

Supplementary material for “Multilabel reductions: what is my loss optimising?”

A Proof of results in body

Proof of Lemma 1. By definition, and linearity of expectation,

$$\begin{aligned}
 \text{Prec}@k(f) &= \mathbb{E}_x \mathbb{E}_{y|x} \left[\frac{1}{k} \cdot |\text{rel}(y) \cap \text{Top}_k(f(x))| \right] \\
 &= \mathbb{E}_x \mathbb{E}_{y|x} \left[\sum_{i \in [L]} \frac{1}{k} \cdot y_i \cdot \mathbb{I}[i \in \text{Top}_k(f(x))] \right] \\
 &= \mathbb{E}_x \left[\sum_{i \in \text{Top}_k(f(x))} \frac{1}{k} \cdot \mathbb{E}_{y|x} [y_i] \right] \\
 &= \mathbb{E}_x \left[\sum_{i \in \text{Top}_k(f(x))} \frac{1}{k} \cdot \mathbb{P}(y_i = 1 \mid x) \right].
 \end{aligned} \tag{12}$$

□

Proof of Lemma 3. By definition, and linearity of expectation,

$$\begin{aligned}
 \text{Rec}@k(f) &= \mathbb{E}_x \mathbb{E}_{y|x} \left[\frac{|\text{rel}(y) \cap \text{Top}_k(f(x))|}{|\text{rel}(y)|} \right] \\
 &= \mathbb{E}_x \mathbb{E}_{y|x} \left[\frac{\sum_{i \in [L]} y_i \cdot \mathbb{I}[i \in \text{Top}_k(f(x))]}{\sum_{j \in [L]} y_j} \right] \\
 &= \mathbb{E}_x \mathbb{E}_{y|x} \left[\sum_{i \in \text{Top}_k(f(x))} \frac{y_i}{\sum_{j \in [L]} y_j} \right] \\
 &= \mathbb{E}_x \left[\sum_{i \in \text{Top}_k(f(x))} \mathbb{E}_{y|x} \left[\frac{y_i}{\sum_{j \in [L]} y_j} \right] \right].
 \end{aligned} \tag{13}$$

Let us now observe that⁵

$$\begin{aligned}
 \mathbb{E}_{y|x} \left[\frac{y_i}{\sum_{j \in [L]} y_j} \right] &= \mathbb{E}_{y_i} \mathbb{E}_{y_{-i}|x, y_i} \left[\frac{y_i}{\sum_{j \in [L]} y_j} \right] \\
 &= \mathbb{P}(y_i = 1 \mid x) \cdot \mathbb{E}_{y_{-i}|x, y_i=1} \left[\frac{1}{1 + \sum_{j \neq i} y_j} \right] + \\
 &\quad (1 - \mathbb{P}(y_i = 1 \mid x)) \cdot \mathbb{E}_{y_{-i}|x, y_i=0} \left[\frac{0}{0 + \sum_{j \neq i} y_j} \right] \\
 &= \mathbb{P}(y_i = 1 \mid x) \cdot \mathbb{E}_{y_{-i}|x, y_i=1} \left[\frac{1}{1 + \sum_{j \neq i} y_j} \right].
 \end{aligned} \tag{14}$$

The result follows. □

Proof of Corollary 4. The fact that the Bayes-optimal scorers have the stated form follows identically to the proof of the optimal scorers for the precision@ k from Wydmuch et al. [2018], except that we now use the probabilities $\mathbb{P}(y'_i = 1 \mid x)$ in place of the marginals.

⁵We have also used here the assumption that $\mathbb{P}(\mathbf{0}_L \mid x) = 0$.

We now demonstrate that the ordering of the two probabilities may be different. In fact, the proof of this result is implicit in [Wydmuch et al. \[2018\]](#), who showed the inconsistency of the pick-one-label reduction for precision; we explicate it here. Consider the case of $L = 3$, and a distribution concentrated on a single $x \in \mathcal{X}$, and just two possible y :

$$\begin{aligned}\mathbb{P}(y = (1, 1, 0) \mid x) &= p \\ \mathbb{P}(y = (0, 0, 1) \mid x) &= 1 - p,\end{aligned}$$

for some $p \in (\frac{1}{2}, 1)$. Clearly, the marginal probabilities are $(p, p, 1 - p)$, and so the first two labels have higher marginal relevance.

Observe now that

$$\mathbb{P}(y'_i = 1 \mid x) = \mathbb{E}_{y|x} \left[\frac{y_i}{\sum_{j \in [L]} y_j} \right] = \begin{cases} p \cdot \frac{1}{2} & \text{if } i = 1 \\ p \cdot \frac{1}{2} & \text{if } i = 2 \\ (1 - p) \cdot \frac{1}{1} & \text{if } i = 3. \end{cases}$$

Now suppose $p \in (\frac{1}{2}, \frac{2}{3})$. Then, we have $1 - p > \frac{p}{2}$. Consequently, the ordering of the two probabilities will not be the same. In particular, the Bayes-optimal predictor for the recall@1 will favour the label y_3 , even though it has lowest marginal probability. \square

Proof of Proposition 5. We provide proofs for each of the reductions in turn.

OVA. For the one-versus-all reduction in Equation 6, observe that for a given $x \in \mathcal{X}$ with $f(x) \in \mathbb{R}^L$,

$$\begin{aligned}\mathbb{E}_{y|x} [\ell_{\text{OVA}}(y, f(x))] &= \mathbb{E}_{y|x} \left[\sum_{i \in [L]} \ell_{\text{BC}}(y_i, f_i(x)) \right] \\ &= \sum_{i \in [L]} \mathbb{E}_{y|x} [\ell_{\text{BC}}(y_i, f_i(x))] \\ &= \sum_{i \in [L]} \mathbb{E}_{y_i|x} [\ell_{\text{BC}}(y_i, f_i(x))].\end{aligned}$$

Consequently,

$$R_{\text{OVA}}(f) = \sum_{i \in [L]} \mathbb{E}_{(x, y_i)} [\ell_{\text{BC}}(y_i, f_i(x))].$$

PAL. For the pick-all-labels reduction, observe that for a given $x \in \mathcal{X}$ with $f(x) \in \mathbb{R}^L$, the loss in Equation 7 has expectation

$$\begin{aligned}\mathbb{E}_{y|x} [\ell_{\text{ML}}(y, f(x))] &= \sum_{i \in [L]} \mathbb{E}_{y|x} [y_i \cdot \ell_{\text{MC}}(i, f(x))] \\ &= \sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \cdot \ell_{\text{MC}}(i, f(x)) \\ &= N(x) \cdot \mathbb{E}_{z|x} \ell_{\text{MC}}(z, f(x)),\end{aligned}$$

where $N(x) \doteq \sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x)$ is a normaliser, being the expected number of labels per instance. Here, z is a discrete random variable taking values in $[L]$, with $\mathbb{P}(z = i \mid x) \doteq \frac{\mathbb{P}(y_i = 1 \mid x)}{N(x)}$. Consequently,

$$R_{\text{PAL}}(f) = \mathbb{E}_{(x, z)} [N(x) \cdot \ell_{\text{MC}}(z, f(x))].$$

OVA-N. For the normalised one-versus-all reduction in Equation 8, observe that

$$\begin{aligned}\mathbb{E}_{y|x} [\ell_{\text{ML}}(y, f(x))] &= \sum_{i \in [L]} \mathbb{E}_{y|x} \left[\frac{y_i}{\sum_{j \in [L]} y_j} \cdot \ell_{\text{BC}}(1, f_i(x)) + \left(1 - \frac{y_i}{\sum_{j \in [L]} y_j} \right) \cdot \ell_{\text{BC}}(0, f_i(x)) \right] \\ &= \sum_{i \in [L]} \mathbb{P}(y'_i = 1 \mid x) \cdot \ell_{\text{BC}}(1, f_i(x)) + (1 - \mathbb{P}(y'_i = 1 \mid x)) \cdot \ell_{\text{BC}}(0, f_i(x))\end{aligned}$$

$$= \sum_{i \in [L]} \mathbb{E}_{y'_i | x} [\ell_{\text{BC}}(y'_i, f_i(x))],$$

where in the third line we used the definition of $\mathbb{P}(y'_i = 1 \mid x)$ (Equation 5) and Equation 14. Consequently,

$$R_{\text{OVA-N}}(f) = \sum_{i \in [L]} \mathbb{E}_{(x, y'_i)} [\ell_{\text{BC}}(y'_i, f_i(x))].$$

PAL-N. For the normalised pick-all-labels reduction observe that for a given $x \in \mathcal{X}$ with $f(x) \in \mathbb{R}^L$, the loss in Equation 9 has expectation

$$\begin{aligned} \mathbb{E}_{y|x} [\ell_{\text{ML}}(y, f(x))] &= \sum_{i \in [L]} \mathbb{E}_{y|x} \left[\frac{y_i}{\sum_{j \in [L]} y_j} \cdot \ell_{\text{MC}}(i, f(x)) \right] \\ &= \sum_{i \in [L]} \mathbb{P}(y'_i = 1 \mid x) \cdot \ell_{\text{MC}}(i, f(x)) \\ &= \sum_{i \in [L]} \mathbb{P}(z' = i \mid x) \cdot \ell_{\text{MC}}(i, f(x)) \\ &= \mathbb{E}_{z'} \ell_{\text{MC}}(z', f(x)), \end{aligned}$$

where in the second line we used the fact that $\mathbb{P}(z' = i \mid x) \doteq \mathbb{P}(y'_i = 1 \mid x)$ is a valid multiclass distribution, as previously noted in Wydmuch et al. [2018] and evident from Equation 14. Consequently,

$$R_{\text{PAL-N}}(f) = \mathbb{E}_{(x, z')} [\ell_{\text{MC}}(z', f(x))].$$

POL. The result here trivially follows from the fact that POL is a stochastic version of PAL. \square

Proof of Lemma 6. By Proposition 5, we have

$$\begin{aligned} R_{\text{PAL}}(f) &= \mathbb{E}_{(x, z)} [N(x) \cdot \ell_{\text{MC}}(z, f(x))] \\ &= \mathbb{E}_x \left[\sum_{i \in [L]} N(x) \cdot \mathbb{P}(z = i \mid x) \cdot \ell_{\text{MC}}(i, f(x)) \right] \\ &= \mathbb{E}_x \left[\sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \cdot \left\{ \ell_{\text{BC}}(1, f_i(x)) + \sum_{j \neq i} \ell_{\text{BC}}(0, f_j(x)) \right\} \right] \\ &= \mathbb{E}_x \left[\sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \cdot \left\{ \ell_{\text{BC}}(1, f_i(x)) - \ell_{\text{BC}}(0, f_i(x)) + \sum_{j \in [L]} \ell_{\text{BC}}(0, f_j(x)) \right\} \right] \\ &= \mathbb{E}_x \left[\sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \cdot \{ \ell_{\text{BC}}(1, f_i(x)) - \ell_{\text{BC}}(0, f_i(x)) \} \right] + \\ &\quad \mathbb{E}_x \left[\left(\sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \right) \cdot \left(\sum_{j \in [L]} \ell_{\text{BC}}(0, f_j(x)) \right) \right] \\ &= \mathbb{E}_x \left[\sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \cdot \{ \ell_{\text{BC}}(1, f_i(x)) - \ell_{\text{BC}}(0, f_i(x)) \} \right] + \\ &\quad \mathbb{E}_x \left[N(x) \cdot \left(\sum_{j \in [L]} \ell_{\text{BC}}(0, f_j(x)) \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_x \left[\sum_{i \in [L]} \mathbb{P}(y_i = 1 \mid x) \cdot \ell_{\text{BC}}(1, f_i(x)) + \mathbb{P}(y_i = 0 \mid x) \cdot \ell_{\text{BC}}(0, f_i(x)) - \ell_{\text{BC}}(0, f_i(x)) \right] + \\
&\quad \mathbb{E}_x \left[N(x) \cdot \left(\sum_{j \in [L]} \ell_{\text{BC}}(0, f_j(x)) \right) \right] \\
&= \sum_{i \in [L]} \mathbb{E}_{x, y_i} [\ell_{\text{BC}}(y_i, f_i(x))] + \mathbb{E}_x \left[(N(x) - 1) \cdot \left(\sum_{j \in [L]} \ell_{\text{BC}}(0, f_j(x)) \right) \right].
\end{aligned}$$

□

Proof of Corollary 7. We provide proofs for each of the reductions in turn.

OVA. For the one-versus-all reduction, we may compute each $f_i^*(x)$ independently. When ℓ_{BC} is a strictly proper loss, these will by definition be equal $\mathbb{P}(y_i = 1 \mid x)$.

PAL. For the pick-all-labels reduction, the factor $N(x)$ only modifies the marginal distribution of x . Since $N(x) > 0$, this weighting factor will not affect the Bayes-optimal solution. Thus, since ℓ_{MC} is strictly proper, the Bayes-optimal prediction will by definition be $\mathbb{P}(z = i \mid x) = \frac{\mathbb{P}(y_i = 1 \mid x)}{N(x)}$.

OVA-N. For the normalised one-versus-all reduction, since ℓ_{BC} is a strictly proper loss, by definition $f_i^*(x) = \mathbb{P}(y_i' = 1 \mid x)$.

PAL-N. For the normalised pick-all-labels reduction, since ℓ_{MC} is a strictly proper loss, and we have a multiclass risk, by definition $f_i^*(x) = \mathbb{P}(z' = i \mid x) = \mathbb{P}(y_i' = 1 \mid x)$.

POL. The result here trivially follows from the fact that POL is a stochastic version of PAL. □

Proof of Corollary 8. Observe first that, by (12),

$$\begin{aligned}
k \cdot \text{Prec@}k(f) &= \sum_{i \in [L]} \mathbb{E}_{(x, y_i)} [y_i \cdot \mathbb{I}[i \in \text{Top}_k(f(x))]] \\
&= \sum_{i \in [L]} \mathbb{E}_{(x, y_i)} [y_i \cdot (1 - \ell_{\text{top-}k}(i, f(x)))] \\
&= \sum_{i \in [L]} \mathbb{E}_{y_i} [y_i] - \mathbb{E}_{(x, y_i)} [y_i \cdot \ell_{\text{top-}k}(i, f(x))] \\
&= \text{constant} - \sum_{i \in [L]} \mathbb{E}_{(x, y_i)} [y_i \cdot \ell_{\text{top-}k}(i, f(x))] \\
&= \text{constant} - \mathbb{E}_{(x, y)} \left[\sum_{i \in [L]} y_i \cdot \ell_{\text{top-}k}(i, f(x)) \right] \\
&= \text{constant} - R_{\text{PAL}}(f).
\end{aligned}$$

Similarly, by (13),

$$\begin{aligned}
\text{Rec@}k(f) &= \mathbb{E}_{(x, y)} \left[\sum_{i \in [L]} \frac{y_i}{\sum_{j \in [L]} y_j} \cdot \mathbb{I}[i \in \text{Top}_k(f(x))]] \right] \\
&= \mathbb{E}_{(x, y)} \left[\sum_{i \in [L]} \frac{y_i}{\sum_{j \in [L]} y_j} \cdot (1 - \ell_{\text{top-}k}(i, f(x))) \right] \\
&= \mathbb{E}_{(x, y)} \left[\sum_{i \in [L]} \frac{y_i}{\sum_{j \in [L]} y_j} \right] - \mathbb{E}_{(x, y)} \left[\sum_{i \in [L]} \frac{y_i}{\sum_{j \in [L]} y_j} \cdot \ell_{\text{top-}k}(i, f(x)) \right]
\end{aligned}$$

$$\begin{aligned}
&= \text{constant} - \mathbb{E}_{(x,y)} \left[\sum_{i \in [L]} \frac{y_i}{\sum_{j \in [L]} y_j} \cdot \ell_{\text{top-}k}(i, f(x)) \right] \\
&= \text{constant} - R_{\text{PAL-N}}(f).
\end{aligned}$$

□

Proof of Proposition 9. For clarity, we will explicate here the dependence of all quantities on the type of risk (multiclass or multilabel reduction), and the base loss. For simplicity, we will assume the existence of a *surrogate regret bound* for the top- k risk: that is, we assume that for every scorer, the surrogate loss ℓ_{MC} satisfies

$$\text{reg}_{\text{MC}}(f; \ell_{\text{top-}k}) \leq \Psi(\text{reg}_{\text{MC}}(f; \ell_{\text{MC}}))$$

for some function Ψ such that $\lim_{z \rightarrow 0^+} \Psi(z) = 0$, where $\text{reg}_{\text{MC}}(f; \ell_{\text{MC}})$ denotes multiclass regret for a scorer using a loss ℓ_{MC} . This clearly implies consistency.

We provide bounds for the precision and recall in turn.

Precision bound. By Corollary 8,

$$\text{Prec@}k(f) = \text{Constant} - \frac{1}{k} \cdot R_{\text{PAL}}(f; \ell_{\text{top-}k}).$$

Thus,

$$\text{reg}(f; \text{P@}k) = \frac{1}{k} \cdot \text{reg}_{\text{PAL}}(f; \ell_{\text{top-}k})$$

where $\text{reg}_{\text{PAL}}(f; \ell_{\text{MC}}) \doteq R_{\text{PAL}}(f; \ell_{\text{MC}}) - \inf_{g: \mathcal{X} \rightarrow \mathbb{R}} R_{\text{PAL}}(g; \ell_{\text{MC}})$ is the regret under the PAL reduction, using a base multiclass loss of ℓ_{MC} .

By Proposition 5, for any multiclass loss ℓ_{MC} ,

$$R_{\text{PAL}}^{\mathbb{P}}(f; \ell_{\text{MC}}) = C \cdot \mathbb{E}_{(\bar{x}, z)} [\ell_{\text{MC}}(z, f(\bar{x}))] = C \cdot R_{\text{MC}}^{\bar{\mathbb{P}}}(f; \ell_{\text{MC}}),$$

where, if x has distribution with density $p(x)$ with respect to base measure μ , \bar{x} has distribution with density $\bar{p}(x) \doteq C^{-1} \cdot p(x) \cdot N(x)$, and $C \doteq \int_{\mathcal{X}} p(x') \cdot N(x') \, d\mu(x)$. Note that $C \leq L < +\infty$ since $N(x) \in [0, L]$ for every $x \in \mathcal{X}$. Note also that the above explicates the dependence of the risks on the underlying distributions, as they are distinct in the LHS and RHS.

The above implies that the multilabel PAL risk given a distribution $\mathbb{P}(x, z)$ is expressible as a constant times a standard multiclass risk over a distribution $\bar{\mathbb{P}}(\bar{x}, z)$. Thus, a Bayes-optimal scorer for the PAL risk must satisfy

$$f^* \in \underset{f: \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} R_{\text{PAL}}^{\mathbb{P}}(f; \ell_{\text{MC}}) = \underset{f: \mathcal{X} \rightarrow \mathbb{R}}{\text{argmin}} R_{\text{MC}}^{\bar{\mathbb{P}}}(f; \ell_{\text{MC}}).$$

Consequently, if f^* denotes the Bayes-optimal scorer for the PAL risk, we have

$$\begin{aligned}
\text{reg}_{\text{PAL}}^{\mathbb{P}}(f; \ell_{\text{MC}}) &= C \cdot \mathbb{E}_{(\bar{x}, z)} [\ell_{\text{MC}}(z, f(\bar{x})) - \ell_{\text{MC}}(z, f^*(\bar{x}))] \\
&= C \cdot \text{reg}_{\text{MC}}^{\bar{\mathbb{P}}}(f; \ell_{\text{MC}}).
\end{aligned}$$

We have thus reduced the *multilabel* regret on the left-hand side to a *multiclass* regret on the right-hand side. Further, in doing so, we have introduced a distribution $\bar{\mathbb{P}}(\bar{x}, z)$ whose marginal density \bar{p} is *distorted* compared to the original p .

Now, since ℓ_{MC} is assumed to be a consistent surrogate to the top- k loss,

$$\begin{aligned}
\text{reg}(f; \text{P@}k) &= \frac{1}{k} \cdot \text{reg}_{\text{PAL}}^{\mathbb{P}}(f; \ell_{\text{top-}k}) \\
&= \frac{C}{k} \cdot \text{reg}_{\text{MC}}^{\bar{\mathbb{P}}}(f; \ell_{\text{top-}k}) \\
&\leq \frac{C}{k} \cdot \Psi(\text{reg}_{\text{MC}}^{\bar{\mathbb{P}}}(f; \ell_{\text{MC}}))
\end{aligned}$$

$$= \frac{C}{k} \cdot \Psi \left(\text{reg}_{\text{PAL}}^{\mathbb{P}}(f; \ell_{\text{MC}}) \right),$$

where we have used the fact that a surrogate regret bound holds for *any* distribution, even one where the marginals are distorted. We thus have a surrogate regret bound for the precision@ k .

Recall bound. By Corollary 8,

$$\text{Rec}@k(f) = \text{Constant} - R_{\text{PAL-N}}(f; \ell_{\text{top-}k}).$$

Thus,

$$\text{reg}(f; \text{R}@k) = \text{reg}_{\text{PAL-N}}(f; \ell_{\text{top-}k}).$$

By Proposition 5,

$$R_{\text{PAL-N}}(f; \ell_{\text{top-}k}) = \mathbb{E}_{(x, z')} [\ell_{\text{top-}k}(z', f(x))]$$

where $\mathbb{P}'(x, z')$ is as per (11). That is, the multilabel PAN risk is a multiclass risk. Now, since ℓ_{MC} is assumed to be a consistent surrogate to the top- k loss,

$$\begin{aligned} \text{reg}(f; \text{R}@k) &= \text{reg}_{\text{PAL-N}}^{\mathbb{P}}(f; \ell_{\text{top-}k}) \\ &= \text{reg}_{\text{MC}}^{\mathbb{P}'}(f; \ell_{\text{top-}k}) \\ &\leq \Psi \left(\text{reg}^{\mathbb{P}'}(f; \ell_{\text{MC}}) \right) \\ &= \Psi \left(\text{reg}_{\text{PAL-N}}^{\mathbb{P}}(f; \ell_{\text{MC}}) \right). \end{aligned}$$

We thus have a surrogate regret bound for the recall@ k . \square

Proof of Proposition 10. Observe that by Lemma 3 and Corollary 7,

$$\begin{aligned} \text{reg}(f; \text{R}@k) &= \mathbb{E}_x \left[\sum_{i \in \text{Top}_k(f(x))} \mathbb{P}(y'_i = 1 \mid x) - \sum_{i \in \text{Top}_k(f^*(x))} \mathbb{P}(y'_i = 1 \mid x) \right] \\ &= \mathbb{E}_x \left[\sum_{i \in \text{Top}_k(f(x))} \mathbb{P}(y'_i = 1 \mid x) - \sum_{i \in \text{Top}_k(\mathbb{P}(y' \mid x))} \mathbb{P}(y'_i = 1 \mid x) \right]. \end{aligned}$$

The first inequality thus follows exactly as per Wydmuch et al. [2018, Theorem 2]: we simply swap $\mathbb{P}(y_i = 1 \mid x)$ with $\mathbb{P}(y'_i = 1 \mid x)$. Thus,

$$\text{reg}(f; \text{R}@k) \leq 2 \cdot \mathbb{E}_x \left[\max_{i \in [L]} |f_i(x) - \mathbb{P}(y'_i = 1 \mid x)| \right]. \quad (15)$$

The second inequality is a standard consequence of existing bounds for strongly proper losses [Agarwal, 2014]. By Proposition 5,

$$\begin{aligned} R_{\text{OVA-N}}(f) &= \sum_{i \in [L]} \mathbb{E}_{(x, y_i)} [\ell_{\text{BC}}(y_i, f_i(x))] \\ &= \sum_{i \in [L]} R_{\text{BC}}(f_i), \end{aligned}$$

where R_{BC} denotes a binary classification risk. Since the risk decomposes across each $i \in [L]$, the Bayes-optimal scorers may be computed separately for i . The regret thus similarly decomposes as:

$$\begin{aligned} \text{reg}(f; \ell_{\text{OVA-N}}) &= \sum_{i \in [L]} \text{reg}(f_i; \ell_{\text{BC}}) \\ &\geq \sum_{i \in [L]} \frac{\lambda}{2} \cdot \mathbb{E}_x [(f_i(x) - \mathbb{P}(y'_i = 1 \mid x))^2] \end{aligned}$$

$$\begin{aligned}
&= \frac{\lambda}{2} \cdot \mathbb{E}_x \left[\sum_{i \in [L]} (f_i(x) - \mathbb{P}(y'_i = 1 \mid x))^2 \right] \\
&\geq \frac{\lambda}{2} \cdot \mathbb{E}_x \left[\max_{i \in [L]} (f_i(x) - \mathbb{P}(y'_i = 1 \mid x))^2 \right] \text{ by non-negativity of each term} \\
&= \frac{\lambda}{2} \cdot \mathbb{E}_x \left[\left(\max_{i \in [L]} |f_i(x) - \mathbb{P}(y'_i = 1 \mid x)| \right)^2 \right] \\
&\geq \frac{\lambda}{2} \cdot \left(\mathbb{E}_x \left[\max_{i \in [L]} |f_i(x) - \mathbb{P}(y'_i = 1 \mid x)| \right] \right)^2 \text{ by Jensen's inequality} \\
&\geq \frac{\lambda}{2} \cdot \left(\max_{i \in [L]} \mathbb{E}_x |f_i(x) - \mathbb{P}(y'_i = 1 \mid x)| \right)^2 \text{ by Jensen's inequality,}
\end{aligned}$$

where in the second line, we used the fact that ℓ_{BC} is strongly proper, and Agarwal [2014, Theorem 13]; and in the fifth line, we used the fact that the maximum is over a sequence of L nonnegative quantities, so that the maximum of their squares is the square of their maximum. We thus have

$$\max_{i \in [L]} \mathbb{E}_x |f_i(x) - \mathbb{P}(y'_i = 1 \mid x)| \leq \sqrt{\frac{2}{\lambda}} \cdot \sqrt{\text{reg}(f; \ell_{\text{OVA-N}})}.$$

Applying this to (15), the claim follows. \square

Proof of Lemma 11. Observe that

$$\begin{aligned}
R_{\text{top-}k}(f) &= \mathbb{P}(z \notin \text{Top}_k(f(x))) \\
&= \mathbb{E}_x \mathbb{E}_{z|x} [\mathbb{1}_{z \notin \text{Top}_k(f(x))}] \\
&= \mathbb{E}_x \mathbb{E}_{z|x} [\ell_{\text{top-}k}(z, f(x))] \\
&= \mathbb{E}_x \left[\sum_{i \in [L]} \mathbb{P}(z = i \mid x) \cdot \ell_{\text{top-}k}(i, f(x)) \right] \\
&= \mathbb{E}_x \left[\sum_{i \notin \text{Top}_k(f(x))} N(x)^{-1} \cdot \mathbb{P}(y_i = 1 \mid x) \right],
\end{aligned}$$

where in the last line we used the definition of $\mathbb{P}(z = i \mid x)$ from Equation 11. \square