

MLT Homework 8

Ana Borovac
Jonas Haslbeck
Bas Haver

November 4, 2018

Question 1

The Aggregating Algorithm plays $w_1^k = 1/K$ and updates as

$$w_{t+1}^k = \frac{w_t^k e^{-l_t^k}}{\sum_{j=1}^K w_t^j e^{-l_t^j}}$$

Let us define the Kullback-Leibler divergence aka relative entropy (notion of distance between probability distributions) from $p \in \Delta_K$ to $q \in \Delta_K$ by

$$KL(p, q) = \sum_{k=1}^K p_k \ln \frac{p_k}{q_k}$$

Fix $w_t \in \Delta_K$ and $l_t \in \mathbb{R}^K$. Consider the minimisation problem

$$\min_{w \in \Delta_K} w^T l_t + KL(w, w_t) \tag{1}$$

Subquestion 1.1

Show that the minimiser of problem (1) is w_{t+1} .

Solution

We would like to show:

$$\min_{w \in \Delta_K} w^T l_t + KL(w, w_t) = w_{t+1}^T l_t; \quad w_{t+1}^k = \frac{w_t^k e^{-l_t^k}}{\sum_{j=1}^K w_t^j e^{-l_t^j}}$$

$$\begin{aligned}
w_2^k &= \frac{w_1^k e^{-l_1^k}}{\sum_{j=1}^K w_1^j e^{-l_1^j}} \\
&= \frac{\frac{1}{K} e^{-l_1^k}}{\sum_{j=1}^K \frac{1}{K} e^{-l_1^j}} \\
&= \frac{e^{-l_1^k}}{K \sum_{j=1}^K e^{-l_1^j}}
\end{aligned}$$

$$\begin{aligned}
\min_{w \in \Delta_K} w^T l_1 + \text{KL}(w, w_1) &= \min_{w \in \Delta_K} \sum_{k=1}^K w^k l_1^k + \sum_{k=1}^K w^k \ln \frac{w^k}{w_1^k} \\
&= \min_{w \in \Delta_K} \sum_{k=1}^K w^k (l_1^k + \ln K w^k) \\
&= \min_{w \in \Delta_K} e^{\sum_{k=1}^K w^k (l_1^k + \ln K w^k)} \\
&= \min_{w \in \Delta_K} \prod_{k=1}^K e^{w^k (l_1^k + \ln K w^k)} \\
&= \min_{w \in \Delta_K} \prod_{k=1}^K e^{w^k l_1^k} e^{\ln K w^k} \\
&= \min_{w \in \Delta_K} \prod_{k=1}^K K w^k e^{w^k l_1^k} \\
&= \max_{w \in \Delta_K} \prod_{k=1}^K \frac{1}{K} w^{-k} e^{-w^k l_1^k}
\end{aligned}$$

Let's calculate what we have on the left side:

$$\begin{aligned}
w^T l_t + \text{KL}(w, w_t) &= w^T l_t + \sum_{k=1}^K w^k \ln \frac{w^k}{w_t^k} \\
&= \sum_{k=1}^K w^k l_t^k + \sum_{k=1}^K w^k \ln \frac{w^k}{w_t^k} \\
&= \sum_{k=1}^K w^k l_t^k + \sum_{k=1}^K w^k \ln \frac{w^k}{\frac{w_{t-1}^k e^{-l_{t-1}^k}}{\sum_{j=1}^K w_{t-1}^j e^{-l_{t-1}^j}}} \\
&= \sum_{k=1}^K w^k l_t^k + \sum_{k=1}^K w^k \ln \frac{w^k \sum_{j=1}^K w_{t-1}^j e^{-l_{t-1}^j}}{w_{t-1}^k e^{-l_{t-1}^k}} \\
&= \sum_{k=1}^K w^k \left(l_t^k + \ln \frac{w^k \sum_{j=1}^K w_{t-1}^j e^{-l_{t-1}^j}}{w_{t-1}^k e^{-l_{t-1}^k}} \right)
\end{aligned}$$

Subquestion 1.2

Show that the value of problem (1) is the mix loss.

Solution

Question 2

We saw in the lecture that the Hedge algorithm (for the dot-loss game) with learning rate $\eta = \sqrt{\frac{8 \ln K}{T}}$ has regret after T rounds bounded by $\sqrt{T/2 \ln K}$. In practice, we may not know T in advance, or we may even desire an algorithm that has good guarantees for all T simultaneously, i.e. that keeps on operating forever.

Consider the following exponential (base 3) restarting schedule to accomplish this. We run Hedge for 1 round, with η tuned for 1 round. After that, we restart Hedge, and run it for 3 rounds with η tuned for 3 rounds. After that, we restart Hedge again for 9 rounds with η tuned for 9 rounds, and so on.

Prove that the overall accumulated regret of Hedge with this scheme is bounded above by a universal constant times $\sqrt{T \ln K}$. (Your argument should work for T that are not a power of 3).

Solution

Question 3

Consider the $K = 2$ expert version of the T -round dot loss game (Definition 2). In this exercise we will prove that the worst-case expected regret is at least of order \sqrt{T} . Consider an adversary that for each $t = 1, \dots, T$ assigns loss vector $l_t = (0, 1)$ or $l_t = (1, 0)$ i.i.d uniformly at random.

Subquestion 3.1

Show that the expected loss of any learner is $T/2$.

Solution

We calculate the dot loss as:

$$\sum_{k=1}^K w_t^k l_t^k$$

Where, in our case:

$$w_t \in \{(0, 0), (0, 1), (1, 0), (1, 1)\} \text{ and } l_t \in \{(0, 1), (1, 0)\}$$

In the table below we can see all the possible values of $L_t = \sum_{k=1}^2 w_t^k l_t^k$:

$\sum_{k=1}^2 w_t^k l_t^k$	$(0, 1)$	$(1, 0)$
$(0, 0)$	0	0
$(0, 1)$	1	0
$(0, 1)$	0	1
$(1, 1)$	1	1

From that it follows:

$$P(L = 0) = P(L = 1) = \frac{1}{2}$$

And we can conclude:

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T L_t \right] &= \sum_{t=1}^T \mathbb{E}[L_t] \\ &= \sum_{t=1}^T \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 1 \\ &= \frac{T}{2} \end{aligned}$$

Subquestion 3.2

Show that $2(1/2 - l_t^k)$ is Rademacher for each $k \in \{1, 2\}$.

Solution

Subquestion 3.3

Show that $\sum_{t=1}^T (1/2 - l_t^2) = -\sum_{t=1}^T (1/2 - l_t^1)$.

Solution

$$\begin{aligned}\sum_{t=1}^T (1/2 - l_t^2) &= -\sum_{t=1}^T (1/2 - l_t^1) \\ \sum_{t=1}^T (1/2 - l_t^2) + \sum_{t=1}^T (1/2 - l_t^1) &= 0 \\ \sum_{t=1}^T (1/2 - l_t^2 + 1/2 - l_t^1) &= 0 \\ \sum_{t=1}^T (1/2 + 1/2 - (l_t^1 + l_t^2)) &= 0 \\ \sum_{t=1}^T (1 - 1) &= 0 \\ 0 &= 0\end{aligned}$$

Subquestion 3.4

Argue that the expected loss of the best expert is bounded above by $\mathbb{E}[\min_k \sum_{t=1}^T l_t^k] \leq T/2 - c\sqrt{T}$ for some $c > 0$. You can use the following fact. Let X_1, \dots, X_T be i.i.d Rademacher random variables. Then

$$\mathbb{E} \left[\sum_{t=1}^T X_t \right] \in \left[\sqrt{\frac{2(T-1)}{\pi}}, \sqrt{\frac{2(T+1)}{\pi}} \right].$$

Solution