

A PROOFS

Proof of Proposition 1. By definition,

$$\sum_{t=1}^T \mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f, x_t) \leq \sum_{t=1}^T \mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_T) .$$

Peeling off the T -th expected loss, we have

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_T) &\leq \sum_{t=1}^{T-1} \mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) + \{\mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_T)\} \\ &\leq \sum_{t=1}^{T-1} \mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{T-1}) \end{aligned}$$

where we used the fact that q_T is an admissible algorithm for this relaxation, and thus the last inequality holds for any choice x_T of the opponent. Repeating the process, we obtain

$$\sum_{t=1}^T \mathbb{E}_{f_t \sim q_t} \ell(f_t, x_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f, x_t) \leq \mathbf{Rel}_T(\mathcal{F}) .$$

We remark that the left-hand side of this inequality is random, while the right-hand side is not. Since the inequality holds for any realization of the process, it also holds in expectation. The inequality

$$\mathcal{V}_T(\mathcal{F}) \leq \mathbf{Rel}_T(\mathcal{F})$$

holds by unwinding the value recursively and using admissibility of the relaxation. The high-probability bound is an immediate consequences of (6) and the Hoeffding-Azuma inequality for bounded martingales. The last statement is immediate. \square

Proof of Proposition 2. Denote $L_t(f) = \sum_{s=1}^t \ell(f, x_s)$. The first step of the proof is an application of the minimax theorem (we assume the necessary conditions hold):

$$\begin{aligned} &\inf_{q_t \in \Delta(\mathcal{F})} \sup_{x_t \sim q_t} \left\{ \mathbb{E} [\ell(f_t, x_t)] + \sup_{\mathbf{x}} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_t(f) \right] \right\} \\ &= \sup_{p_t \in \Delta(\mathcal{X})} \inf_{f_t \in \mathcal{F}} \left\{ \mathbb{E} [\ell(f_t, x_t)] + \mathbb{E} \sup_{x_t \sim p_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_t(f) \right] \right\} \end{aligned}$$

For any $p_t \in \Delta(\mathcal{X})$, the infimum over f_t of the above expression is equal to

$$\begin{aligned} &\mathbb{E} \sup_{x_t \sim p_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_{t-1}(f) + \inf_{f_t \in \mathcal{F}} \mathbb{E} [\ell(f_t, x_t)] - \ell(f, x_t) \right] \\ &\leq \mathbb{E} \sup_{x_t \sim p_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_{t-1}(f) + \mathbb{E} [\ell(f, x_t)] - \ell(f, x_t) \right] \\ &\leq \mathbb{E} \sup_{x_t, x'_t \sim p_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_{t-1}(f) + \ell(f, x'_t) - \ell(f, x_t) \right] \end{aligned}$$

We now argue that the independent x_t and x'_t have the same distribution p_t , and thus we can introduce a random sign ϵ_t . The above expression then equals to

$$\begin{aligned} &\mathbb{E} \mathbb{E} \sup_{x_t, x'_t \sim p_t} \mathbb{E}_{\epsilon_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_{t-1}(f) + \epsilon_t (\ell(f, x'_t) - \ell(f, x_t)) \right] \\ &\leq \sup_{x_t, x'_t \in \mathcal{X}} \mathbb{E} \sup_{\epsilon_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[2 \sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - L_{t-1}(f) + \epsilon_t (\ell(f, x'_t) - \ell(f, x_t)) \right] \end{aligned}$$

where we upper bounded the expectation by the supremum. Splitting the resulting expression into two parts, we arrive at the upper bound of

$$2 \sup_{x_t \in \mathcal{X}} \mathbb{E} \sup_{\epsilon_t} \mathbb{E}_{\epsilon_{t+1:T}} \sup_{f \in \mathcal{F}} \left[\sum_{s=t+1}^T \epsilon_s \ell(f, \mathbf{x}_{s-t}(\epsilon_{t+1:s-1})) - \frac{1}{2} L_{t-1}(f) + \epsilon_t \ell(f, x_t) \right] = \mathfrak{R}_T(\mathcal{F}|x_1, \dots, x_{t-1}) .$$

The last equality is easy to verify, as we are effectively adding a root x_t to the two subtrees, for $\epsilon_t = +1$ and $\epsilon_t = -1$, respectively.

One can see that the proof of admissibility corresponds to one step minimax swap and symmetrization in the proof of [15]. In contrast, in the latter paper, all T minimax swaps are performed at once, followed by T symmetrization steps. \square

Proof of Proposition 3. Let us first prove that the relaxation is admissible with the Exponential Weights algorithm as an admissible algorithm. Let $L_t(f) = \sum_{i=1}^t \ell(f, x_i)$. Let λ^* be the optimal value in the definition of $\mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{t-1})$. Then

$$\begin{aligned} & \inf_{q_t \in \Delta(\mathcal{F})} \sup_{x_t \in \mathcal{X}} \left\{ \mathbb{E}_{f \sim q_t} [\ell(f, x_t)] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \\ & \leq \inf_{q_t \in \Delta(\mathcal{F})} \sup_{x_t \in \mathcal{X}} \left\{ \mathbb{E}_{f \sim q_t} [\ell(f, x_t)] + \frac{1}{\lambda^*} \log \left(\sum_{f \in \mathcal{F}} \exp(-\lambda^* L_t(f)) \right) + 2\lambda^*(T-t) \right\} \end{aligned}$$

Let us upper bound the infimum by a particular choice of q which is the exponential weights distribution

$$q_t(f) = \exp(-\lambda^* L_{t-1}(f)) / Z_{t-1}$$

where $Z_{t-1} = \sum_{f \in \mathcal{F}} \exp(-\lambda^* L_{t-1}(f))$. By [6, Lemma A.1],

$$\begin{aligned} \frac{1}{\lambda^*} \log \left(\sum_{f \in \mathcal{F}} \exp(-\lambda^* L_t(f)) \right) &= \frac{1}{\lambda^*} \log (\mathbb{E}_{f \sim q_t} \exp(-\lambda^* \ell(f, x_t))) + \frac{1}{\lambda^*} \log Z_{t-1} \\ &\leq -\mathbb{E}_{f \sim q_t} \ell(f, x_t) + \frac{\lambda^*}{2} + \frac{1}{\lambda^*} \log Z_{t-1} \end{aligned}$$

Hence,

$$\begin{aligned} & \inf_{q_t \in \Delta(\mathcal{F})} \sup_{x_t \in \mathcal{X}} \left\{ \mathbb{E}_{f \sim q_t} [\ell(f, x_t)] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \\ & \leq \frac{1}{\lambda^*} \log \left(\sum_{f \in \mathcal{F}} \exp(-\lambda^* L_{t-1}(f)) \right) + 2\lambda^*(T-t+1) \\ & = \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{t-1}) \end{aligned}$$

by the optimality of λ^* . The bound can be improved by a factor of 2 for some loss functions, since it will disappear from the definition of sequential Rademacher complexity.

We conclude that the Exponential Weights algorithm is an admissible strategy for the relaxation (9). The final regret bound follows immediately from the bound on sequential Rademacher complexity (which, in this case, is simply the supremum of a martingale difference process indexed by N elements – see e.g. [15]).

Arriving at the relaxation We now show that the Exponential Weights relaxation arises naturally as an upper bound on sequential Rademacher complexity of a finite class. For any $\lambda > 0$,

$$\begin{aligned} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left\{ 2 \sum_{i=1}^{T-t} \epsilon_i \ell(f, \mathbf{x}_i(\epsilon)) - L_t(f) \right\} \right] &\leq \frac{1}{\lambda} \log \left(\mathbb{E} \left[\sup_{f \in \mathcal{F}} \exp \left(2\lambda \sum_{i=1}^{T-t} \epsilon_i \ell(f, \mathbf{x}_i(\epsilon)) - \lambda L_t(f) \right) \right] \right) \\ &\leq \frac{1}{\lambda} \log \left(\mathbb{E} \left[\sum_{f \in \mathcal{F}} \exp \left(2\lambda \sum_{i=1}^{T-t} \epsilon_i \ell(f, \mathbf{x}_i(\epsilon)) - \lambda L_t(f) \right) \right] \right) \\ &= \frac{1}{\lambda} \log \left(\sum_{f \in \mathcal{F}} \exp(-\lambda L_t(f)) \mathbb{E} \left[\prod_{i=1}^{T-t} \exp(2\lambda \epsilon_i \ell(f, \mathbf{x}_i(\epsilon))) \right] \right) \end{aligned}$$

Since, conditioned on $\epsilon_1, \dots, \epsilon_{i-1}$, the random variable $\epsilon_i \ell(f, \mathbf{x}_i(\epsilon))$ is subgaussian, we can upper bound the expected value of the product, peeling one random variable at a time from the end (see

[15] for the proof). We arrive at the upper bound

$$\begin{aligned}
& \frac{1}{\lambda} \log \left(\sum_{f \in \mathcal{F}} \exp(-\lambda L_t(f)) \times \exp \left(2\lambda^2 \max_{\epsilon_1, \dots, \epsilon_{T-t} \in \{\pm 1\}} \sum_{i=1}^{T-t} \ell(f, \mathbf{x}_i(\epsilon))^2 \right) \right) \\
& \leq \frac{1}{\lambda} \log \left(\sum_{f \in \mathcal{F}} \exp \left(-\lambda L_t(f) + 2\lambda^2 \max_{\epsilon_1, \dots, \epsilon_{T-t} \in \{\pm 1\}} \sum_{i=1}^{T-t} \ell(f, \mathbf{x}_i(\epsilon))^2 \right) \right) \\
& \leq \frac{1}{\lambda} \log \left(\sum_{f \in \mathcal{F}} \exp(-\lambda L_t(f)) \right) + 2\lambda \sup_{\mathbf{x}} \sup_{f \in \mathcal{F}} \max_{\epsilon_1, \dots, \epsilon_{T-t} \in \{\pm 1\}} \sum_{i=1}^{T-t} \ell(f, \mathbf{x}_i(\epsilon))^2
\end{aligned}$$

The last term, representing the “worst future”, is upper bounded by $2\lambda(T-t)$, assuming that the losses are bounded by 1. This removes the \mathbf{x} tree and leads to the relaxation (9) and a computationally tractable algorithm. \square

Proof of Proposition 4. The argument can be seen as a generalization of the Euclidean proof in [2] to general smooth norms. Let $\tilde{x}_{t-1} = \sum_{i=1}^{t-1} x_i$. The optimal algorithm for the relaxation (10) is

$$f_t^* = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \sup_{x_t \in \mathcal{X}} \left\{ \langle f, x_t \rangle + \sqrt{\|\tilde{x}_{t-1}\|^2 + \langle \nabla \|\tilde{x}_{t-1}\|^2, x_t \rangle} + C(T-t+1) \right\} \right\} \quad (27)$$

We shall show admissibility instead using

$$f_t = -\frac{\nabla \|\tilde{x}_{t-1}\|^2}{2\sqrt{\|\tilde{x}_{t-1}\|^2 + C(T-t+1)}}. \quad (28)$$

Plugging this choice into the admissibility condition (4), we get

$$\sup_{x_t \in \mathcal{X}} \left\{ -\frac{\langle \nabla \|\tilde{x}_{t-1}\|^2, x_t \rangle}{2\sqrt{A}} + \sqrt{A + \langle \nabla \|\tilde{x}_{t-1}\|^2, x_t \rangle} \right\}$$

where $A = \|\tilde{x}_{t-1}\|^2 + C(T-t+1)$. It can be easily verified that an expression of the form $-\frac{x}{2\sqrt{y}} + \sqrt{y+x}$ is maximized at $x=0$ for a positive y . Hence,

$$\begin{aligned}
& \inf_{f_t \in \mathcal{F}} \left\{ \sup_{x_t \in \mathcal{X}} \left\{ \langle f_t, x_t \rangle + (\|\tilde{x}_{t-1}\|^2 + \langle \nabla \|\tilde{x}_{t-1}\|^2, x_t \rangle + C(T-t+1))^{1/2} \right\} \right\} \leq (\|\tilde{x}_{t-1}\|^2 + C(T-t+1))^{1/2} \\
& \leq (\|\tilde{x}_{t-2}\|^2 + \langle \nabla \|\tilde{x}_{t-2}\|^2, x_{t-1} \rangle + C(T-t+2))^{1/2} = \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{t-1})
\end{aligned}$$

Hence, the choice (28) is an admissible algorithm for the relaxation (10). Evidently, the above proof of admissibility is very simple, but it might seem that we pulled the algorithm (28) out of a hat. We now show that in fact the choice of f_t above is the optimal choice f_t^* . The proof below is not required for admissibility, and we only include it for completeness. The proof uses the fact that for any norm $\|\cdot\|$,

$$\langle \nabla \frac{1}{2} \|x\|^2, x \rangle = \|x\|^2. \quad (29)$$

To prove this, observe that by convexity $\|0\| \geq \|x\| + \langle \nabla \|x\|, 0-x \rangle$ and $\|2x\| \geq \|x\| + \langle \nabla \|x\|, 2x-x \rangle$ implying $\langle \nabla \|x\|, x \rangle = \|x\|$. On the other hand, by the chain rule, $\nabla \frac{1}{2} \|x\|^2 = \|x\| \cdot \nabla \|x\|$, thus implying (29).

Let

$$K \triangleq \operatorname{Kernel}(\nabla \|\tilde{x}_{t-1}\|^2) = \{h : \langle \nabla \|\tilde{x}_{t-1}\|^2, h \rangle = 0\}, \quad K' \triangleq \operatorname{Kernel}(\tilde{x}_{t-1}) = \{h : \langle h, \tilde{x}_{t-1} \rangle = 0\}.$$

We first claim that x_t can always be written as $x_t = \beta \tilde{x}_{t-1} + \gamma y$ for some $y \in K$ and for some scalars β, γ . Indeed, suppose that $x_t = \beta \tilde{x}_{t-1} + \gamma y + z$ for some $y \in K$ and $z \notin K$. Then we can rewrite x_t as

$$x_t = (\beta + \delta) \tilde{x}_{t-1} + (\gamma y - \delta \tilde{x}_{t-1} + z)$$

where $\delta = \frac{\langle \nabla \|\tilde{x}_{t-1}\|^2, z \rangle}{2\|\tilde{x}_{t-1}\|^2}$. It is enough to check that $(\gamma y - \delta \tilde{x}_{t-1} + z) \in K$. Indeed, using (29),

$$\langle \nabla \|\tilde{x}_{t-1}\|^2, -\delta \tilde{x}_{t-1} + z \rangle = -2\delta \|\tilde{x}_{t-1}\|^2 + \langle \nabla \|\tilde{x}_{t-1}\|^2, z \rangle = 0.$$

An analogous proof shows that we may always decompose any f_t as $f_t = -\alpha \nabla \frac{1}{2} \|\tilde{x}_{t-1}\|^2 + g$ for some $g \in K'$ and a scalar α . Hence,

$$\begin{aligned} & \langle f_t, x_t \rangle + (\|\tilde{x}_{t-1}\|^2 + \langle \nabla \|\tilde{x}_{t-1}\|^2, x_t \rangle + C(T-t+1))^{1/2} \\ &= -\alpha \beta \|\tilde{x}_{t-1}\|^2 + \gamma \langle g, y \rangle + (\|\tilde{x}_{t-1}\|^2 + 2\beta \|\tilde{x}_{t-1}\|^2 + C(T-t+1))^{1/2} \end{aligned} \quad (30)$$

Given any $f_t = -\alpha \nabla \frac{1}{2} \|\tilde{x}_{t-1}\|^2 + g$, x_t can be picked with $y \in K$ that satisfies $\langle g, y \rangle \geq 0$. One can always do this because if for some y' , $\langle g, y' \rangle < 0$ by picking $y = -y'$ we can ensure that $\langle g, y \rangle \geq 0$. Hence the minimizer f_t must be once such that $f_t = -\alpha \nabla \frac{1}{2} \|\tilde{x}_{t-1}\|^2$ and thus $\langle g, y \rangle = 0$. Now, it must be that $\alpha \geq 0$ so that x_t either increases the first term or second term but not both. Hence we conclude that $f_t = -\alpha \nabla \frac{1}{2} \|\tilde{x}_{t-1}\|^2$ for some $\alpha \geq 0$. It remains to determine the optimal α . Given such an f_t , the sup over x_t can be written as supremum over β of a concave function, which gives rise to the derivative condition

$$-\alpha \|\tilde{x}_{t-1}\|^2 + \frac{\|\tilde{x}_{t-1}\|^2}{\sqrt{\|\tilde{x}_{t-1}\|^2 + 2\beta \|\tilde{x}_{t-1}\|^2 + C(T-t+1)}} = 0$$

Hence we can conclude that for any $f_t = -\alpha \nabla \frac{1}{2} \|\tilde{x}_{t-1}\|^2$,

$$\begin{aligned} & \sup_{x_t} \left\{ \langle f_t, x_t \rangle + (\|\tilde{x}_{t-1}\|^2 + \langle \nabla \|\tilde{x}_{t-1}\|^2, x_t \rangle + C(T-t+1))^{1/2} \right\} \\ &= \sup_{\beta} \left\{ -\alpha \beta \|\tilde{x}_{t-1}\|^2 + (\|\tilde{x}_{t-1}\|^2 + 2\beta \|\tilde{x}_{t-1}\|^2 + C(T-t+1))^{1/2} \right\} \\ &= \frac{\alpha}{2} (\|\tilde{x}_{t-1}\|^2 + C(T-t+1)) + \frac{1}{2\alpha} \end{aligned}$$

Hence optimizing the above over α we get

$$\alpha = \frac{1}{\sqrt{\|\tilde{x}_{t-1}\|^2 + C(T-t+1)}}.$$

Hence we can conclude that

$$f_t^* = -\frac{\nabla \frac{1}{2} \|\tilde{x}_{t-1}\|^2}{\sqrt{\|\tilde{x}_{t-1}\|^2 + C(T-t+1)}}$$

Plugging back the value of α , we find that $\beta = 0$. Hence we conclude that f_t defined in (28) in fact coincides with the optimal solution (27).

Arriving at the Relaxation The derivation of the relaxation is immediate:

$$\mathfrak{R}_T(\mathcal{F}|x_1, \dots, x_t) = \sup_{\mathbf{x}} \mathbb{E}_{\epsilon_{t+1:T}} \left\| \sum_{s=t+1}^T \epsilon_s \mathbf{x}_{s-t} (\epsilon_{t+1:s-1}) - \sum_{s=1}^t x_s \right\| \quad (31)$$

$$\leq \sup_{\mathbf{x}} \sqrt{\mathbb{E}_{\epsilon_{t+1:T}} \left\| \sum_{s=t+1}^T \epsilon_s \mathbf{x}_{s-t} (\epsilon_{t+1:s-1}) - \sum_{s=1}^t x_s \right\|^2} \quad (32)$$

$$\leq \sup_{\mathbf{x}} \sqrt{\left\| \sum_{s=1}^t x_s \right\|^2 + C \mathbb{E}_{\epsilon_{t+1:T}} \sum_{s=t+1}^T \|\epsilon_s \mathbf{x}_{s-t} (\epsilon_{t+1:s-1})\|^2} \quad (33)$$

where the last step is due to the smoothness of the norm and the fact that the first-order terms disappear under the expectation. The sum of norms is now upper bounded by $T-t$, thus removing the dependence on the “future”, and we arrive at

$$\sqrt{\left\| \sum_{s=1}^t x_s \right\|^2 + C(T-t)} \leq \sqrt{\left\| \sum_{s=1}^{t-1} x_s \right\|^2 + \left\langle \nabla \frac{1}{2} \left\| \sum_{s=1}^{t-1} x_s \right\|^2, x_t \right\rangle + C(T-t+1)}$$

as a relaxation on the sequential Rademacher complexity. \square

Proof of Proposition 5. We would like to show that, with the distribution q_t^* defined in (15),

$$\max_{y_t \in \{\pm 1\}} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t^*} |\hat{y}_t - y_t| + \mathbf{Rel}_T(\mathcal{F}|(x^t, y^t)) \right\} \leq \mathbf{Rel}_T(\mathcal{F}|(x^{t-1}, y^{t-1}))$$

for any $x_t \in \mathcal{X}$. Let $\sigma \in \{\pm 1\}^{t-1}$ and $\sigma_t \in \{\pm 1\}$. We have

$$\begin{aligned} & \mathbf{Rel}_T(\mathcal{F}|(x^t, y^t)) - 2\lambda(T-t) \\ &= \frac{1}{\lambda} \log \left(\sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda|\sigma_t - y_t|\} \right) \\ &\leq \frac{1}{\lambda} \log \left(\sum_{\sigma_t \in \{\pm 1\}} \exp\{-\lambda|\sigma_t - y_t|\} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \right) \end{aligned}$$

Just as in the proof of Proposition 3, we may think of the two choices σ_t as the two experts whose weighting q_t^* is given by the sum involving the Littlestone's dimension of subsets of \mathcal{F} . Introducing the normalization term, we arrive at the upper bound

$$\begin{aligned} & \frac{1}{\lambda} \log(\mathbb{E}_{\sigma_t \sim q_t^*} \exp\{-\lambda|\sigma_t - y_t|\}) + \frac{1}{\lambda} \log \left(\sum_{\sigma_t \in \{\pm 1\}} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \right) \\ &\leq -\mathbb{E}_{\sigma_t \sim q_t^*} |\sigma_t - y_t| + 2\lambda + \frac{1}{\lambda} \log \left(\sum_{\sigma_t \in \{\pm 1\}} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \right) \end{aligned}$$

The last step is due to Lemma A.1 in [6]. It remains to show that the log normalization term is upper bounded by the relaxation at the previous step:

$$\begin{aligned} & \frac{1}{\lambda} \log \left(\sum_{\sigma_t \in \{\pm 1\}} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \right) \\ &\leq \frac{1}{\lambda} \log \left(\sum_{\sigma \in \mathcal{F}|_{x^{t-1}}} \exp\{-\lambda L_{t-1}(\sigma)\} \sum_{\sigma_t \in \{\pm 1\}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \right) \\ &\leq \frac{1}{\lambda} \log \left(\sum_{\sigma \in \mathcal{F}|_{x^{t-1}}} \exp\{-\lambda L_{t-1}(\sigma)\} g(\text{Ldim}(\mathcal{F}_{t-1}(\sigma)), T-t+1) \right) \\ &= \mathbf{Rel}_T(\mathcal{F}|(x^{t-1}, y^{t-1})) \end{aligned}$$

To justify the last inequality, note that $\mathcal{F}_{t-1}(\sigma) = \mathcal{F}_t(\sigma, +1) \cup \mathcal{F}_t(\sigma, -1)$ and at most one of $\mathcal{F}_t(\sigma, +1)$ or $\mathcal{F}_t(\sigma, -1)$ can have Littlestone's dimension $\text{Ldim}(\mathcal{F}_{t-1}(\sigma))$. We now appeal to the recursion

$$g(d, T-t) + g(d-1, T-t) \leq g(d, T-t+1)$$

where $g(d, T-t)$ is the size of the zero cover for a class with Littlestone's dimension d on the worst-case tree of depth $T-t$ (see [15]). This completes the proof of admissibility.

Alternative Method Let us now derive the algorithm. Once again, consider the optimization problem

$$\max_{y_t \in \{\pm 1\}} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t^*} |\hat{y}_t - y_t| + \mathbf{Rel}_T(\mathcal{F}|(x^t, y^t)) \right\}$$

with the relaxation

$$\mathbf{Rel}_T(\mathcal{F}|(x^t, y^t)) = \frac{1}{\lambda} \log \left(\sum_{\sigma \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma)), T-t) \exp\{-\lambda L_t(\sigma)\} \right) + \frac{\lambda}{2}(T-t)$$

The maximum can be written explicitly, as in Section 3:

$$\begin{aligned} & \max \left\{ 1 - q_t^* + \frac{1}{\lambda} \log \left(\sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1 - \sigma_t)\} \right), \right. \\ & \quad \left. 1 + q_t^* + \frac{1}{\lambda} \log \left(\sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1 + \sigma_t)\} \right) \right\} \end{aligned}$$

where we have dropped the $\frac{\lambda}{2}(T-t)$ term from both sides. Equating the two values, we obtain

$$2q_t^* = \frac{1}{\lambda} \log \frac{\sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1-\sigma_t)\}}{\sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1+\sigma_t)\}}$$

The resulting value becomes

$$\begin{aligned} & 1 + \frac{\lambda}{2}(T-t) + \frac{1}{2\lambda} \log \left\{ \sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1-\sigma_t)\} \right\} \\ & \quad + \frac{1}{2\lambda} \log \left\{ \sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1+\sigma_t)\} \right\} \\ & = 1 + \frac{\lambda}{2}(T-t) + \frac{1}{\lambda} \mathbb{E}_\epsilon \log \left\{ \sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \exp\{-\lambda(1-\epsilon\sigma_t)\} \right\} \\ & \leq 1 + \frac{\lambda}{2}(T-t) + \frac{1}{\lambda} \log \left\{ \sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \mathbb{E}_\epsilon \exp\{-\lambda(1-\epsilon\sigma_t)\} \right\} \end{aligned}$$

for a Rademacher random variable $\epsilon \in \{\pm 1\}$. Now,

$$\mathbb{E}_\epsilon \exp\{-\lambda(1-\epsilon\sigma_t)\} = e^{-\lambda} \mathbb{E}_\epsilon e^{\lambda\epsilon\sigma_t} \leq e^{-\lambda} e^{\lambda^2/2}$$

Substituting this into the above expression, we obtain an upper bound of

$$\frac{\lambda}{2}(T-t+1) + \frac{1}{\lambda} \log \left\{ \sum_{(\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}_t(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \right\}$$

which completes the proof of admissibility using the same combinatorial argument as in the earlier part of the proof.

Arriving at the Relaxation Finally, we show that the relaxation we use arises naturally as an upper bound on the sequential Rademacher complexity. Fix a tree \mathbf{x} . Let $\sigma \in \{\pm 1\}^{t-1}$ be a sequence of signs. Observe that given history $x^t = (x_1, \dots, x_t)$, the signs $\epsilon \in \{\pm 1\}^{T-t}$, and a tree \mathbf{x} , the function class \mathcal{F} takes on only a finite number of possible values $(\sigma, \sigma_t, \omega)$ on $(x^t, \mathbf{x}(\epsilon))$. Here, $\mathbf{x}(\epsilon)$ denotes the sequences of values along the path ϵ . We have,

$$\begin{aligned} & \sup_{\mathbf{x}} \mathbb{E}_\epsilon \sup_{f \in \mathcal{F}} \left\{ 2 \sum_{i=1}^{T-t} \epsilon_i f(\mathbf{x}_i(\epsilon)) - \sum_{i=1}^t |f(x_i) - y_i| \right\} \\ & = \sup_{\mathbf{x}} \mathbb{E}_\epsilon \max_{\sigma_t \in \{\pm 1\}} \max_{(\sigma, \omega): (\sigma, \sigma_t, \omega) \in \mathcal{F}|_{(x^t, \mathbf{x}(\epsilon))}} \left\{ 2 \sum_{i=1}^{T-t} \epsilon_i \omega_i - \sum_{i=1}^t |\sigma_i - y_i| \right\} \\ & \leq \sup_{\mathbf{x}} \mathbb{E}_\epsilon \max_{\sigma_t \in \{\pm 1\}} \max_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} \max_{\mathbf{v} \in V(\mathcal{F}(\sigma, \sigma_t), \mathbf{x})} \left\{ 2 \sum_{i=1}^{T-t} \epsilon_i \mathbf{v}_i(\epsilon) - \sum_{i=1}^t |\sigma_i - y_i| \right\} \end{aligned}$$

where $\mathcal{F}|_{(x^t, \mathbf{x}(\epsilon))}$ is the projection of \mathcal{F} onto $(x^t, \mathbf{x}(\epsilon))$, $\mathcal{F}(\sigma, \sigma_t) = \{f \in \mathcal{F} : f(x^t) = (\sigma, \sigma_t)\}$, and $V(\mathcal{F}(\sigma, \sigma_t), \mathbf{x})$ is the zero-cover of the set $\mathcal{F}(\sigma, \sigma_t)$ on the tree \mathbf{x} . We then have the following relaxation:

$$\frac{1}{\lambda} \log \left(\sup_{\mathbf{x}} \mathbb{E}_\epsilon \sum_{\sigma_t \in \{\pm 1\}} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} \sum_{\mathbf{v} \in V(\mathcal{F}(\sigma, \sigma_t), \mathbf{x})} \exp \left\{ 2\lambda \sum_{i=1}^{T-t} \epsilon_i \mathbf{v}_i(\epsilon) - \lambda L_t(\sigma, \sigma_t) \right\} \right)$$

where $L_t(\sigma, \sigma_t) = \sum_{i=1}^t |\sigma_i - y_i|$. The latter quantity can be factorized:

$$\begin{aligned} & \frac{1}{\lambda} \log \left(\sup_{\mathbf{x}} \sum_{\sigma_t \in \{\pm 1\}} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} \exp\{-\lambda L_t(\sigma, \sigma_t)\} \mathbb{E}_\epsilon \sum_{\mathbf{v} \in V(\mathcal{F}(\sigma, \sigma_t), \mathbf{x})} \exp \left\{ 2\lambda \sum_{i=1}^{T-t} \epsilon_i \mathbf{v}_i(\epsilon) \right\} \right) \\ & \leq \frac{1}{\lambda} \log \left(\sup_{\mathbf{x}} \sum_{\sigma_t \in \{\pm 1\}} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} \exp\{-\lambda L_t(\sigma, \sigma_t)\} \text{card}(V(\mathcal{F}(\sigma, \sigma_t), \mathbf{x})) \exp\{2\lambda^2(T-t)\} \right) \\ & \leq \frac{1}{\lambda} \log \left(\sum_{\sigma_t \in \{\pm 1\}} \exp\{-\lambda |\sigma_t - y_t|\} \sum_{\sigma: (\sigma, \sigma_t) \in \mathcal{F}|_{x^t}} g(\text{Ldim}(\mathcal{F}(\sigma, \sigma_t)), T-t) \exp\{-\lambda L_{t-1}(\sigma)\} \right) + 2\lambda(T-t). \end{aligned}$$

This concludes the derivation of the relaxation. \square

Proof of Lemma 6. We first exhibit the proof for the convex loss case. To show admissibility using the particular randomized strategy q_t given in the lemma, we need to show that

$$\sup_{x_t} \left\{ \mathbb{E}_{f \sim q_t} [\ell(f, x_t)] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \leq \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{t-1})$$

The strategy q_t proposed by the lemma is such that we first draw $x_{t+1}, \dots, x_T \sim D$ and $\epsilon_{t+1}, \dots, \epsilon_T$ Rademacher random variables, and then based on this sample pick $f_t = f_t(x_{t+1:T}, \epsilon_{t+1:T})$ as in (17). Hence,

$$\begin{aligned} & \sup_{x_t} \left\{ \mathbb{E}_{f \sim q_t} [\ell(f, x_t)] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \\ &= \sup_{x_t} \left\{ \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \ell(f_t, x_t) + \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\} \\ &\leq \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \sup_{x_t} \left\{ \ell(f_t, x_t) + \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\} \end{aligned}$$

where $L_t(f) = \sum_{i=1}^t \ell(f, x_i)$. Observe that our strategy “matched the randomness” arising from the relaxation! Now, with f_t defined as

$$f_t = \operatorname{argmin}_{g \in \mathcal{F}} \sup_{x_t \in \mathcal{X}} \left\{ \ell(g, x_t) + \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\}$$

for any given $x_{t+1:T}, \epsilon_{t+1:T}$, we have

$$\sup_{x_t} \left\{ \ell(f_t, x_t) + \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\} = \inf_{g \in \mathcal{F}} \sup_{x_t} \left\{ \ell(g, x_t) + \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\}$$

We can conclude that for this choice of q_t ,

$$\begin{aligned} & \sup_{x_t} \left\{ \mathbb{E}_{f \sim q_t} [\ell(f, x_t)] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \leq \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \inf_{g \in \mathcal{F}} \sup_{x_t} \left\{ \ell(g, x_t) + \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\} \\ &= \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \inf_{g \in \mathcal{F}} \sup_{p_t \in \Delta(\mathcal{X})} \mathbb{E}_{x_t \sim p_t} \left[\ell(g, x_t) + \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right] \\ &= \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \sup_{p \in \Delta(\mathcal{X})} \inf_{g \in \mathcal{F}} \left\{ \mathbb{E}_{x_t \sim p} [\ell(g, x_t)] + \mathbb{E}_{x_t \sim p} \left[\sup_{f \in \mathcal{F}} C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_t(f) \right] \right\} \end{aligned}$$

In the last step we appealed to the minimax theorem which holds as loss is convex in g and \mathcal{F} is a compact convex set and the term in the expectation is linear in p_t , as it is an expectation. The last expression can be written as

$$\begin{aligned} & \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \sup_{p \in \Delta(\mathcal{X})} \mathbb{E}_{x_t \sim p} \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_{t-1}(f) + \inf_{g \in \mathcal{F}} \mathbb{E}_{x_t \sim p} [\ell(g, x_t)] - \ell(f, x_t) \right] \\ &\leq \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \sup_{p \in \Delta(\mathcal{X})} \mathbb{E}_{x_t \sim p} \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_{t-1}(f) + \mathbb{E}_{x_t \sim p} [\ell(f, x_t)] - \ell(f, x_t) \right] \\ &\leq \mathbb{E}_{\substack{\epsilon_{t+1:T} \\ x_{t+1:T}}} \mathbb{E}_{x_t \sim D} \mathbb{E}_{\epsilon_t} \sup_{f \in \mathcal{F}} \left[C \sum_{i=t+1}^T \epsilon_i \ell(f, x_i) - L_{t-1}(f) + C \epsilon_t \ell(f, x_t) \right] \\ &= \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{t-1}) \end{aligned}$$

Last inequality is by Assumption 1, using which we can replace a draw from supremum over distributions by a draw from the “equivalently bad” fixed distribution D by suffering an extra factor of C multiplied to that random instance.

The key step where we needed convexity was to use minimax theorem to swap infimum and supremum inside the expectation. In general the minimax theorem need not hold. In the non-convex scenario this is the reason we add the extra randomization through \hat{q}_t . The non-convex case has a similar proof except that we have expectation w.r.t. \hat{q}_t extra on each round which essentially convexifies our loss and thus allows us to appeal to the minimax theorem. \square

Proof of Lemma 7. Let $w \in \mathbb{R}^N$ be arbitrary. We need to show

$$\max_{x \in \{\pm 1\}^N} \mathbb{E}_\epsilon \max_{i \in [N]} |w_i + 2\epsilon x_i| \leq \mathbb{E}_{x \sim D} \mathbb{E}_\epsilon \max_{i \in [N]} |w_i + C\epsilon x_i| \quad (34)$$

Let $i^* = \operatorname{argmax}_i |w_i|$ and $j^* = \operatorname{argmax}_{i \neq i^*} |w_i|$ be the coordinates with largest and second-largest magnitude. If $|w_{i^*}| - |w_{j^*}| \geq 4$, the statement is immediate as the top coordinate stays at the top. It remains to consider the case when $|w_{i^*}| - |w_{j^*}| < 4$. In this case first note that,

$$\max_{x \in \{\pm 1\}^N} \mathbb{E}_\epsilon \max_{i \in [N]} |w_i + 2\epsilon x_i| \leq |w_{i^*}| + 2$$

On the other hand, since the distribution we consider is symmetric, with probability 1/2 its sign is negative and with remaining probability positive. Define $\sigma_{i^*} = \operatorname{sign}(x_{i^*})$, $\sigma_{j^*} = \operatorname{sign}(x_{j^*})$, $\tau_{i^*} = \operatorname{sign}(w_{i^*})$, and $\tau_{j^*} = \operatorname{sign}(w_{j^*})$. Since each coordinate is drawn i.i.d., using conditional expectations we have,

$$\begin{aligned} \mathbb{E}_{x, \epsilon} \max_i |w_i + C\epsilon x_i| &= \mathbb{E}_x \max_i |w_i + Cx_i| \\ &\geq \frac{\mathbb{E}_x [|w_{i^*} + Cx_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{2} + \frac{\mathbb{E}_x [|w_{j^*} + Cx_{j^*}| \mid \sigma_{i^*} \neq \tau_{i^*}, \sigma_{j^*} = \tau_{j^*}]}{4} + \frac{\mathbb{E} [|w_{i^*} + Cx_{i^*}| \mid \sigma_{i^*} \neq \tau_{i^*}, \sigma_{j^*} \neq \tau_{j^*}]}{4} \\ &\geq \frac{\mathbb{E}_x [|w_{i^*}| + C|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{2} + \frac{\mathbb{E}_x [|w_{j^*}| + C|x_{j^*}| \mid \sigma_{i^*} \neq \tau_{i^*}, \sigma_{j^*} = \tau_{j^*}]}{4} + \frac{\mathbb{E} [|w_{i^*}| - C|x_{i^*}| \mid \sigma_{i^*} \neq \tau_{i^*}, \sigma_{j^*} \neq \tau_{j^*}]}{4} \\ &= \frac{\mathbb{E} [|w_{i^*}| + C|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{2} + \frac{\mathbb{E} [|w_{j^*}| + C|x_{j^*}| \mid \sigma_{j^*} = \tau_{j^*}]}{4} + \frac{\mathbb{E} [|w_{i^*}| - C|x_{i^*}| \mid \sigma_{i^*} \neq \tau_{i^*}]}{4} \\ &= \frac{|w_{i^*}| + C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{2} + \frac{|w_{j^*}| + C\mathbb{E} [|x_{j^*}| \mid \sigma_{j^*} = \tau_{j^*}]}{4} + \frac{|w_{i^*}| - C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} \neq \tau_{i^*}]}{4} \\ &= \frac{2|w_{i^*}| + |w_{j^*}| + 3C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{4} + \frac{|w_{i^*}| - C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} \neq \tau_{i^*}]}{4} \\ &= \frac{3|w_{i^*}| + |w_{j^*}| + 2C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{4} \end{aligned}$$

Now since we are in the case when $|w_{i^*}| - |w_{j^*}| < 4$ we see that

$$\mathbb{E}_{x, \epsilon} \max_i |w_i + C\epsilon x_i| \geq \frac{3|w_{i^*}| + |w_{j^*}| + 2C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{4} \geq \frac{4|w_{i^*}| + 2C\mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{4} - 4$$

On the other hand, as we already argued,

$$\max_{x \in \{\pm 1\}^N} \mathbb{E}_\epsilon \max_{i \in [N]} |w_i + 2\epsilon x_i| \leq |w_{i^*}| + 2$$

Hence, as long as

$$\frac{C \mathbb{E} [|x_{i^*}| \mid \sigma_{i^*} = \tau_{i^*}]}{2} - 2 \geq 2$$

or, in other words, as long as

$$C \geq 6/\mathbb{E} [|x_i| \mid \operatorname{sign}(x_i) = \operatorname{sign}(w_i)] = 6/\mathbb{E} [|x|] ,$$

we have that

$$\max_{x \in \{\pm 1\}^N} \mathbb{E}_\epsilon \max_{i \in [N]} |w_i + 2\epsilon x_i| \leq \mathbb{E}_{x, \epsilon} \max_i |w_i + C\epsilon x_i| .$$

This concludes the proof. \square

Lemma 13. Consider the case when \mathcal{X} is the ℓ_∞^N ball and \mathcal{F} is the ℓ_1^N unit ball. Let $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \langle f, R \rangle$, then for any random vector R ,

$$\mathbb{E} \left[\sup_{x \in \mathcal{X}} \{ \langle f^*, x \rangle + \|R + x\|_\infty \} \right] \leq \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sup_x \{ \langle f, x \rangle + \|R + x\|_\infty \} \right] + 4 \mathbf{P}(\|R\|_\infty \leq 4)$$

Proof. Let $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \langle f, R \rangle$. We start by noting that for any $f' \in \mathcal{F}$,

$$\begin{aligned} \sup_{x \in \mathcal{X}} \{ \langle f', x \rangle + \|R + x\|_\infty \} &= \sup_{x \in \mathcal{X}} \left\{ \langle f', x \rangle + \sup_{f \in \mathcal{F}} \langle f, R + x \rangle \right\} \\ &= \sup_{f \in \mathcal{F}} \sup_{x \in \mathcal{X}} \{ \langle f', x \rangle + \langle f, R + x \rangle \} \\ &= \sup_{f \in \mathcal{F}} \left\{ \sup_{x \in \mathcal{X}} \langle f' + f, x \rangle + \langle f, R \rangle \right\} \\ &= \sup_{f \in \mathcal{F}} \{ \|f' + f\|_1 + \langle f, R \rangle \} \end{aligned}$$

Hence note that

$$\inf_{f' \in \mathcal{F}} \sup_{x \in \mathcal{X}} \{ \langle f', x \rangle + \|R + x\|_\infty \} = \inf_{f' \in \mathcal{F}} \sup_{f \in \mathcal{F}} \{ \|f' + f\|_1 + \langle f, R \rangle \} \quad (35)$$

$$\geq \inf_{f' \in \mathcal{F}} \{ \|f' - f^*\|_1 - \langle f^*, R \rangle \} \geq \inf_{f' \in \mathcal{F}} \{ \|f' - f^*\|_1 + \|R\|_\infty \} = \|R\|_\infty \quad (36)$$

On the other hand note that, f^* is the vertex of the ℓ_1 ball (any one which given by $\operatorname{argmin}_{i \in [d]} |R[i]|$ with sign opposite as sign of $R[i]$ on that vertex). Since the ℓ_1 ball is the convex hull of the $2d$ vertices, any vector $f \in \mathcal{F}$ can be written as $f = \alpha h - \beta f^*$ some $h \in \mathcal{F}$ such that $\|h\|_1 = 1$ and $\langle h, R \rangle = 0$ (which means that h is 0 on the maximal co-ordinate of R specified by f^*) and for some $\beta \in [-1, 1]$, $\alpha \in [0, 1]$ s.t. $\|\alpha h - \beta f^*\|_1 \leq 1$. Further note that the constraint on α, β imposed by requiring that $\|\alpha h - \beta f^*\|_1 \leq 1$ can be written as $\alpha + |\beta| \leq 1$. Hence,

$$\begin{aligned} \sup_{x \in \mathcal{X}} \{ \langle f^*, x \rangle + \|R + x\|_\infty \} &= \sup_{f \in \mathcal{F}} \{ \|f^* + f\|_1 + \langle f, R \rangle \} \\ &= \sup_{\alpha \in [0, 1]} \sup_{h \perp f^*, \|h\|_1 = 1} \sup_{\beta \in [-1, 1], \|\alpha h - \beta f^*\|_1 \leq 1} \{ \|(1 - \beta)f^* + \alpha h\|_1 + \beta \langle f^*, R \rangle + \alpha \langle h, R \rangle \} \\ &= \sup_{\alpha \in [0, 1]} \sup_{h \perp f^*, \|h\|_1 = 1} \sup_{\beta \in [-1, 1], \|\alpha h - \beta f^*\|_1 \leq 1} \{ |1 - \beta| \|f^*\|_1 + \alpha \|h\|_1 + \beta \|R\|_\infty \} \\ &= \sup_{\alpha \in [0, 1]} \sup_{\beta \in [-1, 1]: |\beta| + \alpha \leq 1} \{ |1 - \beta| + \alpha + \beta \|R\|_\infty \} \\ &\leq \sup_{\beta \in [-1, 1]} \{ |1 - \beta| + 1 - |\beta| + \beta \|R\|_\infty \} \\ &\leq \sup_{\beta \in [-1, 1]} \{ 2|1 - \beta| + \beta \|R\|_\infty \} \\ &= \sup_{\beta \in [-1, 1]} \{ 2|1 - \beta| + \beta \|R\|_\infty \} \\ &= \max \{ \|R\|_\infty, 4 - \|R\|_\infty \} \\ &\leq \|R\|_\infty + 4 \mathbf{1}_{\{\|R\|_\infty \leq 4\}} \end{aligned}$$

Hence combining with equation 35 we can conclude that

$$\begin{aligned} \mathbb{E} \left[\sup_x \{ \langle f^*, x \rangle + \|R + x\|_\infty \} \right] &\leq \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sup_x \{ \langle f, x \rangle + \|R + x\|_\infty \} \right] + 4 \mathbb{E} [\mathbf{1}_{\{\|R\|_\infty \leq 4\}}] \\ &= \mathbb{E} \left[\inf_{f \in \mathcal{F}} \sup_x \{ \langle f, x \rangle + \|R + x\|_\infty \} \right] + 4 \mathbf{P}(\|R\|_\infty \leq 4) \end{aligned}$$

□

Proof of Lemma 8. On any round t , the algorithm draws $\epsilon_{t+1}, \dots, \epsilon_T$ and $x_{t+1}, \dots, x_T \sim D^N$ and plays

$$f_t = \operatorname{argmin}_{f \in \mathcal{F}} \left\langle f, \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i \right\rangle$$

We shall show that this randomized algorithm is (almost) admissible w.r.t. the relaxation (with some small additional term at each step). We define the relaxation as

$$\mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) = \mathbb{E}_{x_{t+1}, \dots, x_T \sim D} \left[\left\| \sum_{i=1}^t x_i - C \sum_{i=t+1}^T x_i \right\|_{\infty} \right]$$

Proceeding just as in the proof of Lemma 6 note that, for our randomized strategy,

$$\begin{aligned} & \sup_x \left\{ \mathbb{E}_{f \sim q_t} [\langle f, x \rangle] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \\ &= \sup_x \left\{ \mathbb{E}_{x_{t+1}, \dots, x_T \sim D^N} [\langle f_t, x \rangle] + \mathbb{E}_{x_{t+1}, \dots, x_T \sim D^N} \left[\left\| \sum_{i=1}^{t-1} x_i + x - C \sum_{i=t+1}^T x_i \right\|_{\infty} \right] \right\} \\ &\leq \mathbb{E}_{x_{t+1}, \dots, x_T \sim D^N} \left[\sup_x \left\{ \langle f_t, x \rangle + \left\| \sum_{i=1}^{t-1} x_i + x - C \sum_{i=t+1}^T x_i \right\|_{\infty} \right\} \right] \end{aligned} \quad (37)$$

In view of Lemma 13 (with $R = \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T \epsilon_i x_i$) we conclude that

$$\begin{aligned} & \mathbb{E}_{x_{t+1}, \dots, x_T} \left[\sup_{x \in \mathcal{X}} \left\{ \langle f_t, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\} \right] \\ &\leq \mathbb{E}_{x_{t+1}, \dots, x_T} \left[\inf_{f \in \mathcal{F}} \sup_x \left\{ \langle f, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\} \right] + 4 \mathbf{P} \left(\left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i \right\|_{\infty} \leq 4 \right) \\ &= \mathbb{E}_{x_{t+1}, \dots, x_T} \left[\sup_x \left\{ \langle f_t^*, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\} \right] + 4 \mathbf{P} \left(\left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i \right\|_{\infty} \leq 4 \right) \end{aligned}$$

where

$$f_t^* = \operatorname{argmin}_{f \in \mathcal{F}} \sup_x \left\{ \langle f, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\}$$

Combining with Equation (37) we conclude that

$$\begin{aligned} & \sup_x \left\{ \mathbb{E}_{f \sim q_t} [\langle f, x \rangle] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \\ &\leq \mathbb{E}_{x_{t+1}, \dots, x_T} \left[\sup_x \left\{ \langle f_t^*, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\} \right] + 4 \mathbf{P} \left(\left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i \right\|_{\infty} \leq 4 \right) \end{aligned}$$

Now, since

$$4 \mathbf{P} \left(\left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i \right\|_{\infty} \leq 4 \right) \leq 4 \mathbf{P} \left(C \left\| \sum_{i=t+1}^T x_i \right\|_{\infty} \leq 4 \right) \leq 4 \mathbf{P}_{y_{t+1}, \dots, y_T \sim D} \left(C \left| \sum_{i=t+1}^T y_i \right| \leq 4 \right)$$

we have

$$\sup_x \left\{ \mathbb{E}_{f \sim q_t} [\langle f, x \rangle] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \quad (38)$$

$$\leq \mathbb{E}_{x_{t+1}, \dots, x_T} \left[\sup_x \left\{ \langle f_t^*, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\} \right] + 4 \mathbf{P}_{y_{t+1}, \dots, y_T \sim D} \left(C \left| \sum_{i=t+1}^T y_i \right| \leq 4 \right) \quad (39)$$

In view of Lemma 7, Assumption 2 is satisfied by D^N with constant C . Further in the proof of Lemma 6 we already showed that whenever Assumption 2 is satisfied, the randomized strategy specified by f_t^* is admissible. More specifically we showed that

$$\mathbb{E}_{x_{t+1}, \dots, x_T} \left[\sup_x \left\{ \langle f_t^*, x \rangle + \left\| \sum_{i=1}^{t-1} x_i - C \sum_{i=t+1}^T x_i + x \right\|_{\infty} \right\} \right] \leq \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_{t-1})$$

and so using this in Equation (38) we conclude that for the randomized strategy in the statement of the lemma,

$$\begin{aligned} & \sup_x \left\{ \mathbb{E}_{f \sim q_t} [\langle f, x \rangle] + \mathbf{Rel}_T(\mathcal{F}|x_1, \dots, x_t) \right\} \\ & \leq \mathbf{Rel}_T(F|x_1, \dots, x_{t-1}) + 4 \mathbf{P}_{y_{t+1}, \dots, y_T \sim D} \left(C \left| \sum_{i=t+1}^T y_i \right| \leq 4 \right) \end{aligned}$$

Or in other words the randomized strategy proposed is admissible with an additional additive factor of $4 \mathbf{P}_{y_{t+1}, \dots, y_T \sim D} (C |\sum_{i=t+1}^T y_i| \leq 4)$ at each time step t . Hence by Proposition 1 we have that for the randomized algorithm specified in the lemma,

$$\begin{aligned} \mathbb{E}[\mathbf{Reg}_T] & \leq \mathbf{Rel}_T(F) + 4 \sum_{t=1}^T \mathbf{P}_{y_{t+1}, \dots, y_T \sim D} \left(C \left| \sum_{i=t+1}^T y_i \right| \leq 4 \right) \\ & = C \mathbb{E}_{x_1, \dots, x_T \sim D^N} \left[\left\| \sum_{t=1}^T x_t \right\|_\infty \right] + 4 \sum_{t=1}^T \mathbf{P}_{y_{t+1}, \dots, y_T \sim D} \left(C \left| \sum_{i=t+1}^T y_i \right| \leq 4 \right) \end{aligned}$$

This concludes the proof. \square

Proof of Lemma 9. Instead of using $C = 4\sqrt{2}$ and drawing uniformly from surface of unit sphere we can equivalently think of the constant as being 1 and drawing uniformly from surface of sphere of radius $4\sqrt{2}$. Let $\|\cdot\|$ stand for the Euclidean norm. To prove (19), first observe that

$$\sup_{p \in \Delta(\mathcal{X})} \mathbb{E}_{x_t \sim p} \left\| w + \mathbb{E}_{x \sim p} [x] - x_t \right\| \leq \sup_{x \in \mathcal{X}} \mathbb{E}_\epsilon \|w + 2\epsilon x\| \quad (40)$$

for any $w \in B$. Further, using Jensen's inequality

$$\sup_{x \in \mathcal{X}} \mathbb{E}_\epsilon \|w + 2\epsilon x\| \leq \sup_{x \in \mathcal{X}} \sqrt{\mathbb{E}_\epsilon \|w + 2\epsilon x\|^2} \leq \sup_{x \in \mathcal{X}} \sqrt{\|w\|^2 + \mathbb{E}_\epsilon \|2\epsilon x\|^2} = \sqrt{\|w\|^2 + 4}$$

To prove the lemma, it is then enough to show that for $r = 4\sqrt{2}$

$$\mathbb{E}_{x \sim D} \|w + rx\| \geq \sqrt{\|w\|^2 + 4} \quad (41)$$

for any w , where we omitted ϵ since D is symmetric. This fact can be proved with the following geometric argument.

We define quadruplets $(w + z_1, w + z_2, w - z_1, w - z_2)$ of points on the sphere of radius r . Each quadruplets will have the property that

$$\frac{\|w + z_1\| + \|w + z_2\| + \|w - z_1\| + \|w - z_2\|}{4} \geq \sqrt{\|w\|^2 + 4} \quad (42)$$

for any w . We then argue that the uniform distribution can be decomposed into these quadruplets such that each point on the sphere occurs in only one quadruplet (except for a measure zero set when z_1 is aligned with $-w$), thus concluding that (41) holds true.

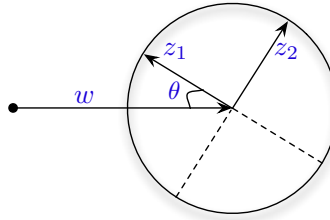


Figure 1: The two-dimensional construction for the proof of Lemma 9.

Pick any direction w^\perp perpendicular to w . A quadruplet is defined by perpendicular vectors z_1 and z_2 which have length r and which lie in the plane spanned by w, w^\perp . Let θ be the angle between $-w$ and z_1 . Since we are now dealing with a two dimensional plane spanned by w and w^\perp , we may as well assume that w is aligned with the positive x -axis, as in Figure 1. We write w for $\|w\|$. The coordinates of the quadruplet are

$$(w-r\cos(\theta), r\sin(\theta)), (w+r\cos(\theta), -r\sin(\theta)), (w+r\sin(\theta), r\cos(\theta)), (w-r\sin(\theta), -r\cos(\theta))$$

For brevity, let $s = \sin(\theta), c = \cos(\theta)$. The desired inequality (42) then reads

$$\sqrt{w^2 - 8wc + r^2} + \sqrt{w^2 + 8wc + r^2} + \sqrt{w^2 + 8ws + r^2} + \sqrt{w^2 - 8ws + r^2} \geq 4\sqrt{w^2 + 4}$$

To prove that this inequality holds, we square both sides, keeping in mind that the terms are non-negative. The sum of four squares on the left hand side gives $4w^2 + 4r^2$. For the six cross terms, we can pass to a lower bound by replacing r^2 in each square root by r^2c^2 or r^2s^2 , whichever completes the square. Then observe that

$$|w + rs| \cdot |w - rs| + |w + rc| \cdot |w - rc| = 2w^2 - r^2$$

while the other four cross terms

$$(|w + rs| \cdot |w - rc| + |w + rs| \cdot |w + rc|) + (|w - rs| \cdot |w + rc| + |w - rs| \cdot |w - rc|) \geq |w + rs| \cdot 2w + |w - rs| \cdot 2w \geq 4w^2$$

Doubling the cross terms gives a contribution of $2(6w^2 - r^2)$, while the sum of squares yielded $4w^2 + 4r^2$. The desired inequality is satisfied as long as $16w^2 + 2r^2 \geq 16(w^2 + 4)$, or $r \geq 4\sqrt{2}$. \square

Proof of Lemma 10. By Lemma 9, Assumption 2 is satisfied by distribution D with constant $C = 4\sqrt{2}$. Hence by Lemma 7 we can conclude that for the randomized algorithm which at round t freshly draws $x_{t+1}, \dots, x_T \sim D$ and picks

$$f_t^* = \operatorname{argmin}_{f \in \mathcal{F}} \sup_{x \in \mathcal{X}} \left\{ \langle f, x \rangle + \left\| -\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T x_i - x \right\|_2 \right\}$$

(we dropped the ϵ 's as the distribution is symmetric to start with) the expected regret is bounded as

$$\mathbb{E}[\mathbf{Reg}_T] \leq 4\sqrt{2} \mathbb{E}_{x_1, \dots, x_T \sim D} \left[\left\| \sum_{t=1}^T x_t \right\|_2 \right] \leq 4\sqrt{2T}$$

We claim that the strategy specified in the lemma that chooses

$$f_t = \frac{-\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T x_i}{\sqrt{\left\| -\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T \epsilon_i x_i \right\|_2^2 + 1}}$$

is the same as choosing f_t^* . To see this let us start by defining

$$\bar{x}_t = -\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T x_i$$

Now note that

$$\begin{aligned} f_t^* &= \operatorname{argmin}_{f \in \mathcal{F}} \sup_{x \in \mathcal{X}} \left\{ \langle f, x \rangle + \left\| -\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T x_i - x \right\|_2 \right\} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sup_{x \in \mathcal{X}} \{ \langle f, x \rangle + \|\bar{x}_t - x\|_2 \} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sup_{x \in \mathcal{X}} \left\{ \langle f, x \rangle + \sqrt{\|\bar{x}_t - x\|_2^2} \right\} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sup_{x: \|x\|_2 \leq 1} \left\{ \langle f, x \rangle + \sqrt{\|\bar{x}_t\|_2^2 - 2\langle \bar{x}_t, x \rangle + \|x\|_2^2} \right\} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sup_{x: \|x\|_2 = 1} \left\{ \langle f, x \rangle + \sqrt{\|\bar{x}_t\|_2^2 - 2\langle \bar{x}_t, x \rangle + 1} \right\} \end{aligned}$$

However this argmin calculation is identical to the one in the proof of Proposition 4 (with $C = 1$ and $T - t = 0$) and the solution is given by

$$f_t^* = f_t = \frac{-\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T x_i}{\sqrt{\left\| -\sum_{i=1}^{t-1} x_i + 4\sqrt{2} \sum_{i=t+1}^T \epsilon_i x_i \right\|_2^2 + 1}}$$

Thus we conclude the proof. \square

Proof of Lemma 11. We shall start by showing that the relaxation is admissible for the game where we pick prediction \hat{y}_t and the adversary then directly picks the gradient $\partial\ell(\hat{y}_t, y_t)$. To this end note that

$$\begin{aligned} & \inf_{\hat{y}_t} \sup_{\partial\ell(\hat{y}_t, y_t)} \{ \partial\ell(\hat{y}_t, y_t) \cdot \hat{y}_t + \mathbf{Rel}_T(\mathcal{F} | \partial\ell(\hat{y}_1, y_1), \dots, \partial\ell(\hat{y}_t, y_t)) \} \\ &= \inf_{\hat{y}_t} \sup_{\partial\ell(\hat{y}_t, y_t)} \left\{ \partial\ell(\hat{y}_t, y_t) \cdot \hat{y}_t + \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} 2L \sum_{i=t+1}^T \epsilon_i f[t] - \sum_{i=1}^t \partial\ell(\hat{y}_i, y_i) \cdot f[i] \right] \right\} \\ &\leq \inf_{\hat{y}_t} \sup_{r_t \in [-L, L]} \left\{ r_t \cdot \hat{y}_t + \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} 2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) - r_t \cdot f[t] \right] \right\} \end{aligned}$$

Let us use the notation $L_{t-1}(f) = \sum_{i=1}^{t-1} \partial\ell(\hat{y}_i, y_i) \cdot f[i]$ for the present proof. The supremum over $r_t \in [-L, L]$ is achieved at the endpoints since the expression is convex in r_t . Therefore, the last expression is equal to

$$\begin{aligned} & \inf_{\hat{y}_t} \sup_{r_t \in \{-L, L\}} \left\{ r_t \cdot \hat{y}_t + \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) - r_t \cdot f[t] \right] \right\} \\ &= \inf_{\hat{y}_t} \sup_{p_t \in \Delta(\{-L, L\})} \mathbb{E} \left[r_t \cdot \hat{y}_t + \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) - r_t \cdot f[t] \right] \right] \\ &= \sup_{p_t \in \Delta(\{-L, L\})} \inf_{\hat{y}_t} \mathbb{E} \left[r_t \cdot \hat{y}_t + \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) - r_t \cdot f[t] \right] \right] \end{aligned}$$

where the last step is due to the minimax theorem. The last quantity is equal to

$$\begin{aligned} & \sup_{p_t \in \Delta(\{-L, L\})} \mathbb{E} \left[\mathbb{E}_{\epsilon} \left[\inf_{r_t \sim p_t} \mathbb{E} [r_t] \cdot \hat{y}_t + \sup_{f \in \mathcal{F}} \left(2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) - r_t \cdot f[t] \right) \right] \right] \\ &\leq \sup_{p_t \in \Delta(\{-L, L\})} \mathbb{E} \left[\mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \left(2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) + \left(\mathbb{E}_{r_t \sim p_t} [r_t] - r_t \right) \cdot f[t] \right) \right] \right] \\ &\leq \sup_{p_t \in \Delta(\{-L, L\})} \mathbb{E} \left[\mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) + (r'_t - r_t) \cdot f[t] \right] \right] \\ &= \sup_{p_t \in \Delta(\{-L, L\})} \mathbb{E} \left[\mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) + \epsilon_t (r'_t - r_t) \cdot f[t] \right] \right] \end{aligned}$$

By passing to the worst-case choice of r_t, r'_t (which is achieved at the endpoints because of convexity), we obtain a further upper bound

$$\begin{aligned} & \sup_{r_t, r'_t \in \{L, -L\}} \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) + \epsilon_t (r'_t - r_t) \cdot f[t] \right] \\ &\leq \sup_{r_t \in \{L, -L\}} \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t+1}^T \epsilon_i f[t] - L_{t-1}(f) + 2\epsilon_t r_t \cdot f[t] \right] \\ &= \sup_{r_t \in \{L, -L\}} \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \left[2L \sum_{i=t}^T \epsilon_i f[t] - L_{t-1}(f) \right] \\ &= \mathbf{Rel}_T(\mathcal{F} | \partial\ell(\hat{y}_1, y_1), \dots, \partial\ell(\hat{y}_{t-1}, y_{t-1})) \end{aligned}$$

Thus we see that the relaxation is admissible. Now the corresponding prediction is given by

$$\begin{aligned}\hat{y}_t &= \operatorname{argmin}_{\hat{y}} \sup_{r_t \in [-L, L]} \left\{ r_t \hat{y} + \mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - \sum_{i=1}^{t-1} \partial \ell(\hat{y}_i, y_i) f[i] - r_t f[t] \right\} \right] \right\} \\ &= \operatorname{argmin}_{\hat{y}} \sup_{r_t \in [-L, L]} \left\{ r_t \hat{y} + \mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - \sum_{i=1}^{t-1} \partial \ell(\hat{y}_i, y_i) f[i] - r_t f[t] \right\} \right] \right\} \\ &= \operatorname{argmin}_{\hat{y}} \sup_{r_t \in \{-L, L\}} \left\{ r_t \hat{y} + \mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - \sum_{i=1}^{t-1} \partial \ell(\hat{y}_i, y_i) f[i] - r_t f[t] \right\} \right] \right\}\end{aligned}$$

The last step holds because of convexity of the term inside the supremum over r_t is convex in r_t and so the supremum is attained at the endpoints of the interval. The \hat{y}_t above is attained when both terms of the supremum are equalized, that is for \hat{y}_t is the prediction that satisfies :

$$\hat{y}_t = \mathbb{E} \left[\sup_{\epsilon} \left\{ \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2L} \sum_{i=1}^{t-1} \partial \ell(\hat{y}_i, y_i) f[i] + \frac{1}{2} f[t] \right\} - \sup_{f \in \mathcal{F}} \left\{ \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2L} \sum_{i=1}^{t-1} \partial \ell(\hat{y}_i, y_i) f[i] - \frac{1}{2} f[t] \right\} \right]$$

Finally since the relaxation is admissible we can conclude that the regret of the algorithm is bounded as

$$\mathbf{Reg}_T \leq \mathbf{Rel}_T(\mathcal{F}) = 2L \mathbb{E} \left[\sup_{\epsilon} \sum_{t=1}^T \epsilon_t f[t] \right].$$

This concludes the proof. \square

Proof of Lemma 12. The proof is similar to that of Lemma 11, with a few more twists. We want to establish admissibility of the relaxation given in (21) w.r.t. the randomized strategy q_t we provided. To this end note that

$$\begin{aligned}& \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_t(f) \right\} \right] \right\} \\ &= \sup_{y_t} \left\{ \mathbb{E}_{\epsilon} [\ell(\hat{y}_t(\epsilon), y_t)] + \mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_t(f) \right\} \right] \right\} \\ &\leq \mathbb{E} \left[\sup_{\epsilon} \left\{ \ell(\hat{y}_t(\epsilon), y_t) + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_t(f) \right\} \right\} \right]\end{aligned}$$

by Jensen's inequality, with the usual notation $L_t(f) = \sum_{i=1}^t \ell(f[i], y_i)$. Further, by convexity of the loss, we may pass to the upper bound

$$\begin{aligned}& \mathbb{E} \left[\sup_{\epsilon} \left\{ \partial \ell(\hat{y}_t(\epsilon), y_t) \hat{y}_t(\epsilon) + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - \partial \ell(\hat{y}_t(\epsilon), y_t) f[t] \right\} \right\} \right] \\ &\leq \mathbb{E} \left[\sup_{\epsilon} \left\{ \mathbb{E}_{r_t} [r_t \cdot \hat{y}_t(\epsilon)] + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - \mathbb{E}_{r_t} [r_t \cdot f[t]] \right\} \right\} \right]\end{aligned}$$

where r_t is a $\{\pm L\}$ -valued random variable with the mean $\partial \ell(\hat{y}_t(\epsilon), y_t)$. With the help of Jensen's inequality, and passing to the worst-case r_t (observe that this is legal for any given ϵ), we have an upper bound

$$\begin{aligned}& \mathbb{E} \left[\sup_{\epsilon} \left\{ \mathbb{E}_{r_t \sim \partial \ell(\hat{y}_t(\epsilon), y_t)} [r_t \cdot \hat{y}_t(\epsilon)] + \mathbb{E}_{r_t \sim \partial \ell(\hat{y}_t(\epsilon), y_t)} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right] \right\} \right] \\ &\leq \mathbb{E} \left[\sup_{r_t \in \{\pm L\}} \left\{ r_t \cdot \hat{y}_t(\epsilon) + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right\} \right] \quad (43)\end{aligned}$$

Now the strategy we defined is

$$\hat{y}_t(\epsilon) = \operatorname{argmin}_{\hat{y}_t} \sup_{r_t \in \{\pm L\}} \left\{ r_t \cdot \hat{y}_t(\epsilon) + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - \sum_{i=1}^{t-1} \ell(f[i], y_i) - r_t \cdot f[t] \right\} \right\}$$

which can be re-written as

$$\hat{y}_t(\epsilon) = \left(\sup_{f \in \mathcal{F}} \left\{ \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2L} L_{t-1}(f) + \frac{1}{2} f[t] \right\} - \sup_{f \in \mathcal{F}} \left\{ \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2L} L_{t-1}(f) - \frac{1}{2} f[t] \right\} \right)$$

By this choice of $\hat{y}_t(\epsilon)$, plugging back in Equation (43) we see that

$$\begin{aligned} & \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_t(f) \right\} \right] \right\} \\ & \leq \mathbb{E}_{\epsilon} \left[\sup_{r_t \in \{\pm L\}} \left\{ r_t \cdot \hat{y}_t(\epsilon) + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right\} \right] \\ & = \mathbb{E}_{\epsilon} \left[\inf_{\hat{y}_t} \sup_{r_t \in \{\pm L\}} \left\{ r_t \cdot \hat{y}_t + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right\} \right] \\ & = \mathbb{E}_{\epsilon} \left[\inf_{\hat{y}_t} \sup_{p_t \in \Delta(\{\pm L\})} \mathbb{E}_{r_t \sim p_t} \left\{ r_t \cdot \hat{y}_t + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right\} \right] \end{aligned}$$

The expression inside the supremum is linear in p_t , as it is an expectation. Also note that the term is convex in \hat{y}_t , and the domain $\hat{y}_t \in [-\sup_{f \in \mathcal{F}} |f[t]|, \sup_{f \in \mathcal{F}} |f[t]|]$ is a bounded interval (hence, compact). We conclude that we can use the minimax theorem, yielding

$$\begin{aligned} & \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \inf_{\hat{y}_t} \mathbb{E}_{r_t \sim p_t} \left[r_t \cdot \hat{y}_t + \sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right] \right] \\ & = \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \left\{ \inf_{\hat{y}_t} \mathbb{E}_{r_t \sim p_t} [r_t \cdot \hat{y}_t] + \mathbb{E}_{r_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right] \right\} \right] \\ & = \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \left\{ \mathbb{E}_{r_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ \inf_{\hat{y}_t} \mathbb{E}_{r_t \sim p_t} [r_t \cdot \hat{y}_t] + 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right] \right\} \right] \\ & \leq \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \left\{ \mathbb{E}_{r_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{r_t \sim p_t} [r_t \cdot f[t]] + 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) - r_t \cdot f[t] \right\} \right] \right\} \right] \end{aligned}$$

In the last step, we replaced the infimum over \hat{y}_t with $f[t]$, only increasing the quantity. Introducing an i.i.d. copy r'_t of r_t ,

$$\begin{aligned} & = \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \left\{ \mathbb{E}_{r_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) + \left(\mathbb{E}_{r_t \sim p_t} [r_t] - r_t \right) \cdot f[t] \right\} \right] \right\} \right] \\ & \leq \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \left\{ \mathbb{E}_{r_t, r'_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) + (r'_t - r_t) \cdot f[t] \right\} \right] \right\} \right] \end{aligned}$$

Introducing the random sign ϵ_t and passing to the supremum over r_t, r'_t , yields the upper bound

$$\begin{aligned} & \mathbb{E}_{\epsilon} \left[\sup_{p_t \in \Delta(\{\pm L\})} \left\{ \mathbb{E}_{r_t, r'_t \sim p_t} \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) + (r'_t - r_t) \cdot f[t] \right\} \right] \right\} \right] \\ & \leq \mathbb{E}_{\epsilon} \left[\sup_{r_t, r'_t \in \{\pm L\}} \left\{ \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) + \epsilon_t (r'_t - r_t) \cdot f[t] \right\} \right] \right\} \right] \\ & \leq \mathbb{E}_{\epsilon} \left[\sup_{r_t, r'_t \in \{\pm L\}} \left\{ \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ L \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2} L_{t-1}(f) + \epsilon_t r'_t \cdot f[t] \right\} \right] \right\} \right] \\ & \quad + \mathbb{E}_{\epsilon} \left[\sup_{r_t, r'_t \in \{\pm L\}} \left\{ \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ L \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2} L_{t-1}(f) - \epsilon_t r_t \cdot f[t] \right\} \right] \right\} \right] \end{aligned}$$

In the above we split the term in the supremum as the sum of two terms one involving r_t and other r'_t (other terms are equally split by dividing by 2), yielding

$$\mathbb{E}_{\epsilon} \left[\sup_{r_t \in \{\pm L\}} \left\{ \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_{t-1}(f) + 2 \epsilon_t r_t \cdot f[t] \right\} \right] \right\} \right]$$

The above step used the fact that the first term only involved r'_t and second only r_t and further ϵ_t and $-\epsilon_t$ have the same distribution. Now finally noting that irrespective of whether r_t in the above supremum is L or $-L$, since it is multiplied by ϵ_t we obtain an upper bound

$$\mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t}^T \epsilon_i f[i] - L_{t-1}(f) \right\} \right]$$

We conclude that the relaxation

$$\mathbf{Rel}_T(\mathcal{F}|y_1, \dots, y_t) = \mathbb{E} \left[\sup_{\epsilon} \left\{ 2L \sum_{i=t+1}^T \epsilon_i f[i] - L_t(f) \right\} \right]$$

is admissible and further the randomized strategy where on each round we first draw ϵ 's and then set

$$\begin{aligned} \hat{y}_t(\epsilon) &= \left(\sup_{f \in \mathcal{F}} \left\{ \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2L} L_{t-1}(f) + \frac{1}{2} f[t] \right\} - \sup_{f \in \mathcal{F}} \left\{ \sum_{i=t+1}^T \epsilon_i f[i] - \frac{1}{2L} L_{t-1}(f) - \frac{1}{2} f[t] \right\} \right) \\ &= \left(\inf_{f \in \mathcal{F}} \left\{ - \sum_{i=t+1}^T \epsilon_i f[i] + \frac{1}{2L} L_{t-1}(f) + \frac{1}{2} f[t] \right\} - \inf_{f \in \mathcal{F}} \left\{ - \sum_{i=t+1}^T \epsilon_i f[i] + \frac{1}{2L} L_{t-1}(f) - \frac{1}{2} f[t] \right\} \right) \end{aligned}$$

is an admissible strategy. Hence, the expected regret under the strategy is bounded as

$$\mathbb{E}[\mathbf{Reg}_T] \leq \mathbf{Rel}_T(\mathcal{F}) = 2L \mathbb{E} \left[\sup_{\epsilon} \sum_{i=1}^T \epsilon_i f[i] \right]$$

which concludes the proof. \square