

Multilabel reductions: what is my loss optimising?

Aditya Krishna Menon, Ankit Singh Rawat, Sashank J. Reddi, Sanjiv Kumar Google Research NYC

Five common multilabel reductions implicitly optimise for precision or recall

Multilabel classification via reductions Multilabel metrics: precision and recall Multilabel classification: predict a binary label vector **Precision** and **recall@k** are standard metrics for retrieval settings. They measure the # of positives in the top-k scoring indices, suitably normalised: # of labels **y** = Given (x, y) for $y \in \{0, 1\}^L$, find $f(x) \in \mathbb{R}^{L}$ **f** = Ideally, want $f_i(x)$ high if $y_i = 1$. The challenge is that L Prec@ may be large; so how do we efficiently find such an *f*? Rec@*k*(*f*) = $\mathbb{E}_{(x, y)}[| \operatorname{top}_{k}(f(x)) \cap \operatorname{pos}(y) | / | \operatorname{pos}(y) |]$ Common algorithms reduce the problem to binary or multiclass learning; e.g., we may create **Fact:** Given a distribution P(x, y), the optimal f^* for these metrics preserve the ordering of: • multiple binary examples, one for each label (Prec@k) P($\mathbf{y}_i = 1 | x$) Marginal relevance Non weig (Rec@k) P($\mathbf{y}'_i = 1 | x$) = P($\mathbf{y}_i = 1 | x$) $\cdot \mathbb{E}[(1 + \Sigma_{j \neq i} \mathbf{y}_j)^{-1}]$ Non-constant • multiple multi-class examples, one for each **positive** label • one multi-class example for a **random positive** label **y** = 0 0 Note that owing to the weighting above, these Multi-class examples, optimal solutions are **incompatible** in general! (*x*, 🔽) fed into e.g. softmax cross-entropy Implicit losses for multilabel reductions Key question of this work To compare different reductions, we explicate their Such reductions can be practically effective, but: implicit multilabel losses, i.e., $\ell(y, f(x))$: what multilabel metric do they optimise? One-versus-all (OVA) $\sum \ell_{\rm BC}(y_i, f_i)$ Answering this helps inform the choice of reduction, $\sum y_i \cdot \ell_{\mathrm{MC}}(i, f)$ Pick-all-labels (PAL) depending on our end goal. $i \in [L]$ Key contribution of this work

five common reductions *implicitly* optimise for either precision <u>or</u> recall@k!

To begin, we study a basic property of these metrics.

weighting

Binary loss Multiclass loss; has label "competition" $\sum_{i\in[L]} \left\{ \frac{y_i}{\sum_{j\in[L]} y_j} \cdot \ell_{\mathrm{BC}}(1,f_i) + \left(1 - \frac{y_i}{\sum_{j\in[L]} y_j}\right) \cdot \ell_{\mathrm{BC}}(0,f_i) \right\}$ OVA normalised $\sum_{i \in [L]} \frac{y_i}{\sum_{j \in [L]} y_j} \cdot \ell_{\mathrm{MC}}(i, f)$ PAL normalised Pick-one-label (POL) For log-loss, KL(label, prediction)

Optimal scores for multilabel reductions

Equipped with these losses, observe that, e.g.,

generally:

each other. Further remarks:

Illustration of differences in reductions

Synthetic problem where normalised reduction gives recall gains, at the expense of precision.



 $\operatorname{Prec}@k(\boldsymbol{f}) = \mathbb{E}_{(x, y)}[\Sigma_{i \in [L]} \ell_{\operatorname{top-k}}(\boldsymbol{y}, \boldsymbol{f}_{i}(x)) / k],$

where ℓ_{top-k} is a "top-k" multiclass loss; i.e., precision is equivalent to PAL for a specific loss! More

Fact: Given a distribution P(x, y), the optimal f_i^* is: (OVA) $P(\mathbf{y}_i = 1 | x)$ Classical (PAL) $P(\mathbf{y}_i = 1 | x) / N(x)$ Expected # of +'ves Doesn't vary with *i* (Rest) $P(\mathbf{y}'_i = 1 | \mathbf{x}) \longrightarrow c.f. recall@k optimal scorer$

Different reductions target either precision or recall! Subtle differences in the loss thus have non-trivial effects; e.g., $\ell_{PAI - Norm}(y, f(x)) = (\Sigma_i y_i)^{-1} \ell_{PAI}(y, f(x))$, but the optimal solutions for the two are <u>not</u> scalings of

 PAL scores are not calibrated across instances! • PAL may \geq OVA since it directly bounds top-k loss PAL with a multiclass OVA loss ≠ multilabel OVA; PAL implicitly places a greater weight on each "negative" label