

Background Data Resampling for Outlier-Aware Classification

Supplementary Material

Yi Li Nuno Vasconcelos
 University of California, San Diego
 {yli1898, nvasconcelos}@ucsd.edu

A. Probabilistic Interpretation

We show that under the training pipeline described in Section 3 of main text, the optimal classifier output is a smoothed version of the posterior class probabilities. Suppose the in-distribution examples are input-label pairs sampled from $X, Y \sim p(x, y)$ with $Y \in \{1, \dots, K\}$, while background examples are unlabeled images that follow $X_b \sim q(x)$. Consider a classifier which predicts K -way class probabilities $f(x; \theta)$ through a softmax mapping from logits $v(x; \theta) \in \mathbb{R}^K$:

$$f_k(x; \theta) = \frac{e^{v_k(x; \theta)}}{\sum_{j=1}^K e^{v_j(x; \theta)}}. \quad (1)$$

The classifier is trained on both in-distribution and background data to minimize the combined loss

$$L(\theta; p, q) = \mathbb{E}_{X, Y \sim p(\cdot, \cdot)}[L_{\text{cls}}(f(X; \theta); Y)] + \alpha \mathbb{E}_{X \sim q(\cdot)}[L_{\text{uni}}(f(X; \theta))]. \quad (2)$$

where we consider the popular instantiation of

$$L_{\text{cls}}(f(x; \theta); y) = -\log f_y(x; \theta); \quad (3)$$

$$L_{\text{uni}}(f(x; \theta)) = -\frac{1}{K} \sum_{k=1}^K \log f_k(x; \theta) - \log K. \quad (4)$$

Theorem 1. *Under the loss functions defined by (3) and (4), the objective (2) is minimized when*

$$f_k^*(x; \theta) = c(x)p_{Y|X}(k | x) + \frac{1 - c(x)}{K}, \quad (5)$$

$$c(x) = \frac{p_X(x)}{p_X(x) + \alpha q(x)}. \quad (6)$$

where $p_X(x) = \sum_{k=1}^K p(x, k)$ and $p_{Y|X}(k | x) = \frac{p(x, k)}{p_X(x)}$.

Proof. The expectations in (2) are expanded into

$$L = - \int p_X(x) \sum_{k=1}^K p_{Y|X}(k | x) \log f_k(x; \theta) dx - \frac{\alpha}{K} \int q(x) \sum_{k=1}^K \log f_k(x; \theta) dx - \alpha \log K. \quad (7)$$

Using the relation $\log f_k = v_k - \log \sum_j e^{v_j}$ that follows the definition of (1), we have $L = \int l(x) dx$ with

$$l(x) = -p_X(x) \left(\sum_k p_{Y|X}(k | x) v_k - \log \sum_k e^{v_k} \right) - \alpha q(x) \left(\frac{1}{K} \sum_k v_k - \log \sum_k e^{v_k} \right) - \alpha \log K; \quad (8)$$

taking its gradient w.r.t. classifier logits $v(x; \theta)$ gives

$$\frac{\partial l(x)}{\partial v_k} = -p_X(x) (p_{Y|X}(k | x) - f_k(x; \theta)) - \alpha q(x) \left(\frac{1}{K} - f_k(x; \theta) \right). \quad (9)$$

Setting the gradient to zero leads the optimal solution of $f_k(x; \theta)$ in the form of (5), which minimizes the integral L . \square

B. Mini-batch Dataset Resampling

We present the mini-batch version of the adversarial resampling algorithm described in Section 4.3 of main text. In practice, because training is done in stochastic batches, we expect that the data sampler would respect the weights by selecting examples with higher w more frequently. This is realized by decomposing the weighted loss into

$$L_{\text{out}} = \sum_{i=1}^{|\mathcal{D}_b|} \frac{w_i^\epsilon}{Z} \left(\frac{Z w_i^{1-\epsilon}}{\sum_j w_j} L_{\text{uni}}(f(x_i; \theta)) \right), \quad (10)$$

where $\epsilon \in (0, 1)$ controls the degree to which the SGD sampler depends on the weights $w^{(t)}$; $Z = \sum_i w_i^\epsilon$ is a normalization factor. Mini-batch selection is performed with probability

$$p_i(w; \epsilon) = \frac{w_i^\epsilon}{Z} = \frac{w_i^\epsilon}{\sum_j w_j^\epsilon}, \quad (11)$$

Algorithm 1: Adversarial resampling, stochastic mini-batch version.

Input: ID dataset \mathcal{D} , background dataset \mathcal{D}_b , batch size n , sampler coefficient ϵ , pre-trained classifier θ , learning rate η_θ , η_w , loss coefficient α , total iterations T

Initialize: $w^{(0)} \leftarrow [1, \dots, 1], \theta^{(0)} \leftarrow \theta$;

for $t = 0, \dots, T - 1$ **do**

 Sample in-dist. mini-batch $\mathcal{B}^{(t)} \leftarrow \{(x_i, y_i)\}_{i=1}^n$ from \mathcal{D} with uniform probabilities;

 Sample background mini-batch $\mathcal{B}_b^{(t)} \leftarrow \{x_i^b\}_{i=1}^n$ from \mathcal{D}_b with probabilities $p_i(w^{(t)}; \epsilon)$ using (11);

 Compute ID loss $l_{in}^{(t)} \leftarrow L_{in}(\theta^{(t)}; \mathcal{B}^{(t)})$;

 Compute OOD loss $l_{out}^{(t)} \leftarrow L_{out}(\theta^{(t)}; \mathcal{B}_b^{(t)}, \tilde{w}^{(t)})$ with weights $\tilde{w}_i^{(t)} = r_i(w^{(t)}; \epsilon)$ using (12);

 Update classifier

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_\theta \nabla_{\theta^{(t)}} \left(l_{in}^{(t)} + \alpha l_{out}^{(t)} \right);$$

 Update weights

$$w^{(t+1)} \leftarrow w^{(t)} + \eta_w \nabla_{w^{(t)}} l_{out}^{(t)};$$

Output: Resampling weights $w^{(T)}$.

and the uniformity loss for each example within the batch is reweighted by

$$r_i(w; \epsilon) = \frac{Z w_i^{1-\epsilon}}{\sum_j w_j} = \frac{\sum_j w_j^\epsilon}{\sum_j w_j} w_i^{1-\epsilon}. \quad (12)$$

C. Implementation Details

Instead of directly optimizing the weights $\{w_i\}$ subject to the constraint $w_i > 0$, we use softplus reparameterization $w_i = \log(1 + e^{w_i})$ and optimize $\{w_i\}$ unconstrained to ensure that all weights are positive. Mini-batch based resampling, as described in Algorithm 1, is implemented using the hyperparameter choices in Table 1 (see Section 5.1 of main text for additional training details).

Input pre-processing is done as follows: All images are resized to 32 pixels on the shorter edge, randomly cropped at 32×32 before zero-padding 4 pixels each side, then horizontally flipped with probability 50% (except for SVHN dataset which contains digits). Pixel values are scaled to $[0, 1]$ and normalized using mean 0.5 and standard deviation 0.25 regardless of datasets, since the source of test examples is always assumed unknown.

This work is implemented using the PyTorch library [2]; the Python code for 1) pre-training classifiers, 2) adversarial resampling, 3) retraining with resampled background data, and 4) evaluating OOD detection of trained classifiers, are provided along with this PDF.

Sampler coefficient ϵ	0.5
Sampler batch size n	128
Reweighting step size η_w	100
Classifier step size η_θ	0.001
Classifier momentum β	0.9
Classifier weight decay λ	0.0005

Table 1. Hyperparameters for adversarial resampling.

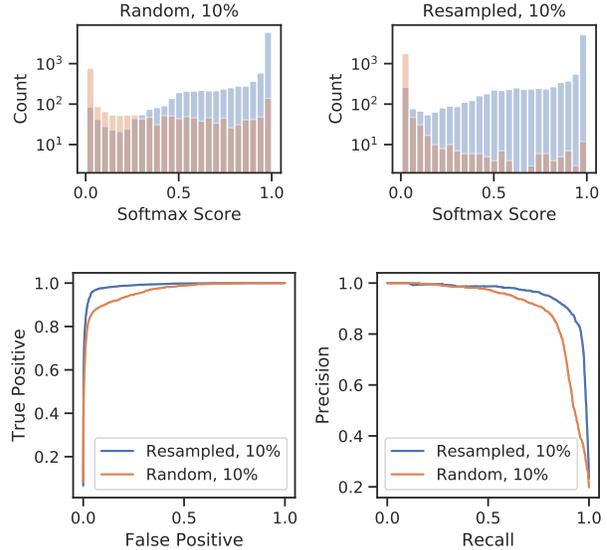


Figure 1. Qualitative result w/ in-distribution $\mathcal{D} = \text{Tiny ImageNet}$ and out-of-distribution $\mathcal{D}_o = \text{SVHN}$, where using resampled data provides significant improvement over randomly selecting background data. **Top:** Separation of ID (blue) and OOD (red) examples by softmax score. **Bottom:** Detection ROC and PR curves.

D. OOD Detection Results by Dataset

Table 2 is a breakdown of the OOD detection results for each out-of-distribution dataset used in the main text. Training on the resampled data provides significant improvement over using random background data, even outperforming full background in detecting challenging OOD examples like SVHN; the most dramatic difference is observed when Tiny ImageNet is used as in-distribution dataset, as visualized in Figure 1. Notice from the histograms that using resampled vs. random background dataset does not affect the softmax scores of ID examples as much as it does for the OOD examples, which is desired as we do not want to sacrifice in-distribution classification performance.

E. Classification Accuracy

Table 3 compares the in-distribution classification accuracy of trained models. Training with background samples improves baseline accuracy on CIFAR-10 and Tiny ImageNet, possibly due to the additional training data that reduces overfitting, but differences are small.

ID data \mathcal{D}	OOD data \mathcal{D}_o	Background data \mathcal{D}_b			
		None [1] ($\gamma = 0$)	Full ($\gamma = 100\%$)	Random ($\gamma = 10\%$)	Resampled ($\gamma = 10\%$)
CIFAR-10	Gaussian	7.13 _{.12} / 96.93 _{.02} / 76.94 _{.15}	1.18 _{.03} / 99.45 _{.00} / 91.82 _{.02}	2.23 _{.09} / 99.07 _{.02} / 88.78 _{.08}	1.78 _{.06} / 99.19 _{.01} / 89.51 _{.05}
	Uniform	7.15 _{.18} / 96.80 _{.05} / 74.99 _{.34}	2.05 _{.09} / 99.15 _{.02} / 89.59 _{.13}	4.29 _{.11} / 98.38 _{.02} / 84.52 _{.09}	2.36 _{.04} / 98.95 _{.01} / 87.16 _{.08}
	Textures	62.41 _{2.23} / 85.54 _{.27} / 53.36 _{.48}	1.16 _{.04} / 99.63 _{.03} / 97.88 _{.08}	1.30 _{.10} / 99.54 _{.02} / 96.96 _{.07}	1.10 _{.03} / 99.71 _{.01} / 98.05 _{.04}
	LSUN	49.10 _{1.79} / 86.54 _{.19} / 52.14 _{.65}	0.58 _{.03} / 99.83 _{.00} / 98.76 _{.01}	0.66 _{.01} / 99.76 _{.01} / 97.99 _{.05}	0.60 _{.03} / 99.84 _{.00} / 98.83 _{.03}
	SVHN	20.46 _{.72} / 92.85 _{.21} / 66.48 _{.72}	7.27 _{.72} / 98.69 _{.09} / 94.03 _{.25}	7.58 _{.70} / 98.43 _{.06} / 91.74 _{.16}	4.90 _{.29} / 98.79 _{.03} / 92.96 _{.11}
	Places	55.07 _{3.15} / 85.65 _{.50} / 52.73 _{.87}	1.03 _{.05} / 99.70 _{.03} / 98.25 _{.07}	1.04 _{.04} / 99.63 _{.02} / 97.55 _{.06}	0.87 _{.07} / 99.76 _{.02} / 98.42 _{.08}
	Average	31.45 / 90.72 / 62.77	2.21 / 99.41 / 95.06	2.85 / 99.14 / 92.92	1.94 / 99.37 / 94.16
CIFAR-100	Gaussian	34.95 _{.33} / 77.98 _{.12} / 28.15 _{.12}	4.04 _{.04} / 97.85 _{.02} / 78.31 _{.17}	8.36 _{.08} / 95.78 _{.03} / 66.12 _{.08}	2.20 _{.03} / 98.64 _{.01} / 84.22 _{.15}
	Uniform	17.73 _{.25} / 91.81 _{.06} / 53.87 _{.32}	5.73 _{.10} / 96.74 _{.02} / 70.16 _{.13}	10.97 _{.30} / 94.30 _{.11} / 59.37 _{.37}	2.39 _{.00} / 98.41 _{.01} / 81.14 _{.14}
	Textures	67.42 _{.65} / 73.61 _{.37} / 32.34 _{.71}	7.07 _{.47} / 98.06 _{.12} / 90.46 _{.22}	9.94 _{.35} / 97.26 _{.06} / 87.76 _{.11}	4.66 _{.15} / 98.56 _{.06} / 90.42 _{.22}
	LSUN	74.16 _{1.36} / 70.54 _{.50} / 27.77 _{.65}	2.15 _{.24} / 99.38 _{.04} / 95.74 _{.16}	2.40 _{.03} / 99.18 _{.04} / 93.00 _{.13}	2.84 _{.26} / 99.04 _{.04} / 91.18 _{.18}
	SVHN	63.39 _{1.16} / 73.43 _{.19} / 30.31 _{.23}	29.04 _{.67} / 91.04 _{.11} / 57.86 _{.45}	31.26 _{.98} / 91.09 _{.09} / 58.56 _{.32}	22.70 _{.87} / 92.97 _{.07} / 64.01 _{.39}
	Places	71.20 _{.94} / 72.89 _{.37} / 31.42 _{.54}	3.04 _{.11} / 99.08 _{.05} / 94.42 _{.09}	3.52 _{.20} / 98.86 _{.08} / 92.23 _{.16}	3.59 _{.29} / 98.92 _{.03} / 91.51 _{.20}
	Average	54.81 / 76.71 / 33.98	8.51 / 97.03 / 81.16	11.08 / 96.08 / 76.17	6.40 / 97.76 / 83.75
Tiny-ImageNet	Gaussian	49.73 _{.21} / 74.38 _{.21} / 26.65 _{.20}	0.30 _{.01} / 99.83 _{.00} / 97.34 _{.03}	0.15 _{.00} / 99.92 _{.00} / 98.93 _{.03}	0.11 _{.01} / 99.98 _{.00} / 99.90 _{.00}
	Uniform	26.63 _{.51} / 90.35 _{.08} / 54.33 _{.33}	0.20 _{.00} / 99.86 _{.00} / 98.03 _{.03}	0.10 _{.00} / 99.96 _{.00} / 99.35 _{.03}	0.09 _{.01} / 99.99 _{.00} / 99.92 _{.00}
	Textures	75.26 _{.26} / 66.00 _{.43} / 25.53 _{.36}	0.85 _{.09} / 99.67 _{.02} / 99.05 _{.04}	0.89 _{.14} / 99.41 _{.06} / 98.63 _{.10}	0.53 _{.06} / 99.81 _{.04} / 99.34 _{.06}
	LSUN	81.30 _{.67} / 61.78 _{.32} / 21.51 _{.20}	0.00 _{.00} / 99.97 _{.02} / 99.93 _{.02}	0.00 _{.00} / 99.96 _{.02} / 99.90 _{.02}	0.00 _{.00} / 99.98 _{.00} / 99.93 _{.01}
	SVHN	63.67 _{1.33} / 74.46 _{.48} / 32.62 _{.64}	21.25 _{2.93} / 97.01 _{.22} / 91.92 _{.39}	47.84 _{1.97} / 89.97 _{.52} / 74.70 _{.76}	6.74 _{.75} / 98.11 _{.14} / 94.16 _{.18}
	Places	77.88 _{1.29} / 65.08 _{.40} / 23.72 _{.37}	0.00 _{.00} / 99.98 _{.00} / 99.91 _{.01}	0.02 _{.00} / 99.92 _{.04} / 99.79 _{.07}	0.01 _{.01} / 99.97 _{.01} / 99.90 _{.03}
	Average	62.41 / 72.01 / 30.73	3.77 / 99.39 / 97.70	8.17 / 98.19 / 95.22	1.25 / 99.64 / 98.86

Table 2. Breakdown of OOD detection performance (FPR95 ↓ / AuROC ↑ / AuPR ↑), reported in Mean_{Std} over 5 runs.

BG data $\mathcal{D}_b \setminus$ ID data \mathcal{D}	C-10	C-100	TIN
None, $\gamma = 0$	94.19	74.77	52.98
Full, $\gamma = 100\%$	94.29	73.76	54.82
Random, $\gamma = 10\%$	94.33	73.81	54.99
Resampled, $\gamma = 10\%$	94.26	73.75	54.21

Table 3. Classification accuracy% of trained models.

References

- [1] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations (ICLR)*, 2017. 3
- [2] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 2