTRL: } \tilde{a}_{t+1} = \nabla F^*(-\eta \sum_{s=1}^t y_s).$$

$$\text{Mirror Descent: } \tilde{a}_{t+1} = \nabla F^*(\nabla F(a_t) - \eta y_t).$$

$$\text{But } \nabla F(a_t) = \nabla F(\tilde{a}_t)$$

$$= \nabla F(\underbrace{\nabla F^*(\nabla F(a_{t-1}) - \eta y_{t-1})}_{\substack{= \nabla F^{-1} \\ \text{inverses.}}})$$

$$= \nabla F(a_{t-1}) - \eta y_{t-1}$$

$$\Rightarrow \tilde{a}_{t+1} = \nabla F^*(\nabla F(a_{t-1}) - \eta y_{t-1} - \eta y_t)$$

\vdots

$$= \nabla F^*(-\eta \sum_{s=1}^t y_s)$$

$\left(\begin{array}{l} \because \nabla F(a_1) = 0 \\ \text{because by defn.,} \\ a_1 = \underset{a}{\operatorname{argmin}} F(a) \end{array} \right)$

Note 5: Exponential weights algorithm.

Suppose $A = \mathbb{P}^{d-1}$, the d -dimensional simplex of probability distributions, and $F(x) = \sum_{i=1}^d (x_i \ln x_i - x_i)$.
 discrete pmf: $\sum_i x_i = 1, x_i \geq 0$.

Here $D_F(x, y) = \sum_{i=1}^d x_i \ln \left(\frac{x_i}{y_i} \right)$ (KL Divergence)

Mirror Descent: extends D_F to \mathbb{R}_+^d $+ \sum_{i=1}^d (y_i - x_i)$

$$\tilde{\alpha}_{t+1} = \underset{\alpha}{\operatorname{argmin}} \left(\eta \langle \alpha, y_t \rangle + D_F(\alpha, \alpha_t) \right)$$

Take derivative and set to 0:

$$0 = \eta y_t + \nabla F(\tilde{\alpha}_{t+1}) - \nabla F(\alpha_t) \rightarrow \textcircled{*}$$

$$\nabla F(x)_i = \frac{\partial F(x)}{\partial x_i} = \ln x_i \Rightarrow \nabla F(x) = \ln(x)$$

Observe $\textcircled{*}$ decouples across coordinates, i.e.,

$$\ln(\tilde{\alpha}_{t+1,i}) = \ln(\alpha_{t,i}) - \eta y_{t,i}$$

$$\Rightarrow \tilde{\alpha}_{t+1,i} = e^{-\eta y_{t,i}} \cdot \alpha_{t,i}, i=1, 2, \dots, d.$$

Project: $a_{t+1} = \underset{a \in A}{\operatorname{argmin}} D_F(a, \tilde{a}_{t+1})$

$$= \underset{a \in A}{\operatorname{argmin}} \sum_{i=1}^d a_i \ln \left(\frac{a_i}{a_{t,i} e^{-\eta y_{t,i}}} \right)$$

constant wrt $a \in A$ ↑
 $\sum_{i=1}^d a_i$ ↑
 $\sum_{i=1}^d a_i = 1$ ↓
 $a_i \in \mathbb{R}^{d-1}$ ↑

Claim: $a_i = \frac{a_{t,i} e^{-\eta y_{t,i}}}{\sum_j a_{t,j} e^{-\eta y_{t,j}}}$ is the minimizer to above. ↓ ∴ this can be ignored.

Pf: Let $C = \sum_{j=1}^d a_{t,j} e^{-\eta y_{t,j}}$

Then observe that $\sum_{i=1}^d a_i \ln \left(\frac{a_i}{a_{t,i} e^{-\eta y_{t,i}}} \right)$

$$= \sum_{i=1}^d a_i \ln \left(a_i \cdot \frac{C}{a_{t,i} e^{-\eta y_{t,i}}} \cdot \frac{1}{C} \right)$$

$$= \sum_{i=1}^d a_i \ln \left(\frac{a_i}{\left(\frac{a_{t,i} e^{-\eta y_{t,i}}}{C} \right)} \right) + \left(\sum_{i=1}^d a_i \right) \ln \left(\frac{1}{C} \right)$$

$$= \sum_{i=1}^d a_i \ln \left(a_i / \left(\frac{a_{t,i} e^{-\eta y_{t,i}}}{C} \right) \right) + \ln \left(\frac{1}{C} \right).$$

≥ 0 for any $a_i \in \mathbb{R}^{d-1}$

with equality when $a_i = \left(\frac{a_{t,i} e^{-\eta y_{t,i}}}{c} \right)$ PQ

$$\therefore \begin{cases} a_{t+1,i} = \frac{a_{t,i} e^{-\eta y_{t,i}}}{\sum_j a_{t,j} e^{-\eta y_{t,j}}} & , i=1, \dots, d \\ \end{cases}$$

results from Mirror Descent + KL projection on
Simplex = Exp 3 Rule PQ

Note 6: Equivalently, we have for FTRL:

FTRL $\tilde{a}_{t+1} = \underset{a}{\operatorname{argmin}} \left(\eta \sum_{s=1}^t \langle a, y_s \rangle + F(a) \right).$

$$a_{t+1} = \underset{a \in \Delta}{\operatorname{argmin}} D_F(a, \tilde{a}_{t+1})$$

$$\Delta = \mathbb{R}^{d-1}, \quad F(x) = \sum_{i=1}^d (x_i \ln x_i - x_i)$$

Then, similar calculation as above shows:

$$a_{t+1,i} = \frac{e^{-\eta \sum_{s=1}^t y_{s,i}}}{\sum_{j=1}^d e^{-\eta \sum_{s=1}^t y_{s,j}}}, \quad i=1, 2, \dots, d.$$

i.e., we get back Exp3 algorithm



Regret Bound

Mirror Descent:

Thm (28.4 in textbook): $\eta > 0$, and F Legendre, with domain D . A a non-empty convex set, with interior $(D) \cap A \neq \emptyset$. \rightarrow (assumed to exist).

Mirror descent actions: $\{a_1, a_2, \dots, a_{n+1}\}$ with

$$a_1 = \arg \min_{a \in A} F(a), \quad \tilde{a}_{t+1} = \arg \min_{a \in D} \eta \langle a, y_t \rangle + D_F(a, a_t)$$

$$a_{t+1} = \arg \min_{a \in A} D_F(a, \tilde{a}_{t+1}).$$

For any $a \in A$,

$$R_n(a) \leq \frac{F(a) - F(a_1)}{\eta} + \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle - \frac{1}{\eta} \sum_{t=1}^n D_F(a_{t+1}, a_t)$$

Further, suppose $\nabla F(a) - \eta y \in \text{interior}(\text{Dom}(F^*)) \neq y \in \mathcal{D}$,
 $a \in A \cap \mathcal{D}$.

Then,

$$R_n(a) \leq \frac{1}{\eta} \left(F(a) - F(a_1) + \sum_{t=1}^n D_F(a_t, \tilde{a}_{t+1}) \right).$$

Prof: Suppose $a \in A$, but $a \notin \mathcal{D}$. Then $F(a) = \infty$, thus
result holds. Thus, assume $a \in A \cap \mathcal{D}$.

Now, by defn, $R_n(a) = \sum_{t=1}^n \langle a_t - a, y_t \rangle$.

Further,

$$\langle a_t - a, y_t \rangle = \langle a_t - a_{t+1}, y_t \rangle + \langle a_{t+1} - a, y_t \rangle \rightarrow ①$$

Recall: $a_{t+1} = \underset{b \in A}{\operatorname{argmin}} \eta \langle b, y_t \rangle + D_F(b, a_t)$.

Further, $D_F(b, a_t) = F(b) - F(a_t) - \langle \nabla F(a_t), b - a_t \rangle$.

Now, $D_F(b, a_t)$ convex in b , and $\eta \langle b, y_t \rangle$ linear in b . Thus, $\eta \langle b, y_t \rangle + D_F(b, a_t)$ convex. From first order optimality, it follows that:

$$\left\langle \nabla \left(\eta \langle b, y_t \rangle + D_F(b, a_t) \right) \Big|_{b=a_{t+1}}, a - a_{t+1} \right\rangle \geq 0$$

for any $a \in A$.

$$\text{i.e., } \left\langle \eta y_t + \nabla F(a_{t+1}) - \nabla F(a_t), a - a_{t+1} \right\rangle \geq 0$$

Reordering this, we have:

$$\begin{aligned} \langle a_{t+1} - a, y_t \rangle &\leq \frac{1}{\eta} \langle a - a_{t+1}, \nabla F(a_{t+1}) - \nabla F(a_t) \rangle \\ &= \frac{1}{\eta} \left(D_F(a, a_t) - D_F(a, a_{t+1}) - D_F(a_{t+1}, a_t) \right) \\ &\quad \begin{matrix} \bullet & \downarrow & \bullet & \downarrow & \bullet & \downarrow & \bullet \\ F(a) - F(a_t) - & & F(a) - F(a_{t+1}) - & & F(a_{t+1}) - F(a_t) - & & \\ & & & & & & \\ & & \langle \nabla F(a_t), a - a_t \rangle & & \langle \nabla F(a_{t+1}), a - a_{t+1} \rangle & & \end{matrix} \end{aligned}$$

$\hookrightarrow \textcircled{2}$

$$\therefore R_n = \sum_{t=1}^n \langle a_t - a, y_t \rangle$$

$$\textcircled{1} = \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle + \sum_{t=1}^n \langle a_{t+1} - a, y_t \rangle$$

$$\stackrel{(2)}{\leq} \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle + \frac{1}{\eta} \sum_{t=1}^n \left(D_F(a, a_t) - D_F(a, a_{t+1}) - D_F(a_{t+1}, a_t) \right)$$

Collect/cancel terms

$$= \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle + \frac{1}{\eta} \left(\underbrace{D_F(a, a_1)}_{\leq F(a) - F(a_1)} - \underbrace{D_F(a, a_{n+1})}_{\geq 0} - \sum_{t=1}^n D_F(a_{t+1}, a_t) \right)$$

(3)

Now, $D_F(a, a_{n+1}) \geq 0$ and $D_F(a, a_1) = F(a) - F(a_1)$

$$- \underbrace{\langle \nabla F(a_1), a - a_1 \rangle}_{\geq 0}$$

However, $a_1 = \underset{a}{\operatorname{argmin}} F(a) \Rightarrow \langle \nabla F(a_1), a - a_1 \rangle \geq 0$
first order optimality.

∴ From (3) + above,

(4)

$$R_n \leq \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle + \frac{F(a) - F(a_1)}{\eta} - \frac{1}{\eta} \sum_{t=1}^n D_F(a_{t+1}, a_t)$$

Next, since $\nabla F(a) - \eta y \in \text{interior}(\text{Dom}(F^*))$, and

$$\tilde{a}_{t+1} = \underset{a \in D}{\operatorname{argmin}} \eta \langle a, y_t \rangle + D_F(a, a_t)$$

gradient at \tilde{a}_{t+1} is zero

$$\Rightarrow \eta y_t + \nabla F(\tilde{a}_{t+1}) - \nabla F(a_t) = 0$$

$$\Rightarrow \langle a_t - a_{t+1}, y_t \rangle = \frac{1}{\eta} \langle a_t - a_{t+1}, \nabla F(a_t) - \nabla F(\tilde{a}_{t+1}) \rangle$$

$$= \frac{1}{\eta} \left(D_F(a_{t+1}, a_t) + D_F(a_t, \tilde{a}_{t+1}) - D_F(a_{t+1}, \tilde{a}_{t+1}) \right).$$

$$\leq \frac{1}{\eta} \left(D_F(a_{t+1}, a_t) + D_F(a_t, \tilde{a}_{t+1}) \right).$$

→ ⑤

Substituting ⑤ in ④,

$$R_n(a) \leq \frac{1}{\eta} \left(F(a) - F(a_1) + \sum_{t=1}^n D_F(a_t, \tilde{a}_{t+1}) \right).$$

□

FTRL Regret Bound:

Thm (28.5 in textbook) : As before: $\eta > 0$, F Legendre, over D . A convex, $A \cap D \neq \emptyset$. $\{a_1, a_2, \dots, a_{n+1}\}$ actions from FTRL. Then, for any $a \in A$, we have:

$$R_n(a) \leq \underbrace{F(a) - F(a_1)}_{\gamma} + \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle - \frac{1}{\gamma} \sum_{t=1}^n D_F(a_{t+1}, a_t)$$

PF: Similar to above.



Applications

① $\mathcal{A} = \text{Unit ball in } \mathbb{R}^d = \{a \in \mathbb{R}^d : \|a\|_2 \leq 1\}$

$$\mathcal{X} = \{y_t \in \mathbb{R}^d : \|y_t\|_2 \leq 1\}$$

(i.e., both player and adversary choose vectors from the unit ℓ_2 -ball $\cong B_2^d$).

$$F(a) = \frac{1}{2} \|a\|_2^2 \quad \gamma = \frac{1}{\sqrt{n}}.$$

Prop (28.6 in book): Setting as above. Then,

$$R_n(a) \leq \sqrt{n}$$

Proof: $D_F(x, y) = \frac{1}{2} \|x - y\|^2$

Further, \tilde{a}_{t+1} s.t. $\gamma y_t + \nabla F(\tilde{a}_{t+1}) - \nabla F(a_t) = 0$

$$\Rightarrow \eta y_t + \tilde{a}_{t+1} - a_t = 0$$

$$\therefore D_F(a_t, \tilde{a}_{t+1}) = \frac{1}{2} \|a_t - \tilde{a}_{t+1}\|^2 = \frac{\eta^2}{2} \|y_t\|^2$$

$$\text{Also, } |F(a) - F(b)| \leq \frac{1}{2} \quad \forall a, b \in B_2^d$$

$$\therefore R_n(\tilde{a}) \leq \frac{1}{2\eta} + \frac{1}{2} \cdot \sum_{t=1}^n \frac{\eta^2}{2} \|y_t\|^2 \stackrel{\leq 1}{\sim}$$

$$\leq \frac{1}{2\eta} + \frac{\eta}{2} \cdot n$$

$$\eta = \frac{1}{\sqrt{n}} \Rightarrow R_n(\tilde{a}) \leq \sqrt{n}$$

② $A = P_{d-1} = d\text{-dimensional simplex of prob. dist.}$

$$y_t \in \Delta = [0, 1]^d \quad F(x) = \sum_{i=1}^d (x_i \ln x_i - x_i),$$

$$D_F(x, y) = \sum_{i=1}^d x_i \ln \left(\frac{x_i}{y_i} \right) + \sum_{i=1}^d (y_i - x_i)$$

$$\eta = \sqrt{2 \ln(d)/n}. \quad x, y \in P_{d-1} \Rightarrow \underbrace{\sum_{i=1}^d (y_i - x_i)}_{= 1} = \frac{\sum x_i}{n} = 1.$$

Prop. (Prop. 28.7 in textbook): $R_n \leq \sqrt{2n \ln(d)}$ ✓

with mirror descent.

Proof:

$$R_n(a) \leq \frac{F(a) - F(a_0)}{2} + \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle$$

$$- \frac{1}{2} \sum_{t=1}^n D_F(a_{t+1}, a_t)$$

↓ Pinsker's inequality.

$$\leq \frac{\ln(d)}{2} + \sum_{t=1}^n \|a_t - a_{t+1}\|_2 \|y_t\|_\infty - \frac{1}{2} \sum_{t=1}^n \|a_{t+1} - a_t\|_1^2$$

z

note ℓ_1 norm.

z^2

$$\text{Note: } z - \frac{z^2}{2} \leq \frac{1}{2}. \quad \forall z \in \mathbb{R}.$$

$$\sqrt{\frac{2 \ln(d)}{n}}$$

$$\Rightarrow R_n(a) \leq \frac{\ln(d)}{2} + \frac{n}{2}$$

A general bound (Corollary 28.8): F Legendre,

$\|\cdot\|_r$ a norm and $\|\cdot\|_{r^*}$ the dual-norm, s.t.

$$D_F(a_{t+1}, a_t) \geq \frac{1}{2} \|a_{t+1} - a_t\|_r^2$$

(i.e., F is 1-strongly convex with the $\|\cdot\|_r$ norm)

Then, with either mirror descent or FTRL, we have

$$R_n \leq \frac{\text{diameter}_F(A)}{\gamma} + \frac{\gamma}{2} \sum_{t=1}^n \|y_t\|_{r^*}^2$$

$$\|y\|_{r^*} = \max_{x: \|x\|_r \leq 1} \langle x, y \rangle$$

Bandit Feedback

Moving from full feedback (i.e. y_t is observed) to the bandit feedback setting (i.e., $\langle A_t, y_t \rangle$ is observed).

Setting: As in the usual adversarial setting, adversary aware of policy of player. Chooses loss sequence $\{y_1, y_2, \dots, y_n\}$, $y_t \in \Delta \subseteq \mathbb{R}^d$.

Player chooses actions : $\{A_1, A_2, \dots, A_n\}$, $A_t \in A$.

The actions now are randomized, hence notation of A_t instead of a_t .

$$\text{Regret: } R_n(a) = \mathbb{E} \left[\sum_{t=1}^n \langle A_t - a, y_t \rangle \right].$$

$$R_n = \max_{a \in A} R_n(a).$$

Algorithm Ideas:

- ① Use importance sampled estimator to determine unbiased estimate of y_t
- ② Use mirror descent / FTRL to evolve the expected action to take at the next step. Actual action is sampled from a suitably chosen dist. whose mean is determined above by Mirror descent / FTRL.

Bandit Mirror Descent / FTRL Algorithm:

Stochastic
M.D. / FTRL

Setting: F is Legendre, Learning rate $\eta > 0$, action set A , $\text{domain}(F) = D$.

"mean" action \rightarrow

Initial: $\bar{A}_1 = \underset{a \in A \cap D}{\operatorname{argmin}} F(a)$

P_t : Distribution over actions at time t , chosen such that

the variance of the IS estimator for \hat{Y}_t does not blow up. (Note that if P_t were deterministic / delta function, the estimator variance will blow up.) P_t a function of past actions / observations.

For each $t=1, 2, \dots, n$:

(a) Sample action $A_t \sim P_t$, with $E_{P_t}[A_t] = \bar{A}_t$, and observe $\langle A_t, y_t \rangle$.

(b) Construct \hat{Y}_t , an unbiased, finite variance estimate for y_t .

(c) Mirror Descent / FTRL:

$$(\text{Mirror Descent}) \quad \bar{A}_{t+1} = \underset{a \in \mathcal{A} \cap D}{\operatorname{argmin}} \quad \eta \langle a, \hat{Y}_t \rangle + D_F(a, \bar{A}_t).$$

$$(\text{FTRL}) \quad \bar{A}_{t+1} = \underset{a \in \mathcal{A} \cap D}{\operatorname{argmin}} \quad \eta \sum_{s=1}^t \langle a, \hat{Y}_s \rangle + F(a).$$

Thm (28.10 in textbook): F Legendre, A convex, $\eta > 0$, and s.t. $E[\hat{Y}_t | \bar{A}_t] = y_t$ $\forall t=1, \dots, n$, and (\hat{Y}_t, \bar{A}_t) well defined.

$$\text{Then, } R_n(a) \leq \mathbb{E} \left[\frac{F(a) - F(\bar{A}_1)}{\eta} + \sum_{t=1}^n \langle \bar{A}_t - \bar{A}_{t+1}, \hat{y}_t \rangle - \frac{1}{\eta} \sum_{t=1}^n D_F(\bar{A}_{t+1}, \bar{A}_t) \right]$$

Further, for Stochastic Mirror Descent where

$$\bar{A}_{t+1} = \underset{a \in A \cap D}{\text{argmin}} \eta \langle a, \hat{y}_t \rangle + D_F(a, \bar{A}_t),$$

and assuming \bar{A}_{t+1} lies in the interior of $A \cap D$,
for all $a \in A$ (a.s.), then

$$R_n \leq \frac{\text{diameter}_F(A)}{\eta} + \frac{1}{\eta} \sum_{t=1}^n \mathbb{E}[D(\bar{A}_t, \bar{A}_{t+1})]$$

Proof:

$$\mathbb{E}[\langle A_t, y_t \rangle] = \mathbb{E}[\langle \bar{A}_t, y_t \rangle]$$

$$\begin{aligned} &= \mathbb{E}\left[\mathbb{E}[\langle A_t, y_t \rangle | \bar{A}_t]\right] \\ &\xrightarrow{\text{unbiased estimator}} = \mathbb{E}\left[\mathbb{E}[\langle \bar{A}_t, \hat{y}_t \rangle | \bar{A}_t]\right] \end{aligned}$$

$$\Rightarrow R_n(a) = \mathbb{E}\left[\sum_{t=1}^n \langle A_t, y_t \rangle - \langle a, y_t \rangle\right]$$

$$= \mathbb{E} \left[\sum_{t=1}^T \langle \bar{A}_t - a, \hat{\gamma}_t \rangle \right].$$

Tracked by
 bandit M.D/FTRL

Rest of the proof is immediate from mirror descent / FTRL proof. \square

Explicit Construction of $\hat{\gamma}_t, p_t$ for $A = B_2^d$:

Linear Bandits on a unit B_2^d .

Setting: $A = B_2^d$ = unit ball in \mathbb{R}^d . = \mathcal{X} .

$S_t \in \{0, 1\}$, a selector r.v.

$V_t \in \{\pm e_1, \pm e_2, \dots, \pm e_d\}$, over $2d$ unit vectors.

s.t.

Note: $\|\bar{A}_t\|_2 < 1$, because we will operate FTRL over the action set $\tilde{A} = \{x : \|x\|_2 \leq r\}$, $r = 1 - 2/\sqrt{d} < 1$

$$\mathbb{E}[S_t | A_1, A_2, \dots, A_{t-1}] = \underbrace{1}_{\approx 0} - \underbrace{\|\bar{A}_t\|_2}_{\approx 1} \quad \text{and } S_t, V_t \text{ cond. indep.}$$

i.e., $P(S_t=1 | A_1, \dots, A_{t-1}) = 1 - \|\bar{A}_t\|_2$ if other r.v.s. given

$$V_t | (A_1, \dots, A_{t-1}) \sim \text{Uniform}\{\pm e_1, \dots, \pm e_d\}. \quad A_1, \dots, A_{t-1}$$

Choose

$$A_t = S_t U_t + \frac{(1-S_t) \bar{A}_t}{\|\bar{A}_t\|_2}$$

i.e., w.p. $1 - \|\bar{A}_t\|_2$, choose a random unit direction, and otherwise choose (normalized) action suggested by M.D./FTRL.

Observe $E[A_t | A_1, \dots, A_{t-1}] = 0 + \bar{A}_t$.

Finally,

$$\hat{\gamma}_t = \left(\frac{d S_t \langle A_t, y_t \rangle}{1 - \|\bar{A}_t\|_2} \right) A_t$$

unbiased estimator.

and $F(a) = -\ln(1 - \|a\|_2) - \|a\|$

Algo: FTRL over $\tilde{\mathcal{B}} = \{a \in \mathbb{R}^d : \|a\|_2 \leq r\}$,

for $\tau = (1 - 2\eta d) < 1$. $\left(\eta = \sqrt{\frac{\ln(n)}{3dn}} \right)$

controls variance of estimator

i.e.,

$$\bar{A}_{t+1} = \underset{a \in \tilde{\mathcal{B}} \cap \mathcal{D}}{\operatorname{argmin}} \eta \sum_{s=1}^t \langle a, \hat{\gamma}_s \rangle + F(a).$$

Facts: With above iteration and regularizer, and with

$$\hat{L}_{t-1} \doteq \sum_{s=1}^{t-1} \hat{\gamma}_s, \text{ we have: } \bar{A}_t = \tilde{\Pi}_2 \left(\frac{\eta \hat{L}_{t-1}}{1 + \eta \|\hat{L}_{t-1}\|} \right)$$

where $\tilde{\Pi}_2(x)$ is the ℓ_2 projection of x onto \tilde{A} .

Summary: At each time $t=1, 2, \dots, n$:

① Compute mean: $\bar{A}_t = \tilde{\Pi}_2 \left(\frac{\eta \hat{L}_{t-1}}{1 + \eta \|\hat{L}_{t-1}\|} \right)$, where

$$\hat{L}_{t-1} = \sum_{s=1}^{t-1} \hat{\gamma}_s$$

② $S_t \in \{0, 1\}$, with $S_t = \begin{cases} 1 & \text{w.p. } 1 - \|\bar{A}_t\|_2 \\ 0 & \text{w.p. } \|\bar{A}_t\|_2 \end{cases}$

Sample:

$$v_t \in \{\pm e_i, i=1, 2, \dots, d\} \text{ uniformly.}$$

③ Play: $A_t = S_t v_t + \frac{(1-S_t) \bar{A}_t}{\|\bar{A}_t\|_2}$

④ Observe | Estimate: Observe $\langle A_t, y_t \rangle$, and update

$$\hat{\gamma}_t = \left(\frac{d S_t \langle A_t, y_t \rangle}{1 - \|\bar{A}_t\|_2} \right) A_t$$

Theorem (28.11 in textbook): $\{y_t, t=1, \dots, n\} \subseteq \mathcal{L} = \mathbb{B}_2^d$.

Then, with $r = 1 - 2\eta d$, $\eta = \sqrt{\frac{\ln(n)}{3dn}}$, we have

$$R_n \leq 2 \sqrt{3nd \ln(n)}$$

Note: This gives \sqrt{d} scaling, whereas previous notes based on Kiefer-Wolfowitz had scaling linearly in d .

Proof: $\hat{a}^* = \arg \min_a \sum_{t=1}^n \langle a, y_t \rangle$, the optimal action.

$$R_n = \mathbb{E} \left[\sum_{t=1}^n \langle A_t - \hat{a}^*, y_t \rangle \right]$$

$$= \mathbb{E} \left[\sum_{t=1}^n \langle A_t - r\hat{a}^*, y_t \rangle \right] + (1-r) \sum_{t=1}^n \langle \hat{a}^*, y_t \rangle$$

$$\leq \mathbb{E} \left[\sum_{t=1}^n \langle A_t - r\hat{a}^*, y_t \rangle \right] + (1-r)n.$$

Now, recall: ψ Legendre, $\eta > 0$, $x, y \in \text{interior } (\mathcal{D})$,

① ψ twice differentiable. Let $z \in [x, y]$ s.t.

$$D_\psi(x, y) = \frac{1}{2} \|x - y\|_{\nabla^2 \psi(z)}^2. \text{ Then, for any } u \in \mathbb{R}^d,$$

$$\langle x-y, u \rangle - \frac{D_\psi(x, y)}{\eta} \leq \frac{\eta}{2} \|u\|^2 \Delta_{\psi(z)}^{-1}$$

(Thm 26.13 in textbook) : Uses Cauchy-Schwartz + Calculus.

② For FTRL, we have from Thm 28.5,

$$R_n(a) \leq \underbrace{F(a) - F(a_1)}_{\eta} + \sum_{t=1}^n \langle a_t - a_{t+1}, y_t \rangle - \frac{1}{\eta} \sum_{t=1}^n D_F(a_{t+1}, a_t)$$

$$\therefore R_n(r\alpha^*) = E \left[\sum_{t=1}^n \langle a_t - r\alpha^*, y_t \rangle \right]$$

look at
page 22/23
in these
notes.

$$= E \left[\sum_{t=1}^n \langle \bar{A}_t - r\alpha^*, \hat{y}_t \rangle \right]$$

$$\leq \frac{\text{Diameter}_F(\bar{A})}{\eta} + \frac{1}{2} E \left[\sum_{t=1}^n \|\hat{y}_t\|^2 \Delta_{F(z_t)}^{-1} \right],$$

$\stackrel{s \ln(\frac{1}{1-r})}{\curvearrowright} \quad \stackrel{E[\cdot] \leq 2d}{\curvearrowright}$

for $z_t \in [\bar{A}_t, \bar{A}_{t+1}]$.

$$\text{Further, } \eta \|\hat{\gamma}_t\|_2 \leq \left(\frac{\eta d}{1-\gamma} \right) = \frac{1}{2} \quad \left(\text{defn. of } \gamma = 1 - 2\eta d \right)$$

$$\curvearrowleft = \alpha \bar{A}_t + (1-\alpha) \bar{A}_{t+1}, \text{ for some } \alpha \in [0, 1]$$

$$\therefore \frac{1 - \|Z_t\|_2}{1 - \|\bar{A}_t\|_2} \leq \sup_{\alpha \in [0, 1]} \frac{1 - \|\alpha \bar{A}_t + (1-\alpha) \bar{A}_{t+1}\|_2}{1 - \|\bar{A}_t\|_2} \quad \begin{cases} \text{(by default),} \\ \|\cdot\| = \|\cdot\|_2, \\ \text{unless otherwise stated} \end{cases}$$

$$? \quad = \max \left\{ 1, \frac{1 - \|\bar{A}_{t+1}\|_2}{1 - \|\bar{A}_t\|_2} \chi_{\{\|\bar{A}_{t+1}\|_2 < \|\bar{A}_t\|_2\}} \right\}$$

$$\leq \max \left\{ 1, \frac{1 + \eta \|\hat{\gamma}_{t-1}\|_2}{1 + \eta \|\hat{\gamma}_t\|_2} \right\} \quad \begin{array}{l} \text{! Recall:} \\ \bar{A}_t = \tilde{T}_2 \left(\frac{\eta \hat{\gamma}_{t-1}}{1 + \eta \|\hat{\gamma}_{t-1}\|_2} \right) \\ \hat{\gamma}_t = \hat{\gamma}_{t-1} + \hat{\gamma}_t \\ \eta \|\hat{\gamma}_t\|_2 \leq \eta \|\hat{\gamma}_{t-1}\|_2 + \frac{1}{2} \end{array}$$

$$\leq \max \left\{ 1, \frac{1 + \eta \|\hat{\gamma}_{t-1}\|_2}{\frac{1}{2} + \eta \|\hat{\gamma}_{t-1}\|_2} \right\} \leq 2.$$

$$\text{Next, } \nabla^2 F(a) = \frac{I}{1 - \|a\|} + \frac{aa^\top}{\|a\|(1 - \|a\|)^2} \geq \frac{I}{1 - \|a\|}$$

$$\Rightarrow \nabla^2 F(a)^{-1} \leq (1 - \|a\|) I$$

$$\Rightarrow \|\hat{\gamma}_t\|_{\nabla^2 F(a)^{-1}}^2 = \hat{\gamma}_t^\top \nabla^2 F(a)^{-1} \hat{\gamma}_t \leq (1 - \|a\|) \|\hat{\gamma}_t\|_2^2$$

$$\therefore E \left[\|\hat{y}_t\|_{\nabla^2 F(z_t)^{-1}}^2 \right] \leq E \left[(1 - \|z_t\|) \|\hat{y}_t\|_2^2 \right]$$

$$= d^2 E \left[\frac{(1 - \|z_t\|) S_t \langle v_t, y_t \rangle^2}{(1 - \|\bar{A}_t\|)^2} \right]$$

From ④ :

$$\left(\frac{1 - \|z_t\|}{1 - \|\bar{A}_t\|} \right) \leq 2$$

$$\Rightarrow E \left[\|\hat{y}_t\|_{\nabla^2 F(z_t)^{-1}}^2 \right] \leq 2d^2 E \left[\frac{S_t}{1 - \|\bar{A}_t\|} \right] E \left[\langle v_t, y_t \rangle^2 \right]$$

$$\leq 2d \cdot = \sum_{i=1}^d \frac{y_i^2}{2d} + \sum_{i=1}^d \frac{(-y_i)^2}{2d}$$

$$\text{Finally, } \text{diameter}(\tilde{\Delta}) \leq \ln \left(\frac{1}{1-\gamma} \right). \quad = \frac{1}{d} \sum_{i=1}^d y_i^2 \leq \frac{1}{d}.$$

$$\therefore R_n \leq (1-\gamma)n + \frac{1}{\eta} \ln \left(\frac{1}{1-\gamma} \right) + \eta n d$$

$$= \frac{1}{\eta} \ln \left(\frac{1}{2\eta d} \right) + 3\eta n d$$

$$\leq 2 \sqrt{3nd \ln(n)}$$

◻