

A. Duality

The following is a version of the Fenchel duality theorem (see (Rockafellar, 1997)).

Theorem 5. *Let X and Y be Banach spaces, and $f: X \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g: Y \rightarrow \mathbb{R} \cup \{+\infty\}$ convex functions. Let $A: X \rightarrow Y$ be a bounded linear map. If g is continuous at some point $y \in A \text{ dom}(f)$, then the following holds:*

$$\inf_{x \in X} (f(x) + g(Ax)) = \sup_{y^* \in Y^*} (-f^*(A^*y^*) - g^*(-y^*)), \quad (16)$$

where f^* and g^* are conjugate functions of f and g respectively, and A^* the adjoint of A . Furthermore, the supremum in (16) is attained if it is finite.

The following lemma gives the expression of the conjugate function of the (extended) relative entropy, which is a standard result (Boyd & Vandenberghe, 2004).

Lemma 6 (Conjugate function of the relative entropy). *Let $f: \mathbb{R}^{\mathcal{X}} \rightarrow \mathbb{R}$ be defined by $f(\mathbf{p}) = D(\mathbf{p} \parallel \mathbf{p}_0)$ if $\mathbf{p} \in \Delta$ and $f(\mathbf{p}) = +\infty$ elsewhere. Then, the conjugate function of f is the function $f^*: \mathbb{R}^{\mathcal{X}} \rightarrow \mathbb{R}$ defined for all $\mathbf{q} \in \mathbb{R}^{\mathcal{X}}$ by*

$$f^*(\mathbf{q}) = \log \left(\sum_{x \in \mathcal{X}} \mathbf{p}_0[x] e^{\mathbf{q}[x]} \right) = \log \left(\mathbb{E}_{x \sim \mathbf{p}_0} [e^{\mathbf{q}[x]}] \right).$$

Proof. By definition of f , for any $\mathbf{q} \in \mathbb{R}^{\mathcal{X}}$, we can write

$$\sup_{\mathbf{p} \in \mathbb{R}^{\mathcal{X}}} (\langle \mathbf{p}, \mathbf{q} \rangle - D(\mathbf{p} \parallel \mathbf{p}_0)) = \sup_{\mathbf{p} \in \Delta} (\langle \mathbf{p}, \mathbf{q} \rangle - D(\mathbf{p} \parallel \mathbf{p}_0)).$$

Fix $\mathbf{q} \in \mathbb{R}^{\mathcal{X}}$ and let $\bar{\mathbf{q}} \in \Delta$ be defined for all $x \in \mathcal{X}$ by

$$\bar{\mathbf{q}}[x] = \frac{\mathbf{p}_0[x] e^{\mathbf{q}[x]}}{\sum_{x \in \mathcal{X}} \mathbf{p}_0[x] e^{\mathbf{q}[x]}} = \frac{\mathbf{p}_0[x] e^{\mathbf{q}[x]}}{\mathbb{E}_{\mathbf{p}_0}[e^{\mathbf{q}}]}. \quad (17)$$

Then, the following holds for all $\mathbf{p} \in \Delta$:

$$\begin{aligned} \langle \mathbf{p}, \mathbf{q} \rangle - D(\mathbf{p} \parallel \mathbf{p}_0) &= \mathbb{E}_{\mathbf{p}}[\log(e^{\mathbf{q}})] - \mathbb{E}_{\mathbf{p}} \left[\log \frac{\mathbf{p}}{\mathbf{p}_0} \right] \\ &= \mathbb{E}_{\mathbf{p}} \left[\log \frac{\mathbf{p}_0 e^{\mathbf{q}}}{\mathbf{p}} \right] \\ &= -D(\mathbf{p} \parallel \bar{\mathbf{q}}) + \log \mathbb{E}_{\mathbf{p}_0}[e^{\mathbf{q}}]. \end{aligned}$$

Since $D(\mathbf{p} \parallel \bar{\mathbf{q}}) \geq 0$ and $D(\mathbf{p} \parallel \bar{\mathbf{q}}) = 0$ for $\mathbf{p} = \bar{\mathbf{q}}$, this shows that $\sup_{\mathbf{p} \in \Delta} (\langle \mathbf{p}, \mathbf{q} \rangle - D(\mathbf{p} \parallel \mathbf{p}_0)) = \log(\mathbb{E}_{\mathbf{p}_0}[e^{\mathbf{q}}])$ and concludes the proof. \square

Theorem 1. *Problem (2) and (3) are equivalent to the dual optimization problem $\sup_{\mathbf{w} \in \mathbb{R}^N} G(\mathbf{w})$:*

$$\sup_{\mathbf{w} \in \mathbb{R}^N} G(\mathbf{w}) = \min_{\mathbf{p}} F(\mathbf{p}). \quad (18)$$

Furthermore, let $\mathbf{p}^* = \arg\min_{\mathbf{p}} F(\mathbf{p})$, then, for any $\epsilon > 0$ and any \mathbf{w} such that $|G(\mathbf{w}) - \sup_{\mathbf{w} \in \mathbb{R}^N} G(\mathbf{w})| < \epsilon$, the following inequality holds: $D(\mathbf{p}^* \parallel \mathbf{p}_{\mathbf{w}}) \leq \epsilon$.

Proof. The proof follows by application of the Fenchel duality theorem (Theorem 5, Appendix A) to the optimization problem (3) with the functions f and g defined for all \mathbf{p} and \mathbf{u} by $f(\mathbf{p}) = D(\mathbf{p} \parallel \mathbf{p}_0) + I_{\Delta}(\mathbf{p})$ and $g(\mathbf{u}) = I_C(\mathbf{u})$ and with A the linear map defined by $A\mathbf{p} = \mathbb{E}_{\mathbf{p}}[\Phi]$.

A is a bounded linear map since $\|A\| \leq \|\Phi\|_{\infty} \leq \Lambda$ and $A^*\mathbf{w} = \mathbf{w} \cdot \Phi$. Furthermore, define $\mathbf{u} \in \mathbb{F}$ by $\mathbf{u}_k = \mathbb{E}_{\mathbf{p}}[\Phi_k]$. Then, \mathbf{u} is in $A \text{ dom} f$ and is in C . Since $\beta_k > 0$ for all k , \mathbf{u} is contained in $\text{int}(C)$. $g = I_C$ equals zero over $\text{int}(C)$ and is therefore continuous over $\text{int}(C)$, thus g is continuous at $\mathbf{u} \in A \text{ dom} f$.

By Lemma 6, the conjugate of f is the function $f^*: \mathbb{R}^{\mathcal{X}} \rightarrow \mathbb{R}$ defined by $f^*(\mathbf{q}) = \log(\sum_{x \in \mathcal{X}} \mathbf{p}_0[x] e^{\mathbf{q}[x]})$ for all $\mathbf{q} \in \mathbb{R}^{\mathcal{X}}$. The conjugate function of $g = I_C$ is the function g^* defined for all $\mathbf{w} \in \mathbb{R}^N$ by

$$\begin{aligned} g^*(\mathbf{w}) &= \sup_{\mathbf{u} \in C} (\mathbf{w} \cdot \mathbf{u} - I_C(\mathbf{u})) \\ &= \sup_{\mathbf{u} \in C} (\mathbf{w} \cdot \mathbf{u}) \\ &= \sup_{\mathbf{u} \in C} \left(\sum_{k=1}^p \mathbf{w}_k \cdot \mathbf{u}_k \right) \\ &= \sum_{k=1}^p \sup_{\|\mathbf{u}_k - \mathbb{E}_S[\Phi_k]\|_{\infty} \leq \beta_k} (\mathbf{w}_k \cdot \mathbf{u}_k) \\ &= \sum_{k=1}^p \mathbf{w}_k \cdot \mathbb{E}_S[\Phi_k] + \sum_{k=1}^p \sup_{\|\mathbf{u}_k\|_{\infty} \leq \beta_k} (\mathbf{w}_k \cdot \mathbf{u}_k) \\ &= \mathbb{E}_S[\mathbf{w} \cdot \Phi] + \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1, \end{aligned}$$

where the penultimate equality holds by definition of the dual norm. In view of these identities, we can write

$$\begin{aligned} &-f^*(A^*\mathbf{w}) - g^*(-\mathbf{w}) \\ &= -\log \left(\sum_{x \in \mathcal{X}} \mathbf{p}_0[x] e^{\mathbf{w} \cdot \Phi(x)} \right) + \mathbb{E}_S[\mathbf{w} \cdot \Phi] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ &= -\log Z_{\mathbf{w}} + \frac{1}{m} \sum_{i=1}^m \mathbf{w} \cdot \Phi(x_i) - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ &= \frac{1}{m} \sum_{i=1}^m \log \frac{e^{\mathbf{w} \cdot \Phi(x_i)}}{Z_{\mathbf{w}}} - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ &= \frac{1}{m} \sum_{i=1}^m \log \left[\frac{\mathbf{p}_{\mathbf{w}}[x_i]}{\mathbf{p}_0[x_i]} \right] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 = G(\mathbf{w}), \end{aligned}$$

which proves that $\sup_{\mathbf{w} \in \mathbb{R}^N} G(\mathbf{w}) = \min_{\mathbf{p}} F(\mathbf{p})$. For any

$\mathbf{w} \in \mathbb{R}^N$, we can write

$$\begin{aligned}
 & G(\mathbf{w}) - D(\mathbf{p}^* \parallel \mathbf{p}_0) \\
 &= \mathbb{E}_{x \sim \hat{\mathbf{p}}} \left[\log \frac{\mathbf{p}_{\mathbf{w}}[x]}{\mathbf{p}_0[x]} \right] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 - \mathbb{E}_{x \sim \mathbf{p}^*} \left[\log \frac{\mathbf{p}^*[x]}{\mathbf{p}_0[x]} \right] \\
 &= \mathbb{E}_{x \sim \hat{\mathbf{p}}} \left[\log \frac{\mathbf{p}_{\mathbf{w}}[x]}{\mathbf{p}_0[x]} \right] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 - \\
 &\quad \mathbb{E}_{x \sim \mathbf{p}^*} \left[\log \frac{\mathbf{p}^*[x] \mathbf{p}_{\mathbf{w}}[x]}{\mathbf{p}_{\mathbf{w}}[x] \mathbf{p}_0[x]} \right] \\
 &= -D(\mathbf{p}^* \parallel \mathbf{p}_{\mathbf{w}}) - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\
 &\quad + \mathbb{E}_{x \sim \hat{\mathbf{p}}} \left[\log \frac{\mathbf{p}_{\mathbf{w}}[x]}{\mathbf{p}_0[x]} \right] - \mathbb{E}_{x \sim \mathbf{p}^*} \left[\log \frac{\mathbf{p}_{\mathbf{w}}(x)}{\mathbf{p}_0(x)} \right].
 \end{aligned}$$

The difference of the last two terms can be expressed as follows

$$\begin{aligned}
 & \mathbb{E}_{x \sim \hat{\mathbf{p}}} \left[\log \frac{\mathbf{p}_{\mathbf{w}}[x]}{\mathbf{p}_0[x]} \right] - \mathbb{E}_{x \sim \mathbf{p}^*} \left[\log \frac{\mathbf{p}_{\mathbf{w}}[x]}{\mathbf{p}_0[x]} \right] \\
 &= \mathbb{E}_{x \sim \hat{\mathbf{p}}} [\mathbf{w} \cdot \Phi(x) - \log Z_{\mathbf{w}}] - \mathbb{E}_{x \sim \mathbf{p}^*} [\mathbf{w} \cdot \Phi(x) - \log Z_{\mathbf{w}}] \\
 &= \mathbb{E}_{x \sim \hat{\mathbf{p}}} [\mathbf{w} \cdot \Phi(x)] - \mathbb{E}_{x \sim \mathbf{p}^*} [\mathbf{w} \cdot \Phi(x)].
 \end{aligned}$$

Plugging back this equality and rearranging yields

$$\begin{aligned}
 D(\mathbf{p}^* \parallel \mathbf{p}_{\mathbf{w}}) &= D(\mathbf{p}^* \parallel \mathbf{p}_0) - G(\mathbf{w}) \\
 &\quad - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 + \mathbf{w} \cdot \left(\mathbb{E}_{x \sim \hat{\mathbf{p}}} [\Phi(x)] - \mathbb{E}_{x \sim \mathbf{p}^*} [\Phi(x)] \right).
 \end{aligned}$$

The solution of the primal optimization, \mathbf{p}^* , verifies the constraint $I_C(\mathbb{E}_{\mathbf{p}^*}[\Phi]) = 0$, that is $\|\mathbb{E}_{x \sim \hat{\mathbf{p}}}[\Phi_k(x)] - \mathbb{E}_{x \sim \mathbf{p}^*}[\Phi_k(x)]\|_{\infty} \leq \beta_k$ for all $k \in [1, p]$. By Hölder's inequality, this implies that

$$\begin{aligned}
 & - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 + \mathbf{w} \cdot \left(\mathbb{E}_{x \sim \hat{\mathbf{p}}} [\Phi(x)] - \mathbb{E}_{x \sim \mathbf{p}^*} [\Phi(x)] \right) \\
 &= - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 + \sum_{k=1}^p \mathbf{w}_k \cdot \left(\mathbb{E}_{x \sim \hat{\mathbf{p}}} [\Phi_k(x)] - \mathbb{E}_{x \sim \mathbf{p}^*} [\Phi_k(x)] \right) \\
 &\leq - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 + \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 = 0.
 \end{aligned}$$

Thus, we can write, for any $\mathbf{w} \in \mathbb{R}^N$,

$$D(\mathbf{p}^* \parallel \mathbf{p}_{\mathbf{w}}) \leq D(\mathbf{p}^* \parallel \mathbf{p}_0) - G(\mathbf{w}).$$

Now, assume that \mathbf{w} verifies $|G(\mathbf{w}) - \sup_{\mathbf{w} \in \mathbb{R}^N} G(\mathbf{w})| \leq \epsilon$ for some $\epsilon > 0$. Then, $D(\mathbf{p}^* \parallel \mathbf{p}_0) - G(\mathbf{w}) = \sup_{\mathbf{w}} G(\mathbf{w}) - G(\mathbf{w}) \leq \epsilon$ implies $D(\mathbf{p}^* \parallel \mathbf{p}_{\mathbf{w}}) \leq \epsilon$. This concludes the proof of the theorem. \square

Theorem 4. Problem (10) is equivalent to dual optimization problem $\sup_{\mathbf{w} \in \mathbb{R}^N} \tilde{G}(\mathbf{w})$:

$$\sup_{\mathbf{w} \in \mathbb{R}^N} \tilde{G}(\mathbf{w}) = \min_{\mathbf{p}} \tilde{F}(\mathbf{p}). \quad (19)$$

Furthermore, let $\mathbf{p}^* = \operatorname{argmin}_{\mathbf{p}} \tilde{F}(\mathbf{p})$. Then, for any $\epsilon > 0$ and any \mathbf{w} such that $|\tilde{G}(\mathbf{w}) - \sup_{\mathbf{w} \in \mathbb{R}^N} \tilde{G}(\mathbf{w})| < \epsilon$, we have $\mathbb{E}_{x \sim \hat{\mathbf{p}}} [D(\mathbf{p}^*[\cdot|x] \parallel \mathbf{p}_0[\cdot|x])] \leq \epsilon$.

Proof. The proof follows by application of the Fenchel duality theorem (Theorem 5, Appendix A) to the optimization problem (11) with the functions \tilde{f} and \tilde{g} defined for all \mathbf{p} and \mathbf{u} by $\tilde{f}(\mathbf{p}) = \mathbb{E}_{x \sim \hat{\mathbf{p}}} [D(\mathbf{p}[\cdot|x] \parallel \mathbf{p}_0[\cdot|x]) + I_{\Delta}(\mathbf{p}[\cdot|x])]$ and $\tilde{g}(\mathbf{u}) = I_C(\mathbf{u})$ and with A the linear map defined by $A\mathbf{p} = \mathbb{E}_{\substack{x \sim \hat{\mathbf{p}} \\ y \sim \mathbf{p}[\cdot|x]}} [\Phi(x, y)]$.

A is a bounded linear map since $\|A\| \leq \|\Phi\|_{\infty} \leq \Lambda$. Note that

$$\begin{aligned}
 A\mathbf{p} &= \mathbb{E}_{\substack{x \sim \hat{\mathbf{p}} \\ y \sim \mathbf{p}[\cdot|x]}} [\Phi(x, y)] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \Phi(x, y) \hat{\mathbf{p}}[x] \mathbf{p}[y|x] \\
 &= \sum_{x \in \operatorname{supp}(\hat{\mathbf{p}})} (\hat{\mathbf{p}}[x] \Phi(x, \cdot)) \cdot (\mathbf{p}[\cdot|x]).
 \end{aligned}$$

Thus, the conjugate of A is defined for all $\mathbf{w} \in \mathbb{R}^N$ by $A^*\mathbf{w} = \mathbf{w} \cdot (\hat{\mathbf{p}}(x) \Phi(x, y))$. Furthermore, define $\mathbf{u} \in \mathbb{F}$ by $\mathbf{u}_k = \mathbb{E}_{(x, y) \sim S} [\Phi_k(x, y)]$. Then, \mathbf{u} is in $A \operatorname{dom} f$ and is in C . Since $\beta_k > 0$ for all k , \mathbf{u} is contained in $\operatorname{int}(C)$. $g = I_C$ equals zero over $\operatorname{int}(C)$ and is therefore continuous over $\operatorname{int}(C)$, thus g is continuous at $\mathbf{u} \in A \operatorname{dom} f$.

The conjugate function of \tilde{f} is defined for all $\mathbf{q} = (\mathbf{q}[\cdot|x_i])_{i \in [1, m]}$ by

$$\begin{aligned}
 \tilde{f}^*(\mathbf{q}) &= \sup_{\mathbf{p}[\cdot|x] \in \Delta} \left\{ \langle \mathbf{p}, \mathbf{q} \rangle - \sum_{x \in \mathcal{X}} \hat{\mathbf{p}}[x] D(\mathbf{p}[\cdot|x] \parallel \mathbf{p}_0[\cdot|x]) \right\} \\
 &= \sup_{\mathbf{p}[\cdot|x] \in \Delta} \left\{ \sum_{x \in \operatorname{supp}(\hat{\mathbf{p}})} \hat{\mathbf{p}}[x] \sum_{y \in \mathcal{Y}} \mathbf{p}[y|x] \mathbf{q}[y|x] (\hat{\mathbf{p}}[x])^{-1} \right. \\
 &\quad \left. - \sum_{x \in \operatorname{supp}(\hat{\mathbf{p}})} \hat{\mathbf{p}}[x] D(\mathbf{p}[\cdot|x] \parallel \mathbf{p}_0[\cdot|x]) \right\} \\
 &= \sum_{x \in \operatorname{supp}(\hat{\mathbf{p}})} \hat{\mathbf{p}}[x] \sup_{\mathbf{p}[\cdot|x]} \left\{ \sum_{y \in \mathcal{Y}} \mathbf{p}[y|x] \left(\frac{\mathbf{q}[y|x]}{\hat{\mathbf{p}}[x]} \right) \right. \\
 &\quad \left. - D(\mathbf{p}[\cdot|x] \parallel \mathbf{p}_0[\cdot|x]) \right\} \\
 &= \sum_{x \in \operatorname{supp}(\hat{\mathbf{p}})} \hat{\mathbf{p}}[x] f_x^* \left(\frac{\mathbf{q}[y|x]}{\hat{\mathbf{p}}[x]} \right)
 \end{aligned}$$

where f_x is defined for all $x \in \mathcal{X}$ and $\mathbf{p}' \in \mathbb{R}^{\mathcal{Y}}$ by $f_x(\mathbf{p}') = D(\mathbf{p}' \parallel \mathbf{p}_0[\cdot|x])$ if $\mathbf{p}' \in \Delta$, $f_x(\mathbf{p}') = +\infty$ otherwise. By

Lemma 6, $f_x^* \left(\frac{q[y|x]}{\hat{p}[x]} \right) = \log \left(\sum_{y \in \mathcal{Y}} p_0[y|x] e^{\frac{q[y|x]}{\hat{p}[x]}} \right)$, thus, \tilde{f}^* is given by

$$\tilde{f}^*(\mathbf{q}) = \mathbb{E}_{x \sim \hat{p}} \left[\log \left(\sum_{y \in \mathcal{Y}} p_0[y|x] e^{\frac{q[y|x]}{\hat{p}[x]}} \right) \right].$$

In view of these identities, we can write

$$\begin{aligned} & -\tilde{f}^*(A^* \mathbf{w}) - \tilde{g}^*(-\mathbf{w}) \\ &= -\mathbb{E}_{x \sim \hat{p}} \left[\log \left(\sum_{y \in \mathcal{Y}} p_0[y|x] e^{\mathbf{w} \cdot \Phi(x,y)} \right) \right] \\ & \quad + \mathbb{E}_S [\mathbf{w} \cdot \Phi] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ &= -\mathbb{E}_{x \sim \hat{p}} [\log Z_{\mathbf{w}}(x)] + \frac{1}{m} \sum_{i=1}^m \mathbf{w} \cdot \Phi(x_i, y_i) - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ &= \frac{1}{m} \sum_{i=1}^m \log \frac{e^{\mathbf{w} \cdot \Phi(x_i, y_i)}}{Z_{\mathbf{w}}(x_i)} - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ &= \frac{1}{m} \sum_{i=1}^m \log \left[\frac{p_{\mathbf{w}}[y_i|x_i]}{p_0[y_i|x_i]} \right] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 = \tilde{G}(\mathbf{w}), \end{aligned}$$

which proves that $\sup_{\mathbf{w} \in \mathbb{R}^N} \tilde{G}(\mathbf{w}) = \min_{\mathbf{p}} \tilde{F}(\mathbf{p})$. The second part of the proof is similar to that of Theorem 1. For any $\mathbf{w} \in \mathbb{R}^N$, we can write

$$\begin{aligned} & \tilde{G}(\mathbf{w}) - \mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_0[\cdot|x])] \\ &= \mathbb{E}_{(x,y) \sim S} \left[\log \frac{p_{\mathbf{w}}[y|x]}{p_0[y|x]} \right] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ & \quad - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} \left[\log \frac{p^*[y|x]}{p_0[y|x]} \right] \\ &= \mathbb{E}_{(x,y) \sim S} \left[\log \frac{p_{\mathbf{w}}[y|x]}{p_0[y|x]} \right] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 - \\ & \quad \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} \left[\log \frac{p^*[y|x] p_{\mathbf{w}}[y|x]}{p_{\mathbf{w}}[y|x] p_0[y|x]} \right] \\ &= -\mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_{\mathbf{w}}[\cdot|x])] - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ & \quad + \mathbb{E}_{(x,y) \sim S} \left[\log \frac{p_{\mathbf{w}}[y|x]}{p_0[y|x]} \right] - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} \left[\log \frac{p_{\mathbf{w}}[y|x]}{p_0[y|x]} \right]. \end{aligned}$$

The difference of the last two terms can be expressed as

follows

$$\begin{aligned} & \mathbb{E}_{(x,y) \sim S} \left[\log \frac{p_{\mathbf{w}}[y|x]}{p_0[y|x]} \right] - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} \left[\log \frac{p_{\mathbf{w}}[y|x]}{p_0[y|x]} \right] \\ &= \mathbb{E}_{(x,y) \sim S} [\mathbf{w} \cdot \Phi(x,y) - \log Z_{\mathbf{w}}(x)] \\ & \quad - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} [\mathbf{w} \cdot \Phi(x,y) - \log Z_{\mathbf{w}}(x)] \\ &= \mathbb{E}_{(x,y) \sim S} [\mathbf{w} \cdot \Phi(x,y)] - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} [\mathbf{w} \cdot \Phi(x,y)]. \end{aligned}$$

Plugging back this equality and rearranging yields

$$\begin{aligned} & \mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_{\mathbf{w}}[\cdot|x])] \\ &= \mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_0[\cdot|x])] - \tilde{G}(\mathbf{w}) - \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ & \quad + \mathbf{w} \cdot \left[\mathbb{E}_{(x,y) \sim S} [\mathbf{w} \cdot \Phi(x,y)] - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} [\mathbf{w} \cdot \Phi(x,y)] \right]. \end{aligned}$$

The solution of the primal optimization, p^* , verifies the constraint $I_C(\mathbb{E}_{x \sim \hat{p}} [\Phi(x,y)]) = 0$, that is $\|\mathbb{E}_{x \sim \hat{p}} [\Phi(x,y)] - \mathbb{E}_{(x,y) \sim S} [\Phi(x,y)]\|_{\infty} \leq \beta_k$ for all $k \in [1, p]$. By Hölder's inequality, this implies that

$$\begin{aligned} & -\sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 \\ & \quad + \mathbf{w} \cdot \left[\mathbb{E}_{(x,y) \sim S} [\mathbf{w} \cdot \Phi(x,y)] - \mathbb{E}_{\substack{x \sim \hat{p} \\ y \sim p^*[\cdot|x]}} [\mathbf{w} \cdot \Phi(x,y)] \right] \\ & \leq -\sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 + \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 = 0. \end{aligned}$$

Thus, we can write, for any $\mathbf{w} \in \mathbb{R}^N$,

$$\begin{aligned} & \mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_{\mathbf{w}}[\cdot|x])] \\ & \leq \mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_0[\cdot|x])] - \tilde{G}(\mathbf{w}). \end{aligned}$$

Now, assume that \mathbf{w} verifies $|\tilde{G}(\mathbf{w}) - \sup_{\mathbf{w} \in \mathbb{R}^N} \tilde{G}(\mathbf{w})| \leq \epsilon$ for some $\epsilon > 0$. Then, $\mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_0[\cdot|x])] - \tilde{G}(\mathbf{w}) = \sup_{\mathbf{w}} \tilde{G}(\mathbf{w}) - \tilde{G}(\mathbf{w}) \leq \epsilon$ implies $\mathbb{E}_{x \sim \hat{p}} [D(p^*[\cdot|x] \parallel p_{\mathbf{w}}[\cdot|x])] \leq \epsilon$. This concludes the proof of the theorem. \square

B. Pseudocode of StructMaxent1

Figure 2 shows the pseudocode of StructMaxent1.

```

STRUCTMAXENT1( $S = (x_1, \dots, x_m)$ )
1  for  $t \leftarrow 1$  to  $T$  do
2      for  $k \leftarrow 1$  to  $p$  and  $j \leftarrow 1$  to  $N_k$  do
3          if  $(w_{t-1,k,j} \neq 0)$  then
4               $d_{k,j} \leftarrow \beta_k \operatorname{sgn}(w_{t-1,k,j}) + \epsilon_{t-1,k,j}$ 
5              elseif  $|\epsilon_{t-1,k,j}| \leq \beta_k$  then
6                   $d_{k,j} \leftarrow 0$ 
7              else  $d_{k,j} \leftarrow -\beta_k \operatorname{sgn}(\epsilon_{t-1,k,j}) + \epsilon_{t-1,k,j}$ 
8               $(k,j) \leftarrow \operatorname{argmax}_{(k,j) \in [1,p] \times [1,N_k]} |d_{k,j}|$ 
9               $\beta \leftarrow \frac{\bar{\Phi}_{t-1,k,j}^+ \bar{\Phi}_{k,j}^- e^{-2w_{k,j}\Lambda} - \bar{\Phi}_{k,j}^+ \bar{\Phi}_{t-1,k,j}^-}{\bar{\Phi}_{t-1,k,j}^+ e^{-2w_{k,j}\Lambda} - \bar{\Phi}_{t-1,k,j}^-}$ 
10             if  $(|\beta| \leq \beta_k)$  then
11                  $\eta \leftarrow -w_{t-1,k,j}$ 
12             elseif  $(\beta > \beta_k)$  then
13                  $\eta \leftarrow \frac{1}{2\Lambda} \log \left[ \frac{\bar{\Phi}_{t-1,k,j}^-(\beta_k - \bar{\Phi}_{k,j}^+)}{\bar{\Phi}_{t-1,k,j}^+(\beta_k - \bar{\Phi}_{k,j}^-)} \right]$ 
14             else  $\eta \leftarrow \frac{1}{2\Lambda} \log \left[ \frac{\bar{\Phi}_{t-1,k,j}^-(\beta_k + \bar{\Phi}_{k,j}^+)}{\bar{\Phi}_{t-1,k,j}^+(\beta_k + \bar{\Phi}_{k,j}^-)} \right]$ 
15              $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} + \eta \mathbf{e}_{k,j}$ 
16              $\rho_{\mathbf{w}_t} \leftarrow \frac{\rho_0[x] e^{\mathbf{w}_t \cdot \Phi(x)}}{\sum_{x \in \mathcal{X}} \rho_0[x] e^{\mathbf{w}_t \cdot \Phi(x)}}$ 
17  return  $\rho_{\mathbf{w}_t}$ 
    
```

Figure 2. Pseudocode of the StructMaxent1 algorithm. For all $(k, j) \in [1, p] \times [1, N_k]$, $\beta_k = 2\mathfrak{R}_m(H_k) + \beta$, $\epsilon_{t-1,k,j} = \mathbb{E}_{\rho_{\mathbf{w}_{t-1}}}[\Phi_{k,j}] - \mathbb{E}_S[\Phi_{k,j}]$ and, for any $s \in \{-1, +1\}$, $\bar{\Phi}_{t-1,k,j}^s = \mathbb{E}_{\rho_{\mathbf{w}_{t-1}}}[\Phi_{k,j}] + s\Lambda$ and $\bar{\Phi}_{k,j}^s = \mathbb{E}_S[\Phi_{k,j}] + s\Lambda$. The closed-form solutions for the step size given here assume that the conditions (8) hold.

C. Algorithm

In this section we derive the step size for the StructMaxent1 and StructMaxent2 algorithms presented in Section 2.4 and Appendix B.

Observe that

$$\begin{aligned}
 & F(\mathbf{w}_{t-1} + \eta \mathbf{e}_{k,j}) - F(\mathbf{w}_{t-1}) \\
 &= \beta_k (|w_{k,j} + \eta| - |w_{k,j}|) - \eta \mathbb{E}_S[\Phi_{k,j}] + \log \left[\frac{\mathbb{E}[e^{\eta \Phi_{k,j}}]}{\rho_{\mathbf{w}_{t-1}}} \right].
 \end{aligned} \tag{20}$$

Since $\Phi_{k,j} \in [-\Lambda, +\Lambda]$, by the convexity of $x \mapsto e^{\eta x}$, we can write

$$e^{\eta \Phi_{k,j}} \leq \frac{\Lambda - \Phi_{k,j}}{2\Lambda} e^{-\eta\Lambda} + \frac{\Phi_{k,j} + \Lambda}{2\Lambda} e^{\eta\Lambda}.$$

Taking the expectation and the log yields

$$\begin{aligned}
 \log \mathbb{E}_{\rho_{\mathbf{w}_{t-1}}} [e^{\eta \Phi_{k,j}}] &\leq \log \left[\frac{\bar{\Phi}_{t-1,k,j}^+ e^{\eta\Lambda} - \bar{\Phi}_{t-1,k,j}^- e^{-\eta\Lambda}}{2\Lambda} \right] \\
 &= -\eta\Lambda + \log \left[\frac{\bar{\Phi}_{t-1,k,j}^+ e^{2\eta\Lambda} - \bar{\Phi}_{t-1,k,j}^-}{2\Lambda} \right],
 \end{aligned}$$

where we used the following notation:

$$\bar{\Phi}_{t-1,k,j}^s = \mathbb{E}_{\rho_{\mathbf{w}_{t-1}}} [\Phi_{k,j}] + s\Lambda \quad \bar{\Phi}_{k,j}^s = \mathbb{E}_S [\Phi_{k,j}] + s\Lambda,$$

for all $(k, j) \in [1, p] \times [1, N_k]$ and $s \in \{-1, +1\}$.

Plugging back this inequality in (20) and ignoring constant terms, minimizing the resulting upper bound on $F(\mathbf{w}_{t-1} + \eta \mathbf{e}_{k,j}) - F(\mathbf{w}_{t-1})$ becomes equivalent to minimizing $\psi(\eta)$ defined for all $\eta \in \mathbb{R}$ by

$$\psi(\eta) = \beta_k |w_{k,j} + \eta| - \eta \bar{\Phi}_{k,j}^+ + \log \left[\bar{\Phi}_{t-1,k,j}^+ e^{2\eta\Lambda} - \bar{\Phi}_{t-1,k,j}^- \right].$$

Let η^* denote the minimizer of $\psi(\eta)$. If $w_{t-1,k,j} + \eta^* = 0$, then the subdifferential of $|w_{t-1,k,j} + \eta|$ at η^* is the set $\{\nu : \nu \in [-1, +1]\}$. Thus, in that case, the subdifferential $\partial\psi(\eta^*)$, contains 0 iff there exists $\nu \in [-1, +1]$ such that

$$\begin{aligned}
 & \beta_k \nu - \bar{\Phi}_{k,j}^+ + \frac{2\Lambda \bar{\Phi}_{t-1,k,j}^+ e^{2\eta^*\Lambda}}{\bar{\Phi}_{t-1,k,j}^+ e^{2\eta^*\Lambda} - \bar{\Phi}_{t-1,k,j}^-} = 0 \\
 \Leftrightarrow & \bar{\Phi}_{k,j}^+ - \frac{2\Lambda \bar{\Phi}_{t-1,k,j}^+ e^{-2w_{t-1,k,j}\Lambda}}{\bar{\Phi}_{t-1,k,j}^+ e^{-2w_{t-1,k,j}\Lambda} - \bar{\Phi}_{t-1,k,j}^-} = \beta_k \nu.
 \end{aligned}$$

Thus, the condition is equivalent to

$$\left| \bar{\Phi}_{k,j}^+ - \frac{2\Lambda \bar{\Phi}_{t-1,k,j}^+ e^{-2w_{t-1,k,j}\Lambda}}{\bar{\Phi}_{t-1,k,j}^+ e^{-2w_{t-1,k,j}\Lambda} - \bar{\Phi}_{t-1,k,j}^-} \right| \leq \beta_k,$$

which can be rewritten as

$$\left| \frac{\bar{\Phi}_{t-1,k,j}^+ \bar{\Phi}_{k,j}^- e^{-2w_{k,j}\Lambda} - \bar{\Phi}_{k,j}^+ \bar{\Phi}_{t-1,k,j}^-}{\bar{\Phi}_{t-1,k,j}^+ e^{-2w_{k,j}\Lambda} - \bar{\Phi}_{t-1,k,j}^-} \right| \leq \beta_k.$$

If $w_{t-1,k,j} + \eta^* > 0$, then ψ is differentiable at η^* and $\psi'(\eta^*) = 0$, that is

$$\begin{aligned}
 & \beta - \bar{\Phi}_{k,j}^+ + \frac{2\Lambda \bar{\Phi}_{t-1,k,j}^+ e^{2\eta^*\Lambda}}{\bar{\Phi}_{t-1,k,j}^+ e^{2\eta^*\Lambda} - \bar{\Phi}_{t-1,k,j}^-} = 0 \\
 \Leftrightarrow & e^{2\eta^*\Lambda} = \frac{\bar{\Phi}_{t-1,k,j}^- (\beta_k - \bar{\Phi}_{k,j}^+)}{\bar{\Phi}_{t-1,k,j}^+ (\beta_k - \bar{\Phi}_{k,j}^-)} \\
 \Leftrightarrow & \eta^* = \frac{1}{2\Lambda} \log \left[\frac{\bar{\Phi}_{t-1,k,j}^- (\beta_k - \bar{\Phi}_{k,j}^+)}{\bar{\Phi}_{t-1,k,j}^+ (\beta_k - \bar{\Phi}_{k,j}^-)} \right].
 \end{aligned}$$

For the step size η^* to be in \mathbb{R} , the following conditions must be met:

$$\begin{aligned}
 & (\bar{\Phi}_{t-1,k,j}^- \neq 0) \wedge (\bar{\Phi}_{t-1,k,j}^+ \neq 0) \wedge \\
 & ((\beta_k - \bar{\Phi}_{k,j}^+) < 0) \wedge ((\beta_k - \bar{\Phi}_{k,j}^-) \neq 0),
 \end{aligned}$$

that is

$$(\mathbb{E}_{\rho_{\mathbf{w}_{t-1}}} [\Phi_{k,j}] \notin \{-\Lambda, +\Lambda\}) \wedge (\mathbb{E}_S [\Phi_{k,j}] > -\Lambda + \beta_k). \tag{21}$$

The condition $w_{t-1,k,j} + \eta^* > 0$ is equivalent to $e^{2\eta^*\Lambda} > e^{-2w_{t-1,k,j}\Lambda}$, which, in view of the expression of $e^{2\eta^*\Lambda}$ given above can be written as

$$\frac{\bar{\Phi}_{t-1,k,j}^+ \bar{\Phi}_{k,j}^- e^{-2w_{k,j}\Lambda} - \bar{\Phi}_{k,j}^+ \bar{\Phi}_{t-1,k,j}^-}{\bar{\Phi}_{t-1,k,j}^+ e^{-2w_{k,j}\Lambda} - \bar{\Phi}_{t-1,k,j}^-} > \beta_k.$$

Similarly, if $w_{t-1,k,j} + \eta^* < 0$, ψ is differentiable at η^* and $\psi'(\eta^*) = 0$, which gives

$$\eta^* = \frac{1}{2\Lambda} \log \left[\frac{\bar{\Phi}_{t-1,k,j}^- (\beta_k + \bar{\Phi}_{k,j}^+)}{\bar{\Phi}_{t-1,k,j}^+ (\beta_k + \bar{\Phi}_{k,j}^-)} \right].$$

Again for the step size η^* to be in \mathbb{R} , the following conditions must be met:

$$\begin{aligned} & (\bar{\Phi}_{t-1,k,j}^- \neq 0) \wedge (\bar{\Phi}_{t-1,k,j}^+ \neq 0) \wedge \\ & ((\beta_k + \bar{\Phi}_{k,j}^+) \neq 0) \wedge ((\beta_k + \bar{\Phi}_{k,j}^-) < 0), \end{aligned}$$

that is

$$(\mathbb{E}_{\mathbf{P}_{\mathbf{w}_{t-1}}} [\Phi_{k,j}] \notin \{-\Lambda, +\Lambda\}) \wedge (\mathbb{E}_S [\Phi_{k,j}] < \Lambda - \beta_k).$$

Combining with condition 21, the following condition on Φ , Λ and β_k must be satisfied:

$$\begin{aligned} & (\mathbb{E}_{\mathbf{P}_{\mathbf{w}_{t-1}}} [\Phi_{k,j}] \notin \{-\Lambda, +\Lambda\}) \wedge \\ & (-\Lambda + \beta_k < \mathbb{E}_S [\Phi_{k,j}] < \Lambda - \beta_k). \end{aligned}$$

Figure 2 shows the pseudocode of our algorithm using the closed-form solution for the step size just presented.

An alternative method consists of using a somewhat looser upper bound for $\log \mathbb{E}_{\mathbf{P}_{\mathbf{w}_{t-1}}} [e^{\eta\Phi_{k,j}}]$ using Hoeffding's lemma and $\Phi_{k,j} \in [-\Lambda, +\Lambda]$:

$$\log \mathbb{E}_{\mathbf{P}_{\mathbf{w}_{t-1}}} [e^{\eta\Phi_{k,j}}] \leq \eta \mathbb{E}_{\mathbf{P}_{\mathbf{w}_{t-1}}} [\Phi_{k,j}] + \frac{\eta^2 \Lambda^2}{2}.$$

Combining this inequality with (20) and disregarding constant terms, minimizing the resulting upper bound on $F(\mathbf{w}_{t-1} + \eta \mathbf{e}_{k,j}) - F(\mathbf{w}_{t-1})$ becomes equivalent to minimizing $\varphi(\eta)$ defined for all $\eta \in \mathbb{R}$ by

$$\varphi(\eta) = \beta_k |w_{k,j} + \eta| + \eta \epsilon_{t-1,k,j} + \frac{\eta^2 \Lambda^2}{2}.$$

Let η^* denote the minimizer of $\varphi(\eta)$. If $w_{t-1,k,j} + \eta^* = 0$, then the subdifferential of $|w_{t-1,k,j} + \eta|$ at η^* is the set $\{\nu: \nu \in [-1, +1]\}$. Thus, in that case, the subdifferential $\partial\varphi(\eta^*)$ contains 0 iff there exists $\nu \in [-1, +1]$ such that

$$\beta_k \nu + \epsilon_{t-1,k,j} + \eta^* \Lambda^2 = 0 \Leftrightarrow w_{t-1,k,j} \Lambda^2 - \epsilon_{t-1,k,j} = \beta_k \nu.$$

The condition is therefore equivalent to

$$|w_{t-1,k,j} \Lambda^2 - \epsilon_{t-1,k,j}| \leq \beta_k.$$

If $w_{t-1,k,j} + \eta^* > 0$, then φ is differentiable at η^* and $\varphi'(\eta^*) = 0$, that is

$$\beta_k + \epsilon_{t-1,k,j} + \eta^* \Lambda^2 = 0 \Leftrightarrow \eta^* = \frac{1}{\Lambda^2} [-\beta_k - \epsilon_{t-1,k,j}].$$

In view of that expression, the condition $w_{t-1,k,j} + \eta^* > 0$ is equivalent to

$$w_{t-1,k,j} \Lambda^2 - \epsilon_{t-1,k,j} > \beta_k.$$

Similarly, if $w_{t-1,k,j} + \eta^* < 0$, φ is differentiable at η^* and $\varphi'(\eta^*) = 0$, which gives

$$\eta^* = \frac{1}{\Lambda^2} [\beta_k - \epsilon_{t-1,k,j}].$$

Figure 1 shows the pseudocode of our algorithm using the closed-form solution for the step size just presented.

D. Convergence analysis

In this section, we give convergence guarantees for both versions of the StructMaxent algorithm.

Theorem 3. *Let $(\mathbf{w}_t)_t$ be the sequence of parameter vectors generated by StructMaxent1 or StructMaxent2. Then, $(\mathbf{w}_t)_t$ converges to the optimal solution \mathbf{w}^* of (6).*

Proof. We begin with the proof for StructMaxent2. Our proof is based on Lemma 19 of (Dudík et al., 2007), which implies that it suffices to show that $F(\mathbf{w}_t)$ admits a finite limit and that there exists a sequence \mathbf{u}_t such that $R(\mathbf{u}_t, \mathbf{w}_t) \rightarrow 0$ as $t \rightarrow \infty$, where R is some auxiliary function. A function R is said to be *auxiliary* if

$$\begin{aligned} R(\mathbf{u}, \mathbf{w}) = & I_C(\mathbf{u}) + \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|_1 + \mathbf{w} \cdot \mathbb{E}_S [\Phi] + \mathbf{w} \cdot \mathbf{u} \\ & + B(\mathbf{u} \parallel \mathbb{E}_{\mathbf{P}_{\mathbf{w}}} [\Phi]), \end{aligned}$$

where B is a Bregman divergences. We will use the Bregman divergence based on the squared difference:

$$B(\mathbf{u} \parallel \mathbb{E}_{\mathbf{P}_{\mathbf{w}}} [\Phi]) = \frac{\|\mathbf{u} - \mathbb{E}_{\mathbf{P}_{\mathbf{w}}} [\Phi]\|_2^2}{2\Lambda^2}.$$

Let $g_0(\mathbf{u}) = I_C(\mathbf{u}) + \mathbf{w} \cdot \mathbf{u}$ and observe that using the same arguments as in the proof of Theorem 1, we can write

$$\begin{aligned} g_0^*(\mathbf{r}) = & \sup_{\mathbf{u} \in C} ((\mathbf{r} - \mathbf{w}) \cdot \mathbf{u} - I_C(\mathbf{u})) \\ = & (\mathbf{r} - \mathbf{w}) \cdot \mathbb{E}_S [\Phi] + \sum_{k=1}^p \beta_k \|\mathbf{r}_k - \mathbf{w}_k\|_1. \end{aligned}$$

Similarly, if $f_0(\mathbf{u}) = B(\mathbf{u} \parallel \mathbb{E}_{\rho_{\mathbf{w}}}[\Phi])$, then

$$\begin{aligned} f_0^*(\mathbf{r}) &= \sup_{\mathbf{u}} (\mathbf{r} \cdot \mathbf{u} - B(\mathbf{u} \parallel \mathbb{E}_{\rho_{\mathbf{w}}}[\Phi])) \\ &= \frac{\Lambda^2 \|\mathbf{r}\|_2}{2} + \mathbf{r} \cdot \mathbb{E}_{\rho_{\mathbf{w}}}[\Phi]. \end{aligned}$$

Therefore, applying Theorem 5 with $A = I$, we obtain

$$\begin{aligned} \inf_{\mathbf{u}} R(\mathbf{u}, \mathbf{w}_t) &= \sup_{\mathbf{r}} \left(-\frac{\Lambda^2 \|\mathbf{r}\|_2}{2} - \mathbf{r} \cdot \mathbb{E}_{\rho_{\mathbf{w}_t}}[\Phi] - \mathbf{r} \cdot \mathbb{E}_S[\Phi] \right. \\ &\quad \left. + \sum_{k=1}^p \beta_k (\|\mathbf{w}_{t,k}\|_1 - \|\mathbf{r}_k + \mathbf{w}_{t,k}\|_1) \right), \end{aligned}$$

and we define \mathbf{u}_t to be the solution of this optimization problem, which, in view of Theorem 5, does exist. We will now argue that $R(\mathbf{u}_t, \mathbf{w}_t) \rightarrow 0$ as $t \rightarrow \infty$. Note that

$$\begin{aligned} R(\mathbf{u}_t, \mathbf{w}_t) &= -\sum_{k=1}^p \sum_{j=1}^{N_k} \inf_r \left(\frac{\lambda^2 r}{2} + r(\mathbb{E}_{\rho_{\mathbf{w}_t}}[\Phi_{k,j}] - \mathbb{E}_S[\Phi_{k,j}]) \right. \\ &\quad \left. + \beta_k |w_{t,k,j}| - \beta_k |r + w_{t,k,j}| \right). \end{aligned}$$

Recall that, by definition of StructMaxent2, the following holds for all $(k, j) \in [1, p] \times [1, N_k]$:

$$\begin{aligned} F(\mathbf{w}_t) - F(\mathbf{w}_{t+1}) & \quad (22) \\ &\geq -\inf_r \left(\frac{\Lambda^2 r}{2} + r(\mathbb{E}_{\rho_{\mathbf{w}_t}}[\Phi_{k,j}] - \mathbb{E}_S[\Phi_{k,j}]) \right. \\ &\quad \left. + \beta_k |w_{t,k,j}| - \beta_k |r + w_{t,k,j}| \right) \\ &\geq 0, \end{aligned}$$

where the last inequality follows by taking $r = 0$. Therefore, to complete the proof, it suffices to show that $\lim_{t \rightarrow \infty} F(\mathbf{w}_t)$ is finite, since then $F(\mathbf{w}_t) - F(\mathbf{w}_{t+1}) \rightarrow 0$ and $R(\mathbf{u}_t, \mathbf{w}_t) \rightarrow 0$. By (22), $F(\mathbf{w}_t)$ is decreasing and it suffices to show that $F(\mathbf{w}_t)$ is bounded below. This is an immediate consequence of the feasibility of the optimization problem $\inf_{\mathbf{w}} F(\mathbf{w})$ which was established in Section 2.2 and the proof for StructMaxent2 is now complete.

The proof for StructMaxent1 requires the use a different Bregman divergence B defined as follows:

$$B(\mathbf{u} \parallel \mathbb{E}_{\rho_{\mathbf{w}}}[\Phi]) = \sum_{k=1}^p \sum_{j=1}^{N_k} D_0(\varphi_{kj}(\mathbf{u}) \parallel \varphi_{kj}(\mathbb{E}_{\rho_{\mathbf{w}}}[\Phi])),$$

where D_0 is unnormalized relative entropy, $\varphi_{kj}(\mathbf{u}) = ((\Lambda - u_{k,j}), (\Lambda + u_{k,j}))$ and $\|\mathbf{u}\|_{\infty} \leq \Lambda$. The rest of the argument remains the same. \square

E. Bounds on Rademacher complexities

In this section, we give the proof of the upper bounds on Rademacher complexities given in (15):

$$\begin{aligned} \mathfrak{R}_m(H_k^{\text{mono}}) &\leq \sqrt{\frac{2k \log d}{m}} \\ \mathfrak{R}_m(H_k^{\text{trees}}) &\leq \sqrt{\frac{(4k+2) \log_2(d+2) \log(m+1)}{m}}. \end{aligned}$$

The first inequality is an immediate consequence of Massart's lemma, which states that

$$\frac{1}{m} \mathbb{E}_{\sigma} \left[\sup_{\mathbf{x} \in A} \sum_{i=1}^m \sigma_i x_i \right] \leq \frac{r \sqrt{2 \log |A|}}{m},$$

where $A \subset \mathbb{R}^n$ is a finite set, $r = \max_{\mathbf{x} \in A} \|\mathbf{x}\|_2$ and σ_i s are Rademacher random variables. If we take A to be the image of the sample under H_k^{mono} then $|A| \leq |H_k^{\text{mono}}| \leq d^k$. Moreover, if the features in H_1^{mono} are normalized to belong to $[-1, 1]$ then $\Lambda = 1$ and $r = \sqrt{m}$. Combining these results with Massart's lemma leads to the desired bound.

Now we derive the second bound of (15). Since each binary decision tree in H_k^{trees} , can be viewed as a binary classifier, Massart's lemma yields that

$$\mathfrak{R}_m(H_k^{\text{trees}}) \leq \sqrt{\frac{2 \log \Pi_{H_k^{\text{trees}}}(m)}{m}},$$

where $\Pi_{H_k^{\text{trees}}}(m)$ is the growth function of H_k^{trees} . We use Sauer's lemma to bound the growth function: $\Pi_{H_k^{\text{trees}}}(m) \leq (em)^{\text{VC-dim}(H_k^{\text{trees}})}$. For the family of binary decision trees in dimension d it is known that $\text{VC-dim}(H_k^{\text{trees}}) \leq (2k+1) \log_2(d+2)$ (Mansour, 1997) and the desired bound follows.