

Towards Calibrated Multi-label Deep Neural Networks

Jiacheng Cheng Nuno Vasconcelos Department of Electrical and Computer Engineering University of California, San Diego {jicheng, nvasconcelos}@ucsd.edu

CVPR 2024

Introduction



Calibration:

A classifier is calibrated if it predicts a posterior class probability of p when the selection of the class is correct $p \times 100\%$ of the time. The importance of calibration has been noted for many applications. For example, in medical diagnosis, probabilities can be used to determine which examples require human inspection, thus avoiding the cost of manually inspecting all images. However, the process can only be trusted if the DNN provides accurate posterior estimates.

A high-accuracy model may not be well-calibrated

Sample ID	Predicted Probability (\hat{p})	Predicted Label	True Label	Correct Prediction?
1	0.9	+1	+1	Yes
2	0.8	+1	+1	Yes
3	0.7	+1	+1	Yes
4	0.6	+1	-1	No
5	0.4	-1	-1	Yes
6	0.3	-1	-1	Yes
7	0.2	-1	-1	Yes
8	0.1	-1	+1	No
9	0.9	+1	+1	Yes
10	0.85	+1	+1	Yes

Model Accuracy:

- Total samples = 10, correctly predicted samples = 8;
- Accuracy = $\frac{8}{10} = 80\%$.

Calibration Issues:

- In the [0.8, 0.9] probability range (Samples 1, 2, 9, 10), the average predicted probability is $\hat{p} = 0.86$, but the actual positive class frequency is 100%.
- In the [0.6, 0.7] probability range (Sample 4), the predicted probability is $\hat{p} = 0.6$, but the actual positive class frequency is 0%.

This indicates that the model, despite being accurate, is poorly calibrated because its predicted probabilities do not align with actual frequencies.

Introduction



Multi-label DNNs can be trained with class probability estimation (CPE) losses that encourage probability calibration, such as the binary cross-entropy (BCE) loss. However, the multi-label setting is highly imbalanced, due to the sparseness of positives, as most tags are absent from any given image. In result, asymmetric losses such as such as the focal loss of or the asymmetric (ASY) loss of tend to produce much higher labeling accuracy than the BCE.

Our preliminary studies reveal that multi-label DNNs trained with existing losses tend to produce poorly-calibrated probabilities. This is illustrated by the calibration curves in Figure 1, where the calibration of the focal and ASY losses are drastically far from perfect. We argue that these popular multi-label losses are poorly suited for class-probability estimation because they are not **strictly proper**. This is a property that denotes the family of **losses uniquely minimized by the true posterior probability**.



Preliminary



Given observation $\mathbf{x} \in \mathcal{X}$, the goal is to estimate the vector $\boldsymbol{\eta}(\mathbf{x}) = \left[\eta^{(1)}(\mathbf{x}), \cdots, \eta^{(T)}(\mathbf{x})\right]^{\top}$ of class-posterior probabilities

$$\eta^{(t)}(\mathbf{x}) = P\left(y^{(t)} = 1 | \mathbf{x}\right), \forall t \in \{1, \cdots, T\}.$$
 (1)

A multi-label DNN typically performs this probability estimation in two steps. First, it maps $\mathbf{x} \in \mathcal{X}$ into a real-valued score vector $\mathbf{v}(\mathbf{x}) = \left[v^{(1)}(\mathbf{x}), \cdots, v^{(T)}(\mathbf{x})\right]^{\top} \in \mathbb{R}^{T}$. The embedding $\mathbf{v} : \mathcal{X} \to \mathbb{R}^{T}$ is composed by a sequence of linear and nonlinear operations. Each $v^{(t)}(\mathbf{x})$ is then mapped into a class-posterior probability estimate with

$$\widehat{\eta}^{(t)}(x) = \widehat{P}\left(y^{(t)} = 1 | \mathbf{x}\right) = [\Psi]^{-1}(v^{(t)}(\mathbf{x})), \quad (2)$$

where $[\Psi]^{-1}(\cdot)$ is an inverse link function. This can be any strictly increasing function $[\Psi]^{-1} : \mathbb{R} \to \Delta$, but is usually the logistic inverse link, or sigmoid activation function

$$\sigma(v) = \frac{1}{1 + e^{(-v)}},$$
(3)

Preliminary



Given a CPE loss $\ell : \Delta \times \{-1, +1\} \to \mathbb{R}$ that assigns a cost $\ell(\widehat{\eta}, \pm 1)$ for predicting $\widehat{\eta}$ as the class-posterior probability of positive class when the true label is $y = \pm 1$,¹ the optimal posterior probability estimator minimizes the *risk*:

$$\mathcal{R}(\widehat{\boldsymbol{\eta}}) = \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[\sum_{t=1}^{T} \ell\left(\widehat{\eta}^{(t)}(\mathbf{x}), y^{(t)}\right) \right], \qquad (4) \quad \text{期望风险}$$

where \mathbb{E} denotes expectation. To train a multi-label probabilistic DNN, this is approximated by the *empirical risk*

$$\widehat{\mathcal{R}}(\widehat{\boldsymbol{\eta}};\mathcal{D}) = \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \ell(\widehat{\boldsymbol{\eta}}^{(t)}(\mathbf{x}_i), y_i^{(t)})$$
(5) 经验风险

on a training dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ of i.i.d. samples from $\mathcal{X} \times \mathcal{Y}$ [74].

Preliminary



BCE Loss

$$\begin{cases} \ell_{+1}^{BCE}(\widehat{\eta}) = -\log(\widehat{\eta}), \\ \ell_{-1}^{BCE}(\widehat{\eta}) = -\log(1-\widehat{\eta}). \end{cases}$$
(6)

ASY Focal Loss

$$\begin{cases} \ell_{+1}^{AsyFocal}(\widehat{\eta}) = -(1-\widehat{\eta})^{\gamma^{+}}\log(\widehat{\eta}), \\ \ell_{-1}^{AsyFocal}(\widehat{\eta}) = -(\widehat{\eta})^{\gamma^{-}}\log(1-(\widehat{\eta})), \end{cases}$$
(7)

ASY Focal Loss

$$\begin{cases} \ell_{+1}^{ASY}(\widehat{\eta}) = -(1-\widehat{\eta})^{\gamma^{+}} \log(\widehat{\eta}), \\ \ell_{-1}^{ASY}(\widehat{\eta}) = -(\widehat{\eta}-m)_{+}^{\gamma^{-}} \log(1-(\widehat{\eta}-m)_{+}), \end{cases}$$
(8)



$$\mathcal{R}(\widehat{\boldsymbol{\eta}}) = \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[\sum_{t=1}^{T} \ell\left(\widehat{\eta}^{(t)}(\mathbf{x}), y^{(t)}\right) \right], \quad (4)$$

Since (4) can be rewritten as

$$\mathcal{R}(\widehat{\boldsymbol{\eta}}) = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathbf{y}|\mathbf{x}} \left[\sum_{t=1}^{T} \ell\left(\widehat{\eta}^{(t)}(\mathbf{x}), y^{(t)}\right) \middle| \mathbf{x} \right] \right]$$
(9)
$$= \mathbb{E}_{\mathbf{x}} \left[\sum_{t=1}^{T} C\left(\eta^{(t)}(\mathbf{x}), \widehat{\eta}^{(t)}(\mathbf{x})\right) \right]$$
(10)

where C is the (pointwise) conditional risk

$$C\left(\eta(\mathbf{x}), \widehat{\eta}(\mathbf{x})\right) = \eta(\mathbf{x})\ell_{+1}(\widehat{\eta}(\mathbf{x})) + (1 - \eta(\mathbf{x}))\ell_{-1}(\widehat{\eta}(\mathbf{x})).$$
(11)

Definition 1. (Strict Properness [4, 19, 68]) The pair of partial losses $\{\ell_{-1}, \ell_{+1}\}$ or $\{\ell_{-1,\Psi}, \ell_{+1,\Psi}\}$ is strictly proper if the conditional risk $C(\eta, \hat{\eta})$ of (11) is uniquely minimized by $\hat{\eta} = \eta$ for all $\eta \in [0, 1]$.

$$C = y(x)[-log(\hat{\eta}(x))] + (l - y(x))(-log(1 - \hat{\eta}(x)))$$

$$C' = -\frac{y(x)}{\hat{\eta}(x)} + \frac{l - y(x)}{l - \hat{\eta}(x)}$$

$$C \gg \bigoplus \inf X$$

$$f(x) = y(x) \text{ if }, \quad C \equiv N$$



Theorem 2. Denote the conditional risk of (11) defined by the ASY loss ℓ^{ASY} of (8) by $C^{ASY}(\eta, \hat{\eta})$. Let

$$\widehat{\eta}^{ASY,*}(\mathbf{x}) = \operatorname*{arg\,min}_{\widehat{\eta} \in [0,1]} C^{ASY}(\eta(\mathbf{x}), \widehat{\eta}(\mathbf{x}))$$
(12)

be the minimizer of this risk for any $\mathbf{x} \in \mathcal{X}$. Then there is a mapping ϕ such that

$$\eta(\mathbf{x}) = \phi\left(\widehat{\eta}^{ASY,*}(\mathbf{x});\gamma^{-},\gamma^{+},m\right),\qquad(13)$$

where

$$\phi(z;\gamma^{-},\gamma^{+},m) = \frac{h((z-m)_{+};\gamma^{-})}{h((z-m)_{+};\gamma^{-}) + h(1-z;\gamma^{+})},$$
(14)

$$h(z;\gamma) = \frac{z^{\gamma} - \gamma z^{\gamma-1}(1-z)\log(1-z)}{1-z}.$$
 (15)

 $\phi: \Delta \to \Delta$ is a bijective map if and only if m = 0.



rather than designing a CPE loss $\ell_{\pm 1}(\hat{\eta})$ explicitly, we consider the separate design of a composite loss $\ell_{\pm 1,\Psi}(v)$ and an inverse link $[\Psi]^{-1}(v)$ as below

$$\begin{cases} \ell_{+1,\Psi}(v) = -\frac{1}{\zeta^+} \log\left(\frac{1}{1+e^{(-k^+(v-b^+))}}\right), \\ \\ \ell_{-1,\Psi}(v) = -\frac{1}{\zeta^-} \log\left(\frac{1}{1+e^{(+k^-(v-b^-))}}\right), \end{cases}$$
(16)

$$[\Psi]^{-1}(v) = \frac{\frac{k^{-}}{\zeta^{-}}\sigma(-k^{-}(v-b^{-}))}{\frac{k^{-}}{\zeta^{-}}\sigma(-k^{-}(v-b^{-})) + \frac{k^{+}}{\zeta^{+}}\sigma(k^{+}(v-b^{+}))}, \quad (17)$$



Figure 2. Example of the SPA losses (left) and the associated inverse link (right) where $(\zeta^+, k^+, b^+, \zeta^-, k^-, b^-) = (1, 1, 0, 5, 3, 1)$.

where $\sigma(\cdot)$ is the sigmoid function of (3). According to the Theorem 3 below, the CPE loss composed of (16) and (17) is strictly proper and thus referred to as the *Strictly Proper* Asymmetric (SPA) loss in this work. In practice, following the practice in the literature [61], we reduce the positive partial loss to the BCE loss by setting $(\zeta^+, k^+, b^+) = (1, 1, 0)$ and only tune the hyperparaeters ζ^{-}, k^{-}, b^{-} of the negative partial loss. The motivation for these hyperparmeters is simple and intuitive: i) k^-, b^- define an affine transformation of the logits $v^{(t)}(\mathbf{x})$ that enables control of the rate at which ℓ_{-1} decays to 0 as $v \to 0$. ii) ζ^- is a scale factor that controls the overall weight of negative examples. Note that unlike prior losses (6)-(8), SPA does not directly operate on the probability estimate $\hat{\eta}(\mathbf{x})$, and that the introduction of these hyperparameters induces the need for the inverse link of (17) to achieve the strict properness, rather than simply using the sigmoid of (3).

Theorem 3. For any ζ^+ , ζ^- , k^+ , $k^- \in \mathbb{R}_{++}$ and b^- , $b^+ \in \mathbb{R}$, the CPE loss composed of composite loss (16) and inverse link function (17) is strictly proper.





Figure 2. Example of the SPA losses (left) and the associated inverse link (right) where $(\zeta^+, k^+, b^+, \zeta^-, k^-, b^-) = (1, 1, 0, 5, 3, 1)$.



Label Pair Regularizer

A strictly proper CPE loss only guarantees perfect label-probability estimates for asymptotically large datasets. In practice, for finite datasets, empirical risk minimization does not guarantee the recovery of true risk minimizer. In this case, probability calibration can usually be improved by adding regularization terms to the loss function. In this work, we propose a regularizer specifically designed for multi-label learning. $\beta^{tt'}(\mathbf{x}) = P\left(y^{(t)} = +1|y^{(t)} \neq y^{(t')}, \mathbf{x}\right)$

Under the independence Assumption 1, the probability estimates of each label in $\{1, \dots, T\}$ are supervised independently, *i.e.* there is no explicit supervision for the joint prediction of multiple classes. This is consistent with the decomposition of the risk of (4) into a sum of t labelspecific risks. However, an example x still provides joint constraints on the probability estimates of different labels. To see this, consider any $x \in \mathcal{X}$ and pair (t, t') of labels with different values. The probability of $y^{(t)} = +1$ is then

Assumption 1. For any $\mathbf{x} \in \mathcal{X}$ and $i \neq j \in \{1, \dots, T\}$, $y^{(i)}$ and $y^{(j)}$ are independent given \mathbf{x} , i.e. $y^{(i)} \perp y^{(j)} | \mathbf{x}$.

$$= P\left(y^{(t)} = +1 | y^{(t')} \neq y^{(t')}, \mathbf{x}\right)$$

$$= \frac{P(y^{(t)} = +1, y^{(t')} = -1 | \mathbf{x})}{P\left(y^{(t)} = +1, y^{(t')} = -1 | \mathbf{x}\right) + P\left(y^{(t)} = -1, y^{(t')} = +1 | \mathbf{x}\right)}$$

$$= \frac{\eta^{(t)}(\mathbf{x})(1 - \eta^{(t')}(\mathbf{x}))}{\eta^{(t)}(\mathbf{x})(1 - \eta^{(t')}(\mathbf{x})) + (1 - \eta^{(t)}(\mathbf{x}))\eta^{(t')}(\mathbf{x})}$$

$$(18)$$

where the last equality follows from the Assumption 1. Let

$$\widehat{\beta}^{tt'} = \frac{\widehat{\eta}^{(t)}(1 - \widehat{\eta}^{(t')})}{\widehat{\eta}^{(t)}(1 - \widehat{\eta}^{(t')}) + \widehat{\eta}^{(t')}(1 - \widehat{\eta}^{(t)})}$$
(19)

be the plugin estimator of $\beta^{tt'}$. The following result shows that the accurate estimation of $\hat{\beta}^{tt'}$ is a necessary condition for the accurate estimation of $\hat{\eta}^{(t)}$ and $\hat{\eta}^{(t')}$.



Label Pair Regularizer

Lemma 1. For any $t \neq t'$, $\widehat{\beta}^{(tt')} = \beta^{(tt')}$ is a necessary condition for $\widehat{\eta}^{(t)} = \eta^{(t)}$ and $\widehat{\eta}^{(t')} = \eta^{(t')}$.

Since the example x can be seen as a calibration constraint for the estimation of $\hat{\beta}^{tt'}$ in addition to those that it already provides for the individual calibration of $\hat{\eta}^{(t)}$ and $\hat{\eta}^{(t')}$, this suggests that introducing calibration supervision on $\hat{\beta}^{tt'}$ can help improve the calibration of $\eta^{(t)}, \eta^{(t')}$. We leverage this observation by introducing a new *Label Pair Regularizer* (LPR) o calibrate the estimate $\hat{\beta}^{ij}$, implemented with the BCE loss

$$\mathcal{L}^{\text{LPR}}(\mathbf{x}, \mathbf{y}) = \frac{2}{T(T-1)} \sum_{t \neq t'} \mathbb{1}_{\substack{y^{(t)} = +1 \\ y^{(t')} = -1}} \left[-\log \widehat{\beta}^{tt'} \right].$$
(20)

Finally, in our proposed training approach, a multi-label DNN is trained by a joint optimization of SPA and LPR:

$$\mathcal{L}^{\text{overall}}(\mathbf{x}, \mathbf{y}) = \frac{1}{T} \sum_{t=1}^{T} \ell_{y^{(t)}, \Psi}(v^{(t)}(\mathbf{x})) + \lambda \mathcal{L}^{\text{LPR}}(\mathbf{x}, \mathbf{y}) \quad (21)$$

where λ is the multiplier balancing the two terms.



Multi-label Image Retrieval

A natural score for image ranking is then the posterior probability

$$s(\mathbf{x}) = P\left(\{y^{(j)} = +1\}_{j \in \mathcal{P}}, \{y^{(k)} = -1\}_{k \in \mathcal{N}} \middle| \mathbf{x}\right)$$
$$= \prod_{j \in \mathcal{P}} \eta^{(j)}(\mathbf{x}) \prod_{k \in \mathcal{N}} (1 - \eta^{(k)}(\mathbf{x}))$$
(22)

where the equality follows from the Assumption 1. This can be estimated by a probabilistic multi-label networks as

$$\widehat{s}(\mathbf{x}) = \prod_{j \in \mathcal{P}} \widehat{\eta}^{(j)}(\mathbf{x}) \prod_{k \in \mathcal{N}} (1 - \widehat{\eta}^{(k)}(\mathbf{x})).$$
(23)

.

C 1

. .. .

Figure 4. Qualitative results of multi-label image retrieval. Correct retrieval results are highlighted in green.

Experiments	
L	

		ECA-ResNet50-T				ViT-B/32			
		Accuracy Calibration		Accuracy		Calib	Calibration		
Dataset	Method	mAP@ \mathbf{y} \uparrow	mÅP@x↑	$\big \operatorname{ACE} \downarrow $	$\text{MCE}\downarrow$	mAP@ y ↑	mÅP@ x ↑	$\text{ACE}\downarrow$	MCE
	BCE	72.2	82.4	8.2	17.0	70.1	81.9	6.0	15.5
	TWL	77.2	87.9	17.2	33.6	76.2	88.0	14.2	31.3
	Focal	74.8	85.9	20.1	34.9	72.0	83.8	23.6	36.2
0000	Focal+ ϕ	-	-	11.1	20.3	-	-	8.5	17.9
COCO	ASY	77.5	88.2	30.6	46.0	76.4	87.7	29.5	48.6
	$ASY + \phi$	77.0	87.8	15.0	26.0	76.2	87.3	16.2	26.2
	SPA	77.8	88.4	5.3	12.0	76.8	87.9	4.8	10.7
	SPA + LPR	<u>77.7</u>	88.6	4.2	9.3	76.6	88.1	2.1	5.3
	BCE	85.1	92.4	6.9	13.8	86.9	92.7	7.8	16.7
	TWL	89.1	93.7	10.8	22.0	<u>90.1</u>	94.5	14.5	28.3
	Focal	87.4	93.3	19.0	35.8	88.4	93.6	16.2	31.7
VOC	Focal + ϕ	-	-	7.0	16.5	-	-	9.3	20.7
VUC	ASY	<u>89.6</u>	94.6	26.4	47.7	90.4	<u>94.7</u>	31.7	52.8
	$ASY + \phi$	89.1	94.3	12.9	24.7	89.8	94.2	15.5	33.0
	SPA	89.5	94.1	5.4	14.9	90.0	93.9	6.1	14.4
	SPA + LPR	89.9	94.3	4.9	11.0	90.4	94.8	5.5	12.7
	BCE	74.9	82.8	11.5	25.2	76.7	83.2	8.6	13.3
	TWL	79.9	85.5	14.3	29.0	81.6	<u>87.6</u>	16.7	32.3
	Focal	78.5	84.1	20.1	34.2	80.8	85.2	18.8	31.2
WIDED A	Focal + ϕ	-	-	9.3	21.0	-	-	7.4	16.5
WIDEK-A	ASY	80.6	86.0	24.2	35.2	82.2	87.8	22.8	36.3
	$ASY + \phi$	79.9	85.4	14.3	27.5	81.8	87.2	13.3	20.7
	SPA	80.1	85.8	5.3	11.2	82.0	<u>87.6</u>	4.3	10.5
	SPA + LPR	<u>80.3</u>	<u>85.9</u>	3.5	8.0	82.7	87.9	2.8	6.2
	BCE	46.3	78.8	9.7	16.6	48.4	80.5	8.2	14.8
	TWL	52.4	84.3	17.9	26.0	53.6	85.4	12.4	26.8
	Focal	48.0	81.0	24.0	40.5	50.1	83.2	24.6	35.7
VISPR	Focal + ϕ	-	-	8.6	18.2	-	-	9.6	14.1
v 151 K	ASY	51.6	84.0	28.8	44.2	53.0	85.0	27.7	45.2
	$ASY + \phi$	51.4	83.7	14.3	23.7	43.9	52.8	16.4	29.0
	SPA	52.4	84.5	5.8	12.1	53.2	85.3	5.9	10.2
	SPA + LPR	52.7	84.9	3.0	8.1	53.4	85.6	2.5	7.4



Experiments



Figure 3. Multi-label image retrieval mAP versus the number of search conditions $|\mathcal{P}| + |\mathcal{N}|$.





Ours









	mAP@ x ↑	mAP@ y ↑	ACE↓	MCE↓
BCE	72.2	82.4	8.2	17.0
BCE + LPR	72.4	82.9	6.5	14.8
Focal	74.8	85.9	20.1	34.9
Focal + LPR	75.1	86.3	19.7	36.4
ASY	77.5	88.2	30.6	46.0
ASY + LPR	77.2	88.1	29.4	44.7

Table 4. Effect of LPR on other CPE losses (COCO, ECA-ResNet50-T).



Thanks