Today

Multiplicative Weight
Update Method
(part II)

CSS.205.1
Toolkit in TCS
— Lecture #16
(12 Apr '21)
Instructor: Prahladh
Harsha

$\underline{MWUM_{\varepsilon}}$  (Parameter: $\eta \in (0, \tfrac{1}{2}]$)

① Initialize: $\forall i \in [n],\ \omega_i^{(1)} \leftarrow 1$

② For $t \leftarrow 1$ to $T$

(a) Choose the probability distribution
$$P^{(t)} = (P_1^{(t)}, \ldots, P_n^{(t)})$$

where $\quad P_i^{(t)} = \dfrac{\omega_i^{(t)}}{\Phi^{(t)}}$ & $\Phi^{(t)} = \sum_{i \in [n]} \omega_i^{(t)}$

(b) Observe the costs
$$M^{(t)} = (m_1^{(t)}, \ldots, m_n^{(t)}) \in [-1, 1]^n$$

(c) Sample according to prob dist $P^{(t)}$

(d) Update the weights
$$\omega_i^{(t+1)} \leftarrow \omega_i^{(t)} \cdot (1 - m_i^{(t)} \cdot \varepsilon)$$

Today:

Expected Performance of MWUM$_\varepsilon$

$\ell(t) :=$ Expected loss at time step $T$

$$= \sum_{c=1}^{n} m_c^{(t)} \cdot p_c^{(t)} = \langle M^{(t)}, p^{(t)} \rangle$$

$L(T) :=$ Expected loss upto time $T$

$$= \sum_{t=1}^{T} \ell(t) = \sum_{t=1}^{T} \langle M^{(t)}, p^{(t)} \rangle.$$

<u>Theorem</u>: Assuming all costs $\in [-1, 1]$

$\varepsilon \in (0, \frac{1}{2}]$, for all experts

$$L(T) \le \sum_{c=1}^{T} m_c^{(t)} + \varepsilon \sum_{c=1}^{T} |m_c^{(t)}| + \frac{\ln n}{\varepsilon}$$

<u>Remarks</u> ① If all $m_c^{(t)} \in [0, 1]$

$$L(T) \le (1+\varepsilon) \sum_{c=1}^{T} m_c^{(t)} + \frac{\ln n}{\varepsilon}.$$

② If you knew beforehand the #steps

can fix $\varepsilon$-learning rate appropriately

**Proof:** As in WM, will be using a potential function.

$$\Phi^{(t)} := \sum_{c \in [n]} \omega_c^{(t)}$$

$$\Phi^{(t+1)} = \sum_{c \in [n]} \omega_c^{(t+1)} = \sum_{c \in [n]} \omega_c^{(t)} \cdot \left(1 - \varepsilon m_c^{(t)}\right)$$

$$= \Phi^{(t)} - \varepsilon \Phi^{(t)} \sum_{c \in [n]} m_c^{(t)} \cdot P_c^{(t)}$$

$$= \Phi^{(t)} \left(1 - \varepsilon \langle M^{(t)}, p^{(t)} \rangle \right)$$

$$\leq \Phi^{(t)} \exp\left(-\varepsilon \langle M^{(t)}, p^{(t)} \rangle \right) \quad \Big/ \quad 1 - \varepsilon x \leq e^{-\varepsilon x}$$

Hence,

$$\Phi^{(T+1)} \leq \Phi^{(1)} \exp\left(-\varepsilon \sum_{t=1}^{T} \langle M^{(t)}, p^{(t)} \rangle \right)$$

$$= n \cdot \exp\left(-\varepsilon L(T)\right)$$

For any expert $i$

$$\Phi^{(t+1)} \geq \omega_i^{(t+1)} = \omega_i^{(t)}\left(1 - \varepsilon m_i^{(t)}\right)$$

$$\Phi^{(T+1)} = \omega_i^{(1)} \prod_{t=1}^{T} \left(1 - \varepsilon m_i^{(t)}\right)$$

$$\geq 1 \prod_{t:\, m_t^{(f)}\geq 0}(1-\varepsilon)^{m_t^{(f)}} \prod_{t:\, m_t^{(f)}<0}(1+\varepsilon)^{m_t^{(f)}} \quad \left| \begin{array}{l} 1-\varepsilon x \geq (1-\varepsilon)^x \\ \qquad \text{if } x\in[0,1] \\ 1-\varepsilon x \geq (1+\varepsilon)^x \text{ if} \\ \qquad x\in[-1,0] \end{array}\right.$$

$$(1-\varepsilon)^{\sum_{\geq 0} m_t^{(f)}} (1+\varepsilon)^{\sum_{<0} m_t^{(f)}} \leq \Phi^{(T+1)} \leq n\cdot\exp(-\varepsilon L(T))$$

$$\ln n - \varepsilon L(T) \geq \sum_{t:\,\geq 0} m_t^{(f)} \ln(1-\varepsilon)$$
$$+ \sum m_t^{(f)} \ln(1+\varepsilon)$$

$$L(T) \leq \frac{\ln n}{\varepsilon} + \frac{1}{\varepsilon}\left( \ln\left(\frac{1}{1-\varepsilon}\right)\sum_{\geq 0} m_t^{(f)} \right.$$
$$\left. - \ln(1+\varepsilon)\sum_{<0} m_t^{(f)} \right)$$

$$\leq \frac{\ln n}{\varepsilon} + \frac{1}{\varepsilon}\left( (\varepsilon+\varepsilon^2)\sum_{\geq 0} m_t^{(f)} \right.$$
$$\left. + (\varepsilon+\varepsilon^2)\sum_{<0} m_t^{(f)} \right) \quad \left| \begin{array}{l} x\in[0,\tfrac{1}{2}] \\ \ln\left(\frac{1}{1-x}\right) \leq x+x^2 \\ \ln(1+x) \geq x+x^2 \end{array}\right.$$

$$= \frac{\ln n}{\varepsilon} + \sum_t m_t^{(f)} + \varepsilon \sum_t |m_t^{(f)}| \qquad \boxtimes$$

Theorem is true for every expert:
(including the expert that makes the minimum loss)

_Cor:_ Let $P$ be any _fixed_ distribution on the $n$ experts.

$$L(T) \leq \sum_{t=1}^{T} \langle M^{(T)}, P \rangle + \varepsilon \sum_{t=1}^{T} \langle M^{(t)}, P \rangle + \frac{\ln n}{\varepsilon}$$

_Pf:_ Apply convex combination w.r.t $P$ of previous theorem.

## Hedge Algorithm

Very similar to the MWUM$_\varepsilon$

Weight Update Rule.

MWUM$_\varepsilon$ : $\omega_i^{(t+1)} \leftarrow \omega_i^{(t)} \left( 1 - \varepsilon \cdot m_i^{(t)} \right)$

Hedge$_\varepsilon$ : $\omega_i^{(t+1)} \leftarrow \omega_i^{(t)} \exp \left( - \varepsilon \cdot m_i^{(t)} \right)$

_Thm:_ (Hedge Alg). $i$- any expert

$$L(T) \leq \sum_{t=1}^{T} m_i^{(t)} + \varepsilon \sum_{t=1}^{T} \langle M^{(t)^2}, P^{(t)} \rangle + \frac{\ln n}{\varepsilon}$$

$$\left( \exp(-\varepsilon x) \le (-\varepsilon x + \varepsilon^2 x^2) \right)$$

for the settings we
have chosen

---

Alg: → Each step is producing a prob
dist $p^{(T)}$ on experts.

→ Modifies $p^{(t)}$ based on
experts loss at the previous
step.

$\overline{P} = \{$ feasible probability dist on
$n$ experts$\}$

$P$ = closed convex set of $[0,1]^n$

So, far $P$ = all possible feasible prob
distributions.

But there could be scenarios where
all prob. dist not feasible.

And the alg must act only
according to some $p \in P$.

Qn: Does same $MWUM_\varepsilon$ work?

Ans: Possibly not, as $P^{(T)} \notin \mathcal{P}$.
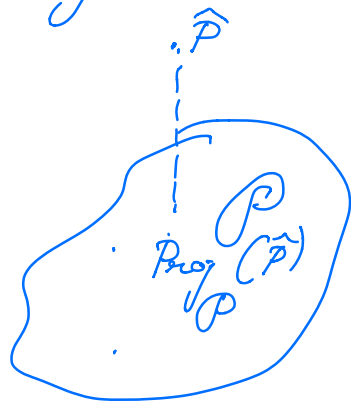
We will "project" $P^{(T)}$ into $\mathcal{P}$

Use relative entropy (Kullback -Leibler divergence) to do this projection

$P, Q$ – 2 distributions on $n$ experts

$P = (P_1, \ldots, P_n) ; Q = (q_1, \ldots, q_n)$

$$RE(P \| Q) = KL(P \| Q) = \sum_{i \in [n]} P_i \ln \frac{P_i}{q_i}$$
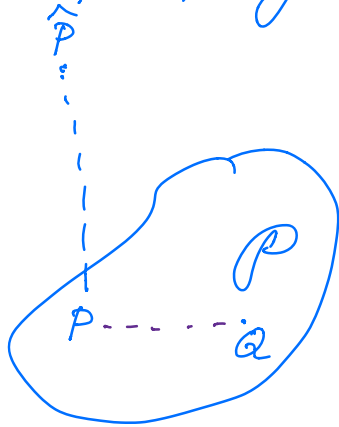
## Bregman Projection:



$\mathcal{P}$ – any closed convex set of prob dist.

$\hat{P}$ – any prob dist $(\hat{P} \notin \mathcal{P})$

$$\text{Proj}_{\mathcal{P}}(\hat{P}) = \underset{P \in \mathcal{P}}{\text{argmin}} \; RE(P \| \hat{P})$$

Convex function - can be obtained
using convex optimization

Property of Projection:



$$P \triangleq \underset{P}{Proj} \left( \hat{P} \right)$$

Q - any prob dist in P.

Generalized Pythagorean Inequality

$$RE(Q||P) + RE(P||\hat{P}) \leq RE(Q||\hat{P})$$

Hence $RE(Q||P) \leq RE(Q||\hat{P})$

(because RE is non-negative)

___

MWUM$_\varepsilon$ (in terms of relative entropy)

① Initialize: $P^{(1)} \leftarrow$ arbitrary distribution in P.

② For $t \leftarrow 1$ to $T$

(a) Observe costs $M^{(t)} = (m_1^{(t)} \ldots, m_n^{(t)})$

(b) Sample accg to $P^{(t)}$ & act appropriately.

(c). Update $P^{(t)}$ to $P^{(t+1)}$ as follows

(i) $P^{(t)} \longrightarrow \hat{P}^{(t+1)}$

$$\hat{P}_i^{(t+1)} \longleftarrow \frac{P_i^{(t)} \left(1 - \varepsilon \cdot m_i^{(t)}\right)}{\hat{\Phi}^{(t)}}$$

where $\hat{\Phi}^{(t)}$ — normalization const to make $\hat{P}^{(t+1)}$ into a dist

(New Step). (ii) $P^{(t+1)} = \text{Proj}_P \left(\hat{P}^{(t+1)}\right)$

_____

## Performance Analysis:

$P$ — be any distribution on $P$.

We will observe

$RE\left(P \| P^{(t)}\right)$ varies as the algorithm progresses

$$RE\left(P \| \hat{P}^{(t+1)}\right) - RE\left(P \| P^{(t)}\right)$$

$$= \sum_{c \in [n]} P_c \ln \frac{P_c^{(t)}}{\hat{P}_c^{(t+1)}}$$

$$= \sum_{c \in [n]} P_c \ln \frac{\bar{\Phi}^{(t)}}{1 - \varepsilon \, m_c^{(t)}}$$

$$\leq \sum_{c: \, m_c^{(t)} \geq 0} P_c \cdot m_c^{(t)} \ln\left(\frac{1}{1-\varepsilon}\right)$$

$\left| \begin{array}{l} 1 - \varepsilon x \geq (1-\varepsilon)^x \\ \qquad \text{if } x \in [0,1] \\ 1 - \varepsilon x \geq (1+\varepsilon)^x \text{ if} \\ \qquad x \in [-1, 0] \end{array} \right.$

$$+ \sum_{c: \, m_c^{(t)} < 0} P_c \cdot m_c^{(t)} \ln(1+\varepsilon)$$

$$+ \sum P_c \ln \bar{\Phi}^{(t)}.$$

$$= \varepsilon \left( \langle M^{(t)}, P \rangle + \varepsilon \langle M^{(t)}, P \rangle \right)$$

$$+ \ln \bar{\Phi}^{(t)}.$$

$\left| \begin{array}{l} x \in [0, \frac{1}{2}] \\ \ln\left(\frac{1}{1-x}\right) \leq x + x^2 \\ \ln(1+x) \geq x + x^2 \end{array} \right.$

$$\ln \bar{\Phi}^{(t)} = \ln \left[ \sum_{c \in [n]} P_c^{(t)} \left(1 - \varepsilon \cdot m_c^{(t)}\right) \right]$$

$$= \ln \left(1 - \varepsilon \sum P_c^{(t)} \cdot m_c^{(t)}\right)$$

$$= \ln \left(1 - \varepsilon \langle P^{(t)}, M^{(t)} \rangle\right)$$

$$\leq -\varepsilon \langle P^{(t)}, M^{(t)} \rangle$$

$$RE(P||\widehat{P}^{(t+1)}) - RE(P||P^{(t)})$$
$$\leq \varepsilon(\langle M^{(t)}, P \rangle + \varepsilon\langle |M^{(t)}|, P \rangle)$$
$$- \varepsilon \langle M^{(t)}, P^{(t)} \rangle \quad \cdots(\#)$$

$\langle M^{(t)}, P \rangle$ — loss of the fixed dist $P$

$\langle M^{(t)}, P^{(A)} \rangle$ — loss of alg at step $t$.

Since $P^{(t+1)} = \text{Proj}_{\mathcal{P}}(\widehat{P}^{(t+1)})$

$$RE(P||P^{(t+1)}) \leq RE(P||\widehat{P}^{(t+1)})$$

$(\#)$ implies

$$RE(P||P^{(t+1)}) - RE(P||P^{(t)})$$
$$\leq \varepsilon(\langle M^{(t)}, P \rangle + \varepsilon\langle |M^{(t)}|, P \rangle)$$
$$\neq \varepsilon \langle M^{(t)}, P^{(t)} \rangle$$

Rewriting: & dividing by $\varepsilon$
$$\langle M^{(t)}, P^{(t)} \rangle \leq \langle M^{(t)}, P \rangle + \varepsilon \langle |M^{(t)}|, P \rangle$$
$$+ \frac{RE(P||P^{(t)}) - RE(P||P^{(t+1)})}{\varepsilon}$$

Summing over $t \leftarrow 1$ to $T$

$$L(T) = \sum_{t=1}^{T} \langle M^{(t)}, P^{(t)} \rangle$$

$$\leq \sum_{t=1}^{T} \langle M^{(t)}, P \rangle + \varepsilon \sum_{t=1}^{T} \langle |M^{(t)}|, P \rangle$$

$$+ \frac{RE(P \| P^{(1)}) - RE(P \| P^{(T+1)})}{\varepsilon}$$

Thm: $MWUM_\varepsilon$ (w.r.t relative entropy) has the following performance.

Let $P^{(1)}$ - any starting dist & $\Big\}$ $\in \mathcal{P}$.

$P$ - any fixed dist

then
$$L(T) \leq \sum_{t=1}^{T} \langle M^{(t)}, P \rangle + \varepsilon \sum_{t=1}^{T} \langle |M^{(t)}|, P \rangle$$

$$+ \frac{RE(P \| P^{(1)})}{\varepsilon}$$