# Convex Calibrated Surrogates for Hierarchical Classification
## Supplementary Material

## A. Additional Notation and Setup

Let $\mu$ be the marginal distribution induced by $D$ over $\mathcal{X}$, and let $\mathbf{p}(x)$ be the distribution over $[n]$ conditioned on $X = x$. For every function $\ell : [n] \times [k] \to \mathbb{R}_+$ and $t \in [k]$ let $\boldsymbol{\ell}_t = [\ell(1, t), \ldots, \ell(n, t)]^\top \in \mathbb{R}_+^n$. For every surrogate $\psi : [n] \times \mathbb{R}^d \to \mathbb{R}_+$ let $\boldsymbol{\psi} : \mathbb{R}^d \to \mathbb{R}_+^n$ be a vector function such that $\psi_y(\mathbf{u}) = \psi(y, \mathbf{u})$ for $y \in [n], \mathbf{u} \in \mathbb{R}^d$. For any integer $d' \in \mathbb{Z}_+$ and pair of vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{d'}$, their inner product is denoted as $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^{d'} u_i v_i$. For a vector $\mathbf{u} \in \mathbb{R}^n$ and a positive integer $a \le n$, the vector $\mathbf{u}\big|_{1:a} \in \mathbb{R}^a$ gives the first $a$ components of $\mathbf{u}$.

Define the conditional regrets $\mathbf{R}_{\mathbf{p}}^{\ell^H}, \mathbf{R}_{\mathbf{p}}^{\ell^{?,n}}$ and $\mathbf{R}_{\mathbf{p}}^{\psi}$ as the regrets incurred for a singleton instance space $\mathcal{X}$, with conditional probability $\mathbf{p} \in \Delta_n$. In particular, we have that

$$
\begin{aligned}
\mathbf{R}_{\mathbf{p}}^{\ell^H}[\widehat{y}] &= \langle \mathbf{p}, \boldsymbol{\ell}_{\widehat{y}}^H \rangle - \inf_{y' \in [n]} \langle \mathbf{p}, \boldsymbol{\ell}_{y'}^H \rangle, & \forall \widehat{y} \in [n] \\
\mathbf{R}_{\mathbf{p}}^{\ell^{?,n}}[\widehat{y}] &= \langle \mathbf{p}, \boldsymbol{\ell}_{\widehat{y}}^{?,n} \rangle - \inf_{y' \in [n] \cup \{\bot\}} \langle \mathbf{p}, \boldsymbol{\ell}_{y'}^{?,n} \rangle, & \forall \widehat{y} \in [n] \cup \{\bot\} \\
\mathbf{R}_{\mathbf{p}}^{\psi}[\mathbf{u}] &= \langle \mathbf{p}, \boldsymbol{\psi}(\mathbf{u}) \rangle - \inf_{\mathbf{u}' \in \mathbb{R}^d} \langle \mathbf{p}, \boldsymbol{\psi}(\mathbf{u}') \rangle, & \forall \mathbf{u} \in \mathbb{R}^d .
\end{aligned}
$$

Let $\mu$ be the marginal distribution induced by $D$ over $\mathcal{X}$, and let $\mathbf{p}(x)$ be the distribution over $[n]$ conditioned on $X = x$. Then we have by linearity of expectation that,

$$
\begin{aligned}
\mathbf{R}_D^{\ell^H}[g] &= \mathbf{E}_{X \sim \mu} \mathbf{R}_{\mathbf{p}(X)}^{\ell^H}[g(X)] & (4) \\
\mathbf{R}_D^{\ell^{?,n}}[g'] &= \mathbf{E}_{X \sim \mu} \mathbf{R}_{\mathbf{p}(X)}^{\ell^{?,n}}[g'(X)] & (5) \\
\mathbf{R}_D^{\psi}[\mathbf{f}] &= \mathbf{E}_{X \sim \mu} \mathbf{R}_{\mathbf{p}(X)}^{\psi}[\mathbf{f}(X)] . & (6)
\end{aligned}
$$

For all $0 \le j \le h$, define $\ell^{H,j} : [n_j] \times [n_j] \to \mathbb{R}_+$ as simply the restriction of $\ell^H$ to $[n_j] \times [n_j]$.

## B. Proofs

### B.1. Proof of Theorem 1

**Theorem.** *Let $H = ([n], E, W)$ and let $\ell^H : [n] \times [n] \to \mathbb{R}_+$ be the tree-distance loss for the tree $H$. Let $\mathbf{p} \in \Delta_n$, and $y \in [n]$. Then there exists a $g^* : \mathcal{X} \to [n]$ such that for all $x \in \mathcal{X}$ the following holds:*

*(a) $S_{g^*(x)}(\mathbf{p}(x)) \ge \frac{1}{2}$*

*(b) $S_y(\mathbf{p}(x)) \le \frac{1}{2}, \forall y \in C(g^*(x))$ .*

*Also, $g^*$ is a Bayes optimal classifier for the tree distance loss, i.e.*

$$
\mathbf{R}_D^{\ell^H}[g^*] = 0 .
$$

*Proof.* We shall simply show for all $\mathbf{p} \in \Delta_n$, there exists a $y^* \in [n]$ such that

$$
\begin{aligned}
S_{y^*}(\mathbf{p}) &\ge \frac{1}{2} & (7) \\
S_y(\mathbf{p}) &\le \frac{1}{2}, & \forall y \in C(y^*), & (8)
\end{aligned}
$$

and is such that

$$\langle \mathbf{p}, \boldsymbol{\ell}_{y^*}^H \rangle = \min_{y \in [n]} \langle \mathbf{p}, \boldsymbol{\ell}_y^H \rangle \ .$$

This would imply $\mathbf{R}_{\mathbf{p}}^{\ell^H}[y^*] = 0$. The theorem then simply follows from linearity of expectation using Equation 4.

Let $\mathbf{p} \in \Delta_n$. We construct a $y^* \in [n]$ satisfying Equations 7 and 8 in the following way. We start at the root node, which always satisfies Equation 7, and keep on moving to the child of the current node that satisfies Equation 7, and terminate when we reach a leaf node, or a node where all of its children fail Equation 7. Clearly the resulting node, $y^*$, satisfies both Equations 7 and 8.

Now we show that $y^*$ indeed minimizes $\langle \mathbf{p}, \boldsymbol{\ell}_y^H \rangle$ over $y \in [n]$.

Let $y' \in \operatorname{argmin}_t \langle \mathbf{p}, \boldsymbol{\ell}_t^H \rangle$. If $y' = y^*$ we are done, hence assume $y' \neq y^*$.

**Case 1:** $y' \notin D(y^*)$

$$
\begin{aligned}
\langle \mathbf{p}, \boldsymbol{\ell}_{y'}^H \rangle - \langle \mathbf{p}, \boldsymbol{\ell}_{y^*}^H \rangle &= \sum_{y \in D(y^*)} p_y(\ell^H(y, y') - \ell^H(y, y^*)) + \sum_{y \in [n] \setminus D(y^*)} p_y(\ell^H(y, y') - \ell^H(y, y^*)) \\
&= \sum_{y \in D(y^*)} p_y(\ell^H(y^*, y')) + \sum_{y \in [n] \setminus D(y^*)} p_y(\ell^H(y, y') - \ell^H(y, y^*)) \\
&\geq \sum_{y \in D(y^*)} p_y(\ell^H(y^*, y')) + \sum_{y \in [n] \setminus D(y^*)} p_y(-\ell^H(y', y^*))) \\
&= \ell^H(y', y^*)(2S_{y^*}(\mathbf{p}) - 1) \\
&\geq 0
\end{aligned}
$$

**Case 2:** $y' \in D(y^*) \setminus C(y^*)$

Let $\hat{y}$ be the child of $y^*$ that is the ancestor of $y'$. Hence we have $S_{\hat{y}}(\mathbf{p}) \leq \frac{1}{2}$.

$$
\begin{aligned}
\langle \mathbf{p}, \boldsymbol{\ell}_{y'}^H \rangle - \langle \mathbf{p}, \boldsymbol{\ell}_{y^*}^H \rangle &= \sum_{y \in D(\hat{y})} p_y(\ell^H(y, y') - \ell^H(y, y^*)) + \sum_{y \in [n] \setminus D(\hat{y})} p_y(\ell^H(y, y') - \ell^H(y, y^*)) \\
&= \sum_{y \in D(\hat{y})} p_y(\ell^H(y, y') - \ell^H(y, y^*)) + \sum_{y \in [n] \setminus D(\hat{y})} p_y(\ell^H(y^*, y')) \\
&\geq \sum_{y \in D(\hat{y})} p_y(-\ell^H(y^*, y')) + \sum_{y \in [n] \setminus D(\hat{y})} p_y(\ell^H(y^*, y')) \\
&= \ell^H(y', y^*)(1 - 2S_{\hat{y}}(\mathbf{p})) \\
&\geq 0
\end{aligned}
$$

**Case 3:** $y' \in C(y^*)$

$$
\begin{aligned}
\langle \mathbf{p}, \boldsymbol{\ell}_{y'}^H \rangle - \langle \mathbf{p}, \boldsymbol{\ell}_{y^*}^H \rangle &= \sum_{y \in D(y')} p_y(\ell^H(y, y') - \ell^H(y, y^*)) + \sum_{y \in [n] \setminus D(y')} p_y(\ell^H(y, y') - \ell^H(y, y^*)) \\
&= \sum_{y \in D(y')} p_y(-\ell^H(y', y^*)) + \sum_{y \in [n] \setminus D(y')} p_y(\ell^H(y', y^*)) \\
&= \ell^H(y', y^*)(1 - 2S_{y'}(\mathbf{p})) \\
&\geq 0
\end{aligned}
$$

Putting all three cases together we have

$$\langle \mathbf{p}, \boldsymbol{\ell}_{y^*}^H \rangle \leq \langle \mathbf{p}, \boldsymbol{\ell}_{y'}^H \rangle = \min_{y \in [n]} \langle \mathbf{p}, \boldsymbol{\ell}_y^H \rangle \ .$$

$\square$

### B.2. Proof of Lemma 3

We first give the proof of a stronger version of Lemma 3.

**Lemma 8.** *For all* $\mathbf{p} \in \Delta_n, \mathbf{u} \in \mathbb{R}^d$ *we have*

$$\mathbf{R}_{\mathbf{p}}^{\ell^H}[\Upsilon^{\mathrm{cas}}(\mathbf{u})] \leq \sum_{j=1}^{h} \gamma_j(\mathbf{u}_j) \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\ell^{?,n_j}}[\Upsilon_j^?(\mathbf{u}_j)],$$

*where* $\gamma_j(\mathbf{u}_j) = \begin{cases} 2\alpha_j & \textit{if } \Upsilon_j^?(\mathbf{u}_j) \neq \perp \\ 2\beta_j & \textit{if } \Upsilon_j^?(\mathbf{u}_j) = \perp \end{cases}.$

*Proof.* For all $j \in [h]$, we will first prove a bound relating the tree-distance regret at level $j$, with the tree-distance regret at level $j-1$ and abstain loss regret at level $j$ as follows:

$$\mathbf{R}_{A^j(\mathbf{p})}^{\ell^{H,j}}[\Upsilon_j^{\mathrm{cas}}(\mathbf{u})] \quad \leq \quad \mathbf{R}_{A^{j-1}(\mathbf{p})}^{\ell^{H,j-1}}[\Upsilon_{j-1}^{\mathrm{cas}}(\mathbf{u})] + \gamma_j(\mathbf{u}_j) \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\ell^{?,n_j}}[\Upsilon_j^?(\mathbf{u}_j)] \ .$$

The theorem would simply follow from applying such a bound recursively and observing that $\mathbf{R}_{A^0(\mathbf{p})}^{\ell^{H,0}}[\Upsilon_0^{\mathrm{cas}}(\mathbf{u})] = 0$.

One observation of the tree-distance loss that will be often of use in the proof is the following:

$$\ell^H(y, y') - \ell^H(P(y), y') = \begin{cases} -\ell^H(y, P(y)) & \text{if } y' \in D(y) \\ \ell^H(y, P(y)) & \text{otherwise} \end{cases}$$

The details of the proof follows: Fix $j \in [h], \mathbf{u} \in \mathbb{R}^d, \mathbf{p} \in \Delta_n$.

Let $y_j^* = \operatorname{argmin}_{y \in [n_j]} \mathbf{R}_{A^j(\mathbf{p})}^{\ell^{H,j}}[y]$.

**Case 1:** $\Upsilon_j^?(\mathbf{u}_j) \neq \perp$

$$\begin{aligned} \mathbf{R}_{A^j(\mathbf{p})}^{\ell^{H,j}}[\Upsilon_j^{\mathrm{cas}}(\mathbf{u})] &= \sum_{y=1}^{n_j} A_y^j(\mathbf{p})(\ell^H(y, \Upsilon_j^{\mathrm{cas}}(\mathbf{u})) - \ell^H(y, y_j^*)) \\ &\leq \ell^H(y_j^*, \Upsilon_j^{\mathrm{cas}}(\mathbf{u}))(1 - 2A_{\Upsilon_j^{\mathrm{cas}}(\mathbf{u})}^j(\mathbf{p})) \end{aligned} \tag{9}$$

We also have,

$$\begin{aligned} \mathbf{R}_{A^j(\mathbf{p})}^{\ell^{?,n_j}}[\Upsilon_j^?(\mathbf{u}_j)] &= 1 - A_{\Upsilon_j^?(\mathbf{u}_j)}^j(\mathbf{p}) - \min_{y \in [n_j] \cup \{\perp\}} \langle A^j(\mathbf{p}), \boldsymbol{\ell}_y^{?,n_j} \rangle \\ &\geq 1 - A_{\Upsilon_j^?(\mathbf{u}_j)}^j(\mathbf{p}) - \langle A^j(\mathbf{p}), \boldsymbol{\ell}_\perp^{?,n_j} \rangle \\ &= \frac{1}{2} - A_{\Upsilon_j^?(\mathbf{u}_j)}^j(\mathbf{p}) \\ &= \frac{1}{2} - A_{\Upsilon_j^{\mathrm{cas}}(\mathbf{u})}^j(\mathbf{p}) \ . \end{aligned} \tag{10}$$

The last inequality above follows because if $\Upsilon_j^?(\mathbf{u}_j) \neq \perp$, then $\Upsilon_j^{\mathrm{cas}}(\mathbf{u}) = \Upsilon_j^?(\mathbf{u}_j)$.

Putting Equations 9 and 10 together, we get

$$
\begin{aligned}
\mathbf{R}^{\ell^{H,j}}_{A^j(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_j(\mathbf{u})] &\leq 2\ell^H(y_j^*, \Upsilon^{\mathrm{cas}}_j(\mathbf{u})) \cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\Upsilon^?_j(\mathbf{u}_j)] \\
&\leq 2\alpha_j \cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\Upsilon^?_j(\mathbf{u}_j)]
\end{aligned}
\tag{11}
$$

**Case 2:** $\Upsilon^?_j(\mathbf{u}_j) = \perp$

In this case $\Upsilon^{\mathrm{cas}}_j(\mathbf{u}) = \Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})$, and hence $\mathrm{lev}(\Upsilon^{\mathrm{cas}}_j(\mathbf{u})) \leq j-1$.

We now have,

$$
\begin{aligned}
\langle A^j(\mathbf{p}), \boldsymbol{\ell}^{H,j}_{\Upsilon^{\mathrm{cas}}_j(\mathbf{u})}\rangle - \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{\Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})}\rangle &= \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{H,j}_{\Upsilon^{\mathrm{cas}}_j(\mathbf{u})}\rangle - \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{\Upsilon^{\mathrm{cas}}_j(\mathbf{u})}\rangle \\
&= \sum_{y \in N_{=j}} S_y(\mathbf{p})\left(\ell^H(y, \Upsilon^{\mathrm{cas}}_j(\mathbf{u})) - \ell^H(P(y), \Upsilon^{\mathrm{cas}}_j(\mathbf{u}))\right) \\
&= \sum_{y \in N_{=j}} S_y(\mathbf{p})\ell^H(y, P(y))
\end{aligned}
\tag{12}
$$

For ease of analysis, we divide case 2, further into two sub-cases.

**Case 2a:** $\mathrm{lev}(y_j^*) < j$

$$
\begin{aligned}
\langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{y_{j-1}^*}\rangle - \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{H,j}_{y_j^*}\rangle &= \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{y_{j-1}^*}\rangle - \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{y_j^*}\rangle \\
&\quad + \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{y_j^*}\rangle - \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{H,j}_{y_j^*}\rangle \\
&\leq \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{y_j^*}\rangle - \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{H,j}_{y_j^*}\rangle \\
&= \sum_{y \in N_{=j}} S_y(\mathbf{p})(\ell^H(P(y), y_j^*) - \ell^H(y, y_j^*)) \\
&= \sum_{y \in N_{=j}} S_y(\mathbf{p})(-\ell^H(y, P(y)))
\end{aligned}
\tag{13}
$$

Adding, Equation 12 and 13, we get

$$
\mathbf{R}^{\ell^{H,j}}_{A^j(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_j(\mathbf{u})] \leq \mathbf{R}^{\ell^{H,j-1}}_{A^{j-1}(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})]
\tag{14}
$$

**Case 2b:** $\mathrm{lev}(y_j^*) = j$

$$
\begin{aligned}
\langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{y_{j-1}^*}\rangle - \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^H_{y_j^*}\big|_{[1:n_{j-1}]}\rangle &\leq \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^{H,j-1}_{P(y_j^*)}\rangle - \langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^H_{y_j^*}\big|_{[1:n_{j-1}]}\rangle \\
&= \sum_{y \in N_{\leq j-1}} A_y^{j-1}(\mathbf{p})(\ell^H(y, P(y_j^*)) - \ell^H(y, y_j^*)) \\
&= \sum_{y \in N_{\leq j-1}} A_y^{j-1}(\mathbf{p})(-\ell^H(y_j^*, P(y_j^*))) \\
&= -\ell^H(y_j^*, P(y_j^*))
\end{aligned}
\tag{15}
$$

Also,

$$
\begin{aligned}
\langle A^{j-1}(\mathbf{p}), \boldsymbol{\ell}^H_{y_j^*}\big|_{[1:n_{j-1}]}\rangle - \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{H,j}_{y_j^*}\rangle &= \sum_{y \in N_{=j}} S_y(\mathbf{p})(\ell^H(P(y), y_j^*) - \ell^H(y, y_j^*)) \\
&= \sum_{y \in N_{=j}\setminus\{y_j^*\}} S_y(\mathbf{p})(-\ell^H(y, P(y))) + S_{y_j^*}(\mathbf{p})(\ell^H(y_j^*, P(y_j^*))) .
\end{aligned}
\tag{16}
$$

Adding Equations 12, 15 and 16, we get

$$
\begin{aligned}
\mathbf{R}^{\ell^{H,j}}_{A^j(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_j(\mathbf{u})] &\leq \mathbf{R}^{\ell^{H,j-1}}_{A^{j-1}(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})] + (2S_{y^*_j}(\mathbf{p}) - 1)\cdot \ell^H(y^*_j, P(y^*_j)) \\
&\leq \mathbf{R}^{\ell^{H,j-1}}_{A^{j-1}(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})] + (2S_{y^*_j}(\mathbf{p}) - 1)\cdot \beta_j .
\end{aligned}
\tag{17}
$$

Inequality 17 follows because by the definitions of $y^*_j$ and Theorem 1 , we have $S_{y^*_j}(\mathbf{p}) \geq \frac{1}{2}$.

Also, we have that

$$
\begin{aligned}
\mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\Upsilon^?_j(\mathbf{u}_j)] &= \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\bot] \\
&= \frac{1}{2} - \min_{y \in [n] \cup \{\bot\}} \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{?,n_j}_y \rangle \\
&\geq \frac{1}{2} - \langle A^j(\mathbf{p}), \boldsymbol{\ell}^{?,n_j}_{y^*_j} \rangle \\
&= \frac{1}{2} - (1 - S_{y^*_j}(\mathbf{p})) \\
&= S_{y^*_j}(\mathbf{p}) - \frac{1}{2} .
\end{aligned}
\tag{18}
$$

Putting Equations 17 and 18 together, we have that

$$
\mathbf{R}^{\ell^{H,j}}_{A^j(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_j(\mathbf{u})] \leq \mathbf{R}^{\ell^{H,j-1}}_{A^{j-1}(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})] + 2\beta_j \cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\Upsilon^?_j(\mathbf{u}_j)].
\tag{19}
$$

Putting the results for case 1, case 2a and case 2b, from Equations 11, 14 and 19 respectively, we have

$$
\mathbf{R}^{\ell^{H,j}}_{A^j(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_j(\mathbf{u})] \leq \mathbf{R}^{\ell^{H,j-1}}_{A^{j-1}(\mathbf{p})}[\Upsilon^{\mathrm{cas}}_{j-1}(\mathbf{u})] + \gamma_j(\mathbf{u}_j)\cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\Upsilon^?_j(\mathbf{u}_j)] .
$$

$\square$

Now the proof of Lemma 3 follows from certain simple considerations.

**Lemma.** *For any distribution $D$ over $\mathcal{X} \times [n]$, let $A^j(D)$ be the distribution over $\mathcal{X} \times [n_j]$ given by the distribution of $(X, \mathrm{anc}_j(Y))$ with $(X, Y) \sim D$. For all $j \in [h]$, let $\mathbf{f}_j : \mathcal{X} \to \mathbb{R}^{d_j}$ be such that $\mathbf{f}(x) = [\mathbf{f}_1(x)^\top, \ldots, \mathbf{f}_h(x)^\top]^\top$. Then for all distributions $D$ over $\mathcal{X} \times [n]$ and all functions $\mathbf{f} : \mathcal{X} \to \mathbb{R}^d$*

$$
\mathbf{R}^{\ell^H}_D[\Upsilon^{\mathrm{cas}} \circ \mathbf{f}] \leq \sum_{j=1}^h 2\alpha_j \cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(D)}[\Upsilon^?_j \circ \mathbf{f}_j] .
$$

*Proof.* Using Lemma 8 and by the observation that $\beta_j \leq \alpha_j$, we have for all $\mathbf{p} \in \Delta_n, \mathbf{u} \in \mathbb{R}^d$ that

$$
\mathbf{R}^{\ell^H}_{\mathbf{p}}[\Upsilon^{\mathrm{cas}}(\mathbf{u})] \leq \sum_{j=1}^h 2\alpha_j \cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p})}[\Upsilon^?_j(\mathbf{u}_j)] .
$$

Let $\mathbf{f} : \mathcal{X} \to \mathbb{R}^d$ be a function. Then for all $x \in \mathcal{X}$,

$$
\mathbf{R}^{\ell^H}_{\mathbf{p}(x)}[\Upsilon^{\mathrm{cas}}(\mathbf{f}(x))] \leq \sum_{j=1}^h 2\alpha_j \cdot \mathbf{R}^{\ell^{?,n_j}}_{A^j(\mathbf{p}(x))}[\Upsilon^?_j(\mathbf{u}_j)] .
$$

Observe that the the marginal distribution over $\mathcal{X}$ for $A^j(D)$ is exactly the same as for $D$, while the conditional probability distribution for the distribution $A^j(D)$ at $x$ is exactly equal to $A^j(\mathbf{p}(x))$. The Lemma now immediately follows from linearity of expectation. $\square$

### B.3. Proof of Theorem 4

**Theorem.** *For all $j \in [h]$, let $\psi^j : [n_j] \times \mathbb{R}^{d_j}$ and $\Upsilon_j^? : \mathbb{R}^{d_j} \to n_j$ be such that for all $\mathbf{f}_j : \mathcal{X} \to \mathbb{R}^{d_j}$, and all distributions $D$ over $\mathcal{X} \times [n_j]$ we have*

$$\mathbf{R}_D^{\ell^?, n_j}[\Upsilon_j^? \circ \mathbf{f}_j] \leq C \cdot \mathbf{R}_D^{\psi^j}[\mathbf{f}_j],$$

*for some constant $C > 0$. Then for all $\mathbf{f} : \mathcal{X} \to \mathbb{R}^d$ and distributions $D$ over $\mathcal{X} \times [n]$,*

$$\mathbf{R}_D^{\ell^H}[\Upsilon^{\mathrm{cas}} \circ \mathbf{f}] \leq 2\alpha_h C \cdot \mathbf{R}_D^{\psi^{\mathrm{cas}}}[\mathbf{f}].$$

*Proof.* Fix $\mathbf{u} \in \mathbb{R}^d, \mathbf{p} \in \Delta_n$. From Lemma 3, we have that

$$
\begin{aligned}
\mathbf{R}_{\mathbf{p}}^{\ell^H}[\Upsilon^{\mathrm{cas}}(\mathbf{u})] &\leq \sum_{j=1}^h 2\alpha_j \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\ell^?, n_j}[\Upsilon_j^?(\mathbf{u}_j)] \\
&\leq 2\alpha_h \cdot \sum_{j=1}^h \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\ell^?, n_j}[\Upsilon_j^?(\mathbf{u}_j)] \\
&\leq 2\alpha_h C \cdot \sum_{j=1}^h \mathbf{R}_{A^j(\mathbf{p})}^{\psi^j}[\mathbf{u}_j] \\
&= 2\alpha_h C \cdot \mathbf{R}_{\mathbf{p}}^{\psi^{\mathrm{cas}}}[\mathbf{u}].
\end{aligned}
$$

The proof now simply follows from linearity of expectation. $\qquad\square$

### B.4. Proof of Theorem 6

While Theorem 2 from Ramaswamy et al. (2015), gives an excess risk bound for the abstain loss excess risk in terms of the OvA-surrogate risk, one can easily get a more refined bound as well from the results of Ramaswamy et al. (2015).

**Lemma 9** ((Ramaswamy et al., 2015)). *Let $\tau \in (-1, 1)$. For all $\mathbf{u} \in \mathbb{R}^n, \mathbf{p} \in \Delta_n$, and $A = \mathbf{1}(\Upsilon_\tau^{\mathrm{OvA}, n}(\mathbf{u}) = n+1)$. Then for all $\mathbf{p} \in \Delta_n$*

$$\mathbf{R}_{\mathbf{p}}^{\ell^?, n}[\Upsilon_\tau^{\mathrm{OvA}, n}(\mathbf{u})] \leq \left( \frac{\mathbf{1}(\Upsilon^{\mathrm{OvA}, n}(\mathbf{u}) = \perp)}{2(1-\tau)} + \frac{\mathbf{1}(\Upsilon^{\mathrm{OvA}, n}(\mathbf{u}) \neq \perp)}{2(1+\tau)} \right) \mathbf{R}_{\mathbf{p}}^{\psi^{\mathrm{OvA}, n}}[\mathbf{u}].$$

We are now ready to prove Theorem 6.

**Theorem.** *For $1 \leq j \leq h$, let $\tau_j = \frac{\alpha_j - \beta_j}{\alpha_j + \beta_j}$. Let the component surrogates and predictors of $\psi^{\mathrm{cas}}$ and $\Upsilon^{\mathrm{cas}}$ be $\psi^j = \psi^{\mathrm{OvA}, n_j}$ and $\Upsilon^j = \Upsilon_{\tau_j}^{\mathrm{OvA}, n_j}$. Then, for all distributions $D$ and functions $\mathbf{f} : \mathcal{X} \to \mathbb{R}^d$,*

$$\mathbf{R}_D^{\ell^H}[\Upsilon^{\mathrm{cas}} \circ \mathbf{f}] \leq \frac{1}{2} \max_{j \in [h]}(\alpha_j + \beta_j) \cdot \mathbf{R}_D^{\psi^{\mathrm{cas}}}[\mathbf{f}]$$

*Proof.* Let $\mathbf{u} \in \mathbb{R}^d, \mathbf{p} \in \Delta_n$, From Lemmas 8 and 9, we have that

$$
\begin{aligned}
\mathbf{R}_{\mathbf{p}}^{\ell^H}[\Upsilon^{\mathrm{cas}}(\mathbf{u})] &\leq \sum_{j=1}^h \gamma_j(\mathbf{u}_j) \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\ell^?, n_j}[\Upsilon_j^?(\mathbf{u}_j)] \\
&\leq \sum_{j=1}^h \gamma_j(\mathbf{u}_j) \left( \frac{\mathbf{1}(\Upsilon^{\mathrm{OvA}, n_j}(\mathbf{u}_j) = \perp)}{2(1-\tau_j)} + \frac{\mathbf{1}(\Upsilon^{\mathrm{OvA}, n_j}(\mathbf{u}_j) \neq \perp)}{2(1+\tau_j)} \right) \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\psi^{\mathrm{OvA}, n_j}}[\mathbf{u}_j] \\
&= \sum_{j=1}^h \left( \frac{\beta_j \cdot \mathbf{1}(\Upsilon^{\mathrm{OvA}, n_j}(\mathbf{u}_j) = \perp)}{(1-\tau_j)} + \frac{\alpha_j \cdot \mathbf{1}(\Upsilon^{\mathrm{OvA}, n_j}(\mathbf{u}_j) \neq \perp)}{(1+\tau_j)} \right) \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\psi^{\mathrm{OvA}, n_j}}[\mathbf{u}_j]
\end{aligned}
$$

For each $j \in [h]$, the coefficients of both the terms within parantheses (i.e. $\frac{\alpha_j}{1+\tau_j}$ and $\frac{\beta_j}{1-\tau_j}$) both evaluate to $\frac{\alpha_j + \beta_j}{2}$ when the thresholds $\tau_j$ are set as $\tau_j = \frac{\alpha_j - \beta_j}{\alpha_j + \beta_j}$. In fact it can easily be seen that this value of $\tau_j$ minimizes the worst-case coefficient of $\mathbf{R}_{A^j(\mathbf{p})}^{\psi^{\mathrm{OvA}, n_j}}[\mathbf{u}_j]$ in the bound. Thus, we have

$$
\begin{aligned}
\mathbf{R}_{\mathbf{p}}^{\ell^H}[\Upsilon^{\mathrm{cas}}(\mathbf{u})] & \leq \sum_{j=1}^{h} \frac{1}{2}(\alpha_j + \beta_j) \cdot \mathbf{R}_{A^j(\mathbf{p})}^{\psi^{\mathrm{OvA}, n_j}}[\mathbf{u}_j] \\
& \leq \frac{1}{2} \max_{j \in [h]}(\alpha_j + \beta_j) \cdot \sum_{j=1}^{h} \mathbf{R}_{A^j(\mathbf{p})}^{\psi^{\mathrm{OvA}, n_j}}[\mathbf{u}_j] \\
& = \frac{1}{2} \max_{j \in [h]}(\alpha_j + \beta_j) \cdot \mathbf{R}_{\mathbf{p}}^{\psi^{\mathrm{cas}}}[\mathbf{u}] .
\end{aligned}
$$

The Theorem now follows from linearity of expectation. $\qquad\square$