

CMPT 409/981: Optimization for Machine Learning

Lecture 15

Sharan Vaswani

October 31, 2024

Recap: Online Optimization

Generic Online Optimization (w_0 , Algorithm \mathcal{A} , Convex set $\mathcal{C} \subseteq \mathbb{R}^d$)

- 1: **for** $k = 1, \dots, T$ **do**
 - 2: Algorithm \mathcal{A} chooses point (decision) $w_k \in \mathcal{C}$
 - 3: Environment chooses and reveals the (potentially adversarial) loss function $f_k : \mathcal{C} \rightarrow \mathbb{R}$
 - 4: Algorithm suffers a cost $f_k(w_k)$
 - 5: **end for**
-

Application: Prediction from Expert Advice: Given d experts,
 $\mathcal{C} = \Delta_d = \{w_i | w_i \geq 0 ; \sum_{i=1}^d w_i = 1\}$ and $f_k(w_k) = \langle c_k, w_k \rangle$ where $c_k \in \mathbb{R}^d$ is the loss vector.

Application: Imitation Learning: Given access to an expert that knows what action $a \in [A]$ to take in each state $s \in [S]$, learn a policy $\pi : [S] \rightarrow [A]$ that imitates the expert, i.e. we want that $\pi(a|s) \approx \pi_{\text{expert}}(a|s)$. Here, $w = \pi$ and $\mathcal{C} = \Delta_A \times \Delta_A \dots \Delta_A$ (simplex for each state) and f_k is a measure of discrepancy between π_k and π_{expert} .

Online Optimization

- Recall that the sequence of losses $\{f_k\}_{k=1}^T$ is potentially adversarial and can also depend on w_k .
- **Objective:** Do well against the *best fixed decision in hindsight*, i.e. if we knew the entire sequence of losses beforehand, we would choose $w^* := \arg \min_{w \in \mathcal{C}} \sum_{k=1}^T f_k(w)$.
- **Regret:** For any fixed decision $u \in \mathcal{C}$,

$$R_T(u) := \sum_{k=1}^T [f_k(w_k) - f_k(u)]$$

When comparing against the best decision in hindsight,

$$R_T := \sum_{k=1}^T [f_k(w_k)] - \min_{w \in \mathcal{C}} \sum_{k=1}^T f_k(w).$$

- We want to design algorithms that achieve a *sublinear regret* (that grows as $o(T)$). A sublinear regret implies that the performance of our sequence of decisions is approaching that of w^* .

- **Online Convex Optimization (OCO):** When the losses f_k are (strongly) convex loss functions.

Example 1: In prediction with expert advice, $f_k(w) = \langle c_k, w \rangle$ is a linear function.

Example 2: In imitation learning, $f_k(\pi) = \mathbb{E}_{s \sim d^{\pi_k}} [\text{KL}(\pi(\cdot|s) || \pi_{\text{expert}}(\cdot|s))]$ where d^{π_k} is a distribution over the states induced by running policy π_k .

Example 3: In online control such as LQR (linear quadratic regulator) with unknown costs/perturbations, f_k is quadratic.

- In Examples 2-3, the loss at iteration $k + 1$ depends on the *learner's* decision at iteration k .

Online Convex Optimization

- **Online-to-Batch conversion:** If the sequence of loss functions is i.i.d from some fixed distribution, we can convert the regret guarantees into the traditional convergence guarantees for the resulting algorithm.

Formally, if f_k are convex and $R(T) = O(\sqrt{T})$, then taking the expectation w.r.t the distribution generating the losses,

$$\mathbb{E} \left[\frac{R_T}{T} \right] = \mathbb{E} \left[\frac{\sum_{k=1}^T [f_k(w_k)] - \sum_{k=1}^T f_k(w^*)}{T} \right] \geq \sum_{k=1}^T [f(\bar{w}_T) - f(w^*)] = O \left(\frac{1}{\sqrt{T}} \right)$$

where $f(w) := \mathbb{E}[f_k(w)]$ (since the losses are i.i.d) and $\bar{w}_T := \frac{\sum_{k=1}^T w_k}{T}$ (since the losses are convex, we used Jensen's inequality).

- If the distribution generating the losses is a uniform discrete distribution on n fixed data-points, then $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$ and we are back in the finite-sum minimization setting.
- Hence, algorithms that attain $R(T) = O(\sqrt{T})$ can result in an $O \left(\frac{1}{\sqrt{T}} \right)$ convergence (in terms of the function values) for convex losses.

Questions?

Online Gradient Descent

The simplest algorithm that results in sublinear regret for OCO is *Online Gradient Descent*.

Online Gradient Descent (OGD): At iteration k , the algorithm chooses the point w_k . After the loss function f_k is revealed, OGD suffers a cost $f_k(w_k)$ and uses the function to compute

$$w_{k+1} = \Pi_C[w_k - \eta_k \nabla f_k(w_k)]$$

where $\Pi_C[x] = \arg \min_{y \in \mathcal{C}} \frac{1}{2} \|y - x\|^2$.

Claim: If the convex set \mathcal{C} has a diameter D i.e. for all $x, y \in \mathcal{C}$, $\|x - y\| \leq D$, for an arbitrary sequence of losses such that each f_k is convex and differentiable, OGD with a non-increasing sequence of step-sizes i.e. $\eta_k \leq \eta_{k-1}$ and $w_1 \in \mathcal{C}$ has the following regret for all $u \in \mathcal{C}$,

$$R_T(u) \leq \frac{D^2}{2\eta_T} + \sum_{k=1}^T \frac{\eta_k}{2} \|\nabla f_k(w_k)\|^2$$

Online Gradient Descent - Convex functions

Proof: Using the update $w_{k+1} = \Pi_{\mathcal{C}}[w_k - \eta_k \nabla f_k(w_k)]$. Since $u \in \mathcal{C}$,

$$\|w_{k+1} - u\|^2 = \|\Pi_{\mathcal{C}}[w_k - \eta_k \nabla f_k(w_k)] - u\|^2 = \|\Pi_{\mathcal{C}}[w_k - \eta_k \nabla f_k(w_k)] - \Pi_{\mathcal{C}}[u]\|^2$$

Since projections are non-expansive i.e. for all x, y , $\|\Pi_{\mathcal{C}}[y] - \Pi_{\mathcal{C}}[x]\| \leq \|y - x\|$,

$$\begin{aligned} &\leq \|w_k - \eta_k \nabla f_k(w_k) - u\|^2 \\ &= \|w_k - u\|^2 - 2\eta_k \langle \nabla f_k(w_k), w_k - u \rangle + \eta_k^2 \|\nabla f_k(w_k)\|^2 \\ &\leq \|w_k - u\|^2 - 2\eta_k [f_k(w_k) - f_k(u)] + \eta_k^2 \|\nabla f_k(w_k)\|^2 \\ &\hspace{25em} (\text{Since } f_k \text{ is convex}) \end{aligned}$$

$$\begin{aligned} \implies 2\eta_k [f_k(w_k) - f_k(u)] &\leq [\|w_k - u\|^2 - \|w_{k+1} - u\|^2] + \eta_k^2 \|\nabla f_k(w_k)\|^2 \\ \implies R_T(u) &\leq \sum_{k=1}^T \left[\frac{\|w_k - u\|^2 - \|w_{k+1} - u\|^2}{2\eta_k} \right] + \sum_{k=1}^T \frac{\eta_k}{2} \|\nabla f_k(w_k)\|^2 \end{aligned}$$

Online Gradient Descent - Convex functions

Recall that $R_T(u) \leq \sum_{k=1}^T \left[\frac{\|w_k - u\|^2 - \|w_{k+1} - u\|^2}{2\eta_k} \right] + \sum_{k=1}^T \frac{\eta_k}{2} \|\nabla f_k(w_k)\|^2$.

$$\begin{aligned} & \sum_{k=1}^T \left[\frac{\|w_k - u\|^2 - \|w_{k+1} - u\|^2}{2\eta_k} \right] \\ &= \sum_{k=2}^T \left[\|w_k - u\|^2 \underbrace{\left(\frac{1}{2\eta_k} - \frac{1}{2\eta_{k-1}} \right)}_{\text{Non-negative since } \eta_k \leq \eta_{k-1}} \right] + \frac{\|w_1 - u\|^2}{2\eta_1} - \frac{\|w_{T+1} - u\|^2}{2\eta_T} \\ &\leq D^2 \sum_{k=2}^T \left[\frac{1}{2\eta_k} - \frac{1}{2\eta_{k-1}} \right] + \frac{D^2}{2\eta_1} = D^2 \left[\frac{1}{2\eta_T} - \frac{1}{2\eta_1} \right] + \frac{D^2}{2\eta_1} = \frac{D^2}{2\eta_T} \\ &\hspace{15em} (\text{Since } \|x - y\| \leq D \text{ for all } x, y \in \mathcal{C}) \end{aligned}$$

Putting everything together,

$$R_T(u) \leq \frac{D^2}{2\eta_T} + \sum_{k=1}^T \frac{\eta_k}{2} \|\nabla f_k(w_k)\|^2$$

Online Gradient Descent - Convex, Lipschitz functions

Claim: If the convex set \mathcal{C} has a diameter D i.e. for all $x, y \in \mathcal{C}$, $\|x - y\| \leq D$, for an arbitrary sequence of losses such that each f_k is convex, differentiable and G -Lipschitz, OGD with $\eta_k = \frac{\eta}{\sqrt{k}}$ and $w_1 \in \mathcal{C}$ has the following regret for all $u \in \mathcal{C}$,

$$R_T(u) \leq \frac{D^2 \sqrt{T}}{2\eta} + G^2 \sqrt{T} \eta$$

Proof: Since the step-size is decreasing, we can use the general result from the previous slide,

$$R_T(u) \leq \frac{D^2}{2\eta_T} + \sum_{k=1}^T \frac{\eta_k}{2} \|\nabla f_k(w_k)\|^2 \leq \frac{D^2}{2\eta_T} + \frac{G^2}{2} \sum_{k=1}^T \eta_k \quad (\text{Since } f_k \text{ is } G\text{-Lipschitz})$$

$$\implies R_T(u) \leq \frac{D^2 \sqrt{T}}{2\eta} + \frac{G^2 \eta}{2} \sum_{k=1}^T \frac{1}{\sqrt{k}} \leq \frac{D^2 \sqrt{T}}{2\eta} + G^2 \sqrt{T} \eta \quad (\text{Since } \sum_{k=1}^T \frac{1}{\sqrt{k}} \leq 2\sqrt{T})$$

- In order to find the “best” η , set it such that $D^2/2\eta = G^2\eta$, implying that $\eta = D/\sqrt{2}G$ and $R_T(u) \leq \sqrt{2} DG \sqrt{T}$. Hence, OGD with a decreasing step-size attains sublinear $\Theta(\sqrt{T})$ regret for convex, Lipschitz functions.

Questions?

Online Mirror Descent

- The OGD update at iteration k can also be written as:

$$w_{k+1} = \arg \min_{w \in \mathcal{C}} \left[\langle \nabla f_k(w_k), w \rangle + \frac{1}{2\eta_k} \|w - w_k\|_2^2 \right]$$

- Online Mirror Descent (OMD) generalizes gradient descent by choosing a strictly convex, differentiable function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ (referred to as the *mirror map*) to induce a distance measure.
- ϕ induces the *Bregman divergence* $D_\phi(\cdot, \cdot)$, a distance measure between points x, y ,

$$D_\phi(y, x) := \phi(y) - \phi(x) - \langle \nabla \phi(x), y - x \rangle.$$

Geometrically, $D_\phi(y, x)$ is the distance between the function $\phi(y)$ and the line $\phi(x) + \langle \nabla \phi(x), y - x \rangle$ which is tangent to the function at x .

- Using D_ϕ as the distance measure results in the mirror descent update:

$$w_{k+1} = \arg \min_{w \in \mathcal{C}} \left[\langle \nabla f_k(w_k), w \rangle + \frac{1}{\eta_k} D_\phi(w, w_k) \right]$$

- Setting $\phi(x) = \frac{1}{2} \|x\|^2$ results in $D_\phi(y, x) = \frac{1}{2} \|y - x\|^2$ and recovers OGD.

Online Mirror Descent – Example

- For prediction with expert advice, $\mathcal{C} = \Delta_d = \{w_i | w_i \geq 0 ; \sum_{i=1}^d w_i = 1\}$ and we want a distance metric between probabilities.
- Typically use the *negative-entropy mirror map* i.e. $\phi(w) = \sum_{i=1}^d w_i \ln(w_i)$.
- For $u, v \in \mathcal{C}$, the corresponding Bregman divergence $D_\phi(u, v)$ can be calculated as:

$$D_\phi(u, v) = \phi(u) - \phi(v) - \langle \nabla \phi(v), u - v \rangle = \phi(u) - \phi(v) - \langle \log(v) + 1, u - v \rangle$$

($\nabla \phi(u) = \log(u) + 1$, where $\log(\cdot)$ is element-wise)

$$\begin{aligned} &= \sum_{i=1}^d u_i \log(u_i) - \sum_{i=1}^d v_i \log(v_i) - \left[\sum_{i=1}^d u_i \log(v_i) - \sum_{i=1}^d v_i \log(v_i) \right] - \sum_{i=1}^d (u_i - v_i) \\ &= \sum_{i=1}^d u_i \log\left(\frac{u_i}{v_i}\right) = \text{KL}(u||v). \end{aligned} \quad (\sum_{i=1}^d u_i = \sum_{i=1}^d v_i = 1)$$

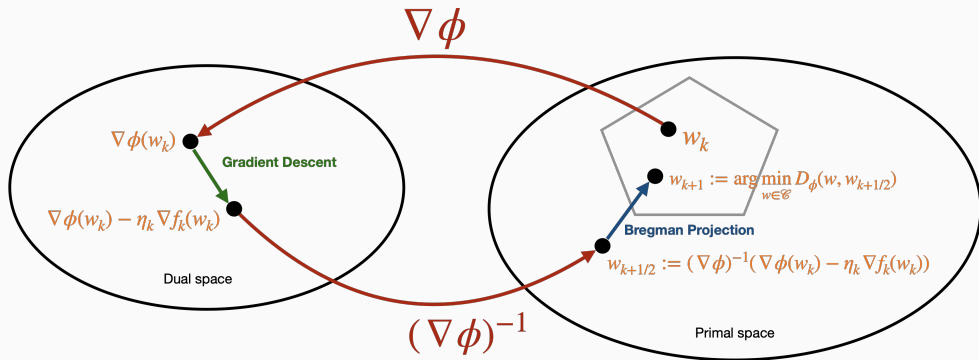
- The KL-divergence is a standard way to measure the distance between probability distributions. For distributions u, v , $\text{KL}(u||v) := \sum_{i=1}^d u_i \log\left(\frac{u_i}{v_i}\right)$ is non-negative and equal to zero iff $u = v$.

Online Mirror Descent

The OMD update can be equivalently written as:

GD in dual space: $w_{k+1/2} = (\nabla\phi)^{-1}(\nabla\phi(w_k) - \eta_k \nabla f_k(w_k))$

Bregman projection: $w_{k+1} = \arg \min_{w \in \mathcal{C}} D_\phi(w, w_{k+1/2})$



Prove in Assignment 3!

Online Mirror Descent – Example

For prediction with expert advice, $\mathcal{C} = \Delta_d = \{w_i | w_i \geq 0 ; \sum_{i=1}^d w_i = 1\}$, $\phi(w) = \sum_{i=1}^d w_i \ln(w_i)$ is the negative-entropy mirror map and $g_k := \nabla f_k(w_k)$, then the OMD update can be written as: (prove in Assignment 3!)

- **GD in dual space:** $w_{k+1/2}[i] = w_k[i] \exp(-\eta_k g_k[i])$
- **Bregman projection:** $w_{k+1}[i] = \frac{w_{k+1/2}[i]}{\|w_{k+1/2}\|_1}$
- **Multiplicative weights update:**

$$w_{k+1}[i] = \frac{w_k[i] \exp(-\eta_k g_k[i])}{\sum_{j=1}^d w_k[j] \exp(-\eta_k g_k[j])}$$

If $w_0[i] = \frac{1}{d}$ for all $i \in [d]$, then, for all k ,

$$w_{k+1}[i] = \frac{\exp\left(-\sum_{m=1}^k \eta_m g_m[i]\right)}{\sum_{j=1}^d \exp\left(-\sum_{m=1}^k \eta_m g_m[j]\right)}$$

Online Mirror Descent – Convex, Lipschitz functions

In order to analyze OMD, we will make some assumptions about \mathcal{C} , f_k and ϕ .

- **Assumption 1:** \mathcal{C} is a convex set and $\forall k$, f_k is a convex function.
- **Assumption 2:** $\forall k$, f_k is G -Lipschitz in the ℓ_p norm (for $p \geq 1$), implying that $\forall w \in \mathcal{C}$,

$$\|\nabla f_k(w)\|_p \leq G$$

- **Assumption 3:** ϕ is ν strongly-convex in the ℓ_q norm (for $q \geq 1$ s.t. $\frac{1}{p} + \frac{1}{q} = 1$) i.e.

$$\phi(y) \geq \phi(x) + \langle \nabla \phi(x), y - x \rangle + \frac{\nu}{2} \|y - x\|_q^2$$

- *Example:* For prediction from expert advice,
 - $\mathcal{C} = \Delta_d$ is a convex set and $f_k(w_k) = \langle c_k, w_k \rangle$ is a convex function.
 - If the costs are bounded by M , then, $\|\nabla f_k(w)\|_\infty = \|c_k\|_\infty \leq M$. Hence, $p = \infty$, $G = M$.
 - If $\phi(w)$ is negative-entropy, then by Pinsker's inequality, $q = 1$ and $\nu = 1$ i.e.

$$\phi(y) - \phi(x) - \langle \nabla \phi(x), y - x \rangle = D_\phi(y, x) = \text{KL}(y||x) \geq \frac{1}{2} \|y - x\|_1^2.$$

Online Mirror Descent – Convex, Lipschitz functions

Claim: For an arbitrary sequence of losses such that each f_k is convex, G -Lipschitz and differentiable, then OMD with a ν strongly-convex mirror map ϕ , $\eta_k = \eta = \sqrt{\frac{2\nu}{T}} \frac{D}{G}$ where $D^2 := \max_{u \in \mathcal{C}} D_\phi(u, w_1)$ has the following regret for all $u \in \mathcal{C}$,

$$R_T(u) \leq \frac{\sqrt{2} DG}{\sqrt{\nu}} \sqrt{T},$$

Proof: Recall the mirror descent update: $\nabla\phi(w_{k+1/2}) = \nabla\phi(w_k) - \eta_k \nabla f_k(w_k)$. Setting $\eta_k = \eta$ and using the definition of regret,

$$\begin{aligned} R_T(u) &= \sum_{k=1}^T f_k(w_k) - f_k(u) \leq \sum_{k=1}^T [\langle g_k, w_k - u \rangle] && \text{(Convexity of } f_k \text{ and } g_k := \nabla f_k(w_k)) \\ &= \sum_{k=1}^T \frac{1}{\eta} \langle \nabla\phi(w_k) - \nabla\phi(w_{k+1/2}), w_k - u \rangle && \text{(Using the OMD update)} \end{aligned}$$

Online Mirror Descent – Convex, Lipschitz functions

Recall that $R_T(u) = \sum_{k=1}^T \frac{1}{\eta} \langle \nabla \phi(w_k) - \nabla \phi(w_{k+1/2}), w_k - u \rangle$

Three point property: for any 3 points x, y, z ,

$$\langle \nabla \phi(z) - \nabla \phi(y), z - x \rangle = D_\phi(x, z) + D_\phi(z, y) - D_\phi(x, y)$$

$$\langle \nabla \phi(w_k) - \nabla \phi(w_{k+1/2}), w_k - u \rangle = D_\phi(u, w_k) + D_\phi(w_k, w_{k+1/2}) - D_\phi(u, w_{k+1/2})$$

$$\implies R_T(u) = \sum_{k=1}^T \frac{1}{\eta} [D_\phi(u, w_k) + D_\phi(w_k, w_{k+1/2}) - D_\phi(u, w_{k+1/2})]$$

From the OMD update, we know that, $w_{k+1} = \arg \min_{w \in \mathcal{W}} D_\phi(w, w_{k+1/2})$. Recall the optimality condition: for a convex function f and a convex set \mathcal{C} , if $x^* = \arg \min_{x \in \mathcal{C}} f(x)$, then $\forall x \in \mathcal{C}$, $\langle \nabla f(x^*), x^* - x \rangle \leq 0$. Using this condition for $D_\phi(w, w_{k+1/2})$, for $u \in \mathcal{C}$,

$$\langle \nabla \phi(w_{k+1}) - \nabla \phi(w_{k+1/2}), w_{k+1} - u \rangle \leq 0$$

$$\implies -D_\phi(u, w_{k+1/2}) \leq -D_\phi(u, w_{k+1}) - D_\phi(w_{k+1}, w_{k+1/2}) \quad (3 \text{ point property})$$

$$\implies R_T(u) \leq \sum_{k=1}^T \frac{1}{\eta} [D_\phi(u, w_k) - D_\phi(u, w_{k+1})] + \frac{1}{\eta} [D_\phi(w_k, w_{k+1/2}) - D_\phi(w_{k+1}, w_{k+1/2})]$$

Online Mirror Descent – Convex, Lipschitz functions

Telescoping we conclude that $R_T(u) \leq \frac{1}{\eta} D_\phi(u, w_1) + \frac{1}{\eta} \sum_{k=1}^T [D_\phi(w_k, w_{k+1/2}) - D_\phi(w_{k+1}, w_{k+1/2})]$.

$$\begin{aligned} D_\phi(w_k, w_{k+1/2}) - D_\phi(w_{k+1}, w_{k+1/2}) &= \phi(w_k) - \phi(w_{k+1}) - \langle \nabla \phi(w_{k+1/2}), w_k - w_{k+1} \rangle \\ &\leq \langle \nabla \phi(w_k) - \nabla \phi(w_{k+1/2}), w_k - w_{k+1} \rangle - \frac{\nu}{2} \|w_k - w_{k+1}\|_q^2 \\ &\quad \text{(Using strong-convexity of } \phi \text{ with } y = w_{k+1} \text{ and } x = w_k) \end{aligned}$$

$$= \eta \langle g_k, w_k - w_{k+1} \rangle - \frac{\nu}{2} \|w_k - w_{k+1}\|_q^2 \quad \text{(Using the OMD update)}$$

$$\leq \eta G \|w_k - w_{k+1}\|_q - \frac{\nu}{2} \|w_k - w_{k+1}\|_q^2$$

$$\text{(Holder's inequality: } \langle x, y \rangle \leq \|x\|_p \|y\|_q \text{ s.t. } \frac{1}{p} + \frac{1}{q} = 1 \text{ and since } \|g_k\|_p \leq G)$$

$$\leq \frac{\eta^2 G^2}{2\nu} \quad \text{(For all } z, a z - b z^2 \leq \frac{a^2}{4b})$$

$$\implies R_T(u) \leq \frac{1}{\eta} D_\phi(u, w_1) + \frac{\eta G^2 T}{2\nu} \leq \frac{D^2}{\eta} + \frac{\eta G^2 T}{2\nu} \quad \text{(Since } D_\phi(u, w_1) \leq D^2)$$

$$\implies R_T(u) \leq \frac{\sqrt{2} D G}{\sqrt{\nu}} \sqrt{T} \quad \text{(Setting } \eta = \sqrt{\frac{2\nu}{T}} \frac{D}{G})$$

Online Mirror Descent – Example

We have proved that for any fixed comparator u , $R_T(u) \leq \frac{\sqrt{2}DG}{\sqrt{\nu}} \sqrt{T}$ where,

(i) $\|\nabla f_k(w)\|_p \leq G$, (ii) $\phi(y) \geq \phi(x) + \langle \nabla \phi(x), y - x \rangle + \frac{\nu}{2} \|y - x\|_q^2$ and (iii) $D_\phi(u, w_1) \leq D^2$.

- Using OMD with negative-entropy for prediction with expert advice, $p = \infty$, $q = 1$, $\nu = 1$.

Since $\|c_k\|_\infty \leq M$, $G = M$. If $\forall i \in [d]$, $w_1[i] = \frac{1}{d}$, $D_\phi(u, w_1) = \sum_{i=1}^d u_i \ln(u_i d) \leq \ln(d)$.

$$\implies R_T(u) \leq \sqrt{2}M \sqrt{\ln(d)} \sqrt{T}$$

- Since OGD is a special case of OMD with $\phi(w) = \frac{1}{2} \|w\|^2$, using OGD for prediction with expert advice, $p = 2$, $q = 2$, $\nu = 1$. Since $\|c_k\|_\infty \leq M$, using the relation between norms, $G = M\sqrt{d}$. If $\forall i \in [d]$, $w_1[i] = \frac{1}{d}$, $D_\phi(u, w_1) = \frac{1}{2} \|u - w_1\|^2 \leq \sqrt{2}$

$$\implies R_T(u) \leq 2M \sqrt{d} \sqrt{T}$$

- Hence, using multiplicative weights results in $O(\sqrt{\ln(d)}\sqrt{T})$ regret which is better than the $O(\sqrt{d} \sqrt{T})$ regret obtained by OGD. For prediction with expert advice, when the number of experts is large, this can be a substantial advantage.

Questions?

Online Gradient Descent - Strongly-convex, Lipschitz functions

Claim: If the convex set \mathcal{C} has a diameter D , for an arbitrary sequence of losses such that each f_k is μ_k strongly-convex (s.t. $\mu := \min_{k \in [T]} \mu_k > 0$), G -Lipschitz and differentiable, then OGD with $\eta_k = \frac{1}{\sum_{i=1}^k \mu_i}$ and $w_1 \in \mathcal{C}$ has the following regret for all $u \in \mathcal{C}$,

$$R_T(u) \leq \frac{G^2}{2\mu} (1 + \log(T))$$

Proof: Similar to the convex proof, use the update $w_{k+1} = \Pi_{\mathcal{C}}[w_k - \eta_k \nabla f_k(w_k)]$. Since $u \in \mathcal{C}$,

$$\begin{aligned} \|w_{k+1} - u\|^2 &= \|\Pi_{\mathcal{C}}[w_k - \eta_k \nabla f_k(w_k)] - u\|^2 = \|\Pi_{\mathcal{C}}[w_k - \eta_k \nabla f_k(w_k)] - \Pi_{\mathcal{C}}[u]\|^2 \\ &\leq \|w_k - u\|^2 - 2\eta_k \langle \nabla f_k(w_k), w_k - u \rangle + \eta_k^2 \|\nabla f_k(w_k)\|^2 \\ &\leq \|w_k - u\|^2 (1 - \mu_k \eta_k) - 2\eta_k [f_k(w_k) - f_k(u)] + \eta_k^2 \|\nabla f_k(w_k)\|^2 \\ &\hspace{15em} (\text{Since } f_k \text{ is } \mu_k \text{ strongly-convex}) \end{aligned}$$

$$\begin{aligned} \Rightarrow R_T(u) &\leq \sum_{k=1}^T \left[\frac{\|w_k - u\|^2 (1 - \mu_k \eta_k) - \|w_{k+1} - u\|^2}{2\eta_k} \right] + \frac{G^2}{2} \sum_{k=1}^T \eta_k \\ &\hspace{15em} (\text{Since } f_k \text{ is } G\text{-Lipschitz}) \end{aligned}$$

Online Gradient Descent - Strongly-convex, Lipschitz functions

Recall that $R_T(u) \leq \sum_{k=1}^T \left[\frac{\|w_k - u\|^2 (1 - \mu_k \eta_k) - \|w_{k+1} - u\|^2}{2\eta_k} \right] + \frac{G^2}{2} \sum_{k=1}^T \eta_k$.

$$\begin{aligned} & \sum_{k=1}^T \left[\frac{\|w_k - u\|^2 (1 - \mu_k \eta_k) - \|w_{k+1} - u\|^2}{2\eta_k} \right] \\ &= \sum_{k=2}^T \left[\|w_k - u\|^2 \underbrace{\left(\frac{1}{2\eta_k} - \frac{1}{2\eta_{k-1}} - \frac{\mu_k}{2} \right)}_{=0} \right] + \|w_1 - u\|^2 \underbrace{\left[\frac{1}{2\eta_1} - \frac{\mu_1}{2} \right]}_{=0} - \frac{\|w_{T+1} - u\|^2}{2\eta_T} \leq 0 \end{aligned}$$

(Since $\eta_k = \frac{1}{\sum_{i=1}^k \mu_i}$)

Putting everything together,

$$R_T(u) \leq \frac{G^2}{2} \sum_{k=1}^T \frac{1}{\mu k} \leq \frac{G^2}{2\mu} (1 + \log(T))$$

(Since $\mu := \min_{k \in [T]} \mu_k$ and $\sum_{k=1}^T 1/k \leq 1 + \log(T)$)

Lower Bound: There is an $\Omega(\log(T))$ lower-bound on the regret for strongly-convex, Lipschitz functions and hence OGD is optimal (in terms of T) for this setting!

Follow the Leader

Common algorithm that achieves logarithmic regret for strongly-convex losses.

Follow the Leader (FTL): At iteration k , the algorithm chooses the point w_k . After the loss function f_k is revealed, FTL suffers a cost $f_k(w_k)$ and uses it to compute

$$w_{k+1} = \arg \min_{w \in \mathcal{C}} \sum_{i=1}^k f_i(w).$$

- × Needs to solve a deterministic optimization sub-problem which can be expensive.
- × Needs to store all the previous loss functions and requires $O(T)$ memory.
- ✓ Does not require any step-size and is hyper-parameter free.
 - In applications such Imitation Learning (IL), interacting with the environment and getting access to f_k is expensive. FTL allows multiple policy updates (when solving the sub-problem) and helps better reuse the collected data. FTL is a standard method to solve online IL problems and the resulting algorithm is known as DAGGER [RGB11].
 - Compared to FTL, OGD requires an environment interaction for each policy update.

Follow the Leader and OGD

To connect FTL and OGD, consider the case when $\mathcal{C} = \mathbb{R}^d$.

$$w_{k+1} = \arg \min_{w \in \mathbb{R}} \sum_{i=1}^k [f_i(w)] \implies \sum_{i=1}^k \nabla f_i(w_{k+1}) = 0$$

- If we define $\tilde{f}_i(w)$ to be a lower-bound on the original μ_i strongly-convex function as $\tilde{f}_i(w) := f_i(w_i) + \langle \nabla f_i(w_i), w - w_i \rangle + \frac{\mu_i}{2} \|w - w_i\|^2$, then $\nabla \tilde{f}_i(w) = \nabla f_i(w_i) + \mu_i[w - w_i]$.
- Using FTL on \tilde{f}_k instead and using that $\sum_{i=1}^k \nabla \tilde{f}_i(w_{k+1}) = 0$ and $\sum_{i=1}^{k-1} \nabla \tilde{f}_i(w_k) = 0$,

$$\sum_{i=1}^k \nabla f_i(w_i) + w_{k+1} \left[\sum_{i=1}^k \mu_i \right] = \sum_{i=1}^k \mu_i w_i \quad ; \quad \sum_{i=1}^{k-1} \nabla f_i(w_i) + w_k \left[\sum_{i=1}^{k-1} \mu_i \right] = \sum_{i=1}^{k-1} \mu_i w_i$$

$$\nabla f_k(w_k) + (w_{k+1} - w_k) \left[\sum_{i=1}^k \mu_i \right] = 0 \implies w_{k+1} = w_k - \eta_k \nabla f_k(w_k). \quad (\text{where } \eta_k := 1/\sum_{i=1}^k \mu_i)$$

(Adding $\mu_k w_k$ to the second equation, and subtracting the two equations)

Hence, for the strongly-convex setting, running FTL on \tilde{f}_k recovers OGD on f_k .

Follow the Leader

Claim: If the convex set \mathcal{C} has a diameter D , for an arbitrary sequence of losses such that each f_k is μ_k strongly-convex (s.t. $\mu := \min_{k \in [T]} \mu_k > 0$), G -Lipschitz and differentiable, FTL with $w_1 \in \mathcal{C}$ has the following regret for all $u \in \mathcal{C}$,

$$R_T(u) \leq \frac{G^2}{2\mu} (1 + \log(T))$$


Hence, FTL achieves the same regret as OGD when the sequence of losses is strongly-convex and Lipschitz (we will prove this later)

- What about when the losses are convex but not strongly-convex?

Consider running FTL on the following problem. $\mathcal{C} = [-1, 1]$ and $f_k(w) = \langle z_k, w \rangle$ where

$$z_1 = -0.5; \quad z_k = 1 \quad \text{for } k = 2, 4, \dots; \quad z_k = -1 \quad \text{for } k = 3, 5, \dots$$

In round 1, FTL suffers $-0.5w_1$ cost and will compute $w_2 = 1$. It will suffer cost of 1 in round 2 and compute $w_3 = -1$. In round 3, it will thus suffer a cost of 1 and so on. Hence, FTL will suffer $O(T)$ regret if the losses are not strongly-convex.

-  Stéphane Ross, Geoffrey Gordon, and Drew Bagnell, *A reduction of imitation learning and structured prediction to no-regret online learning*, Proceedings of the fourteenth international conference on artificial intelligence and statistics, JMLR Workshop and Conference Proceedings, 2011, pp. 627–635.