

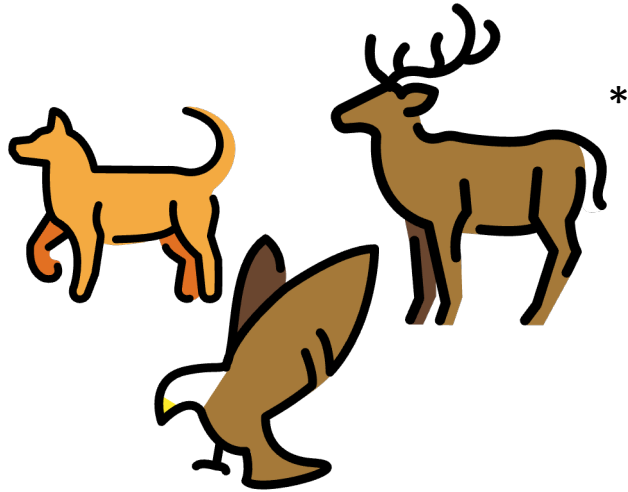
WeiPer: OOD Detection using Weight Perturbations of Class Projections

Maximilian Granz, Manuel Heurich, Tim Landgraf

Out of Distribution Detection

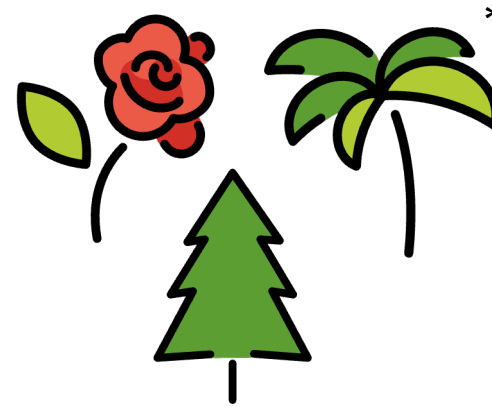
...

$x_{ID} \in X_{ID}$



Animals

$x_{OOD} \in X_{OOD}$



Plants

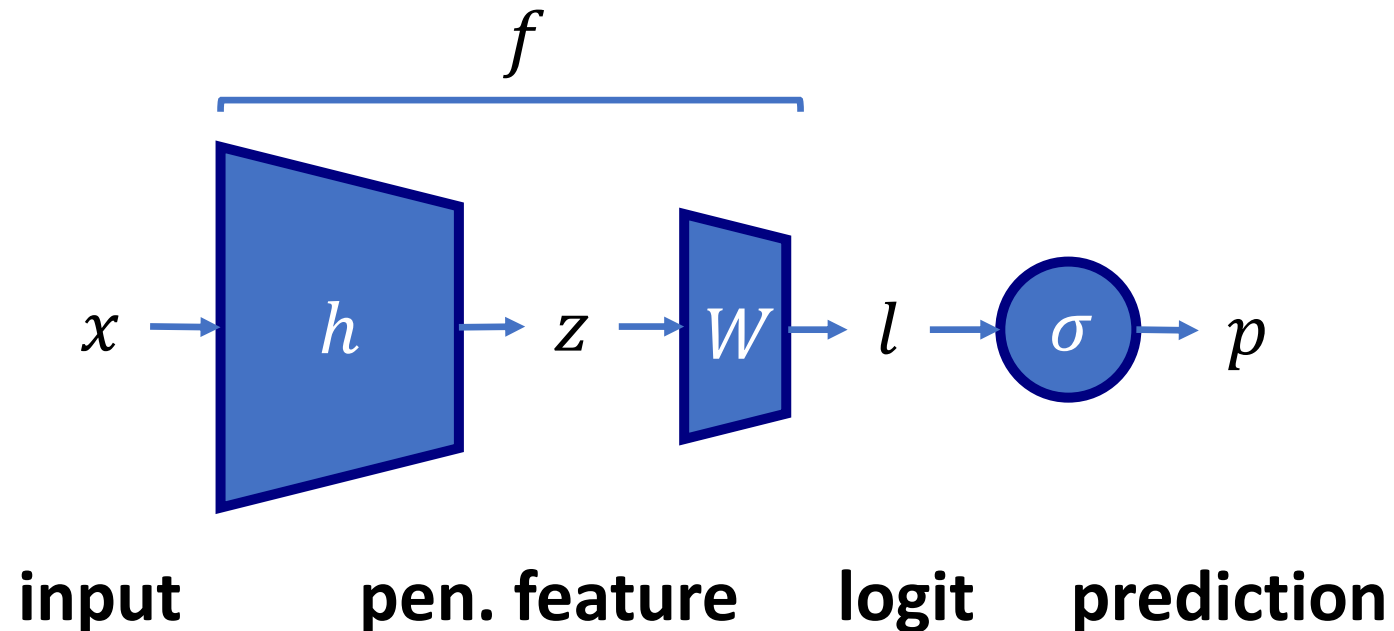
Out of Distribution Detection

Define a **score function** $S(x) \rightarrow \mathbb{R}$

$$O(x) = \begin{cases} \text{ID}, & \text{if } S(x) > \lambda \\ \text{OOD}, & \text{otherwise} \end{cases}$$

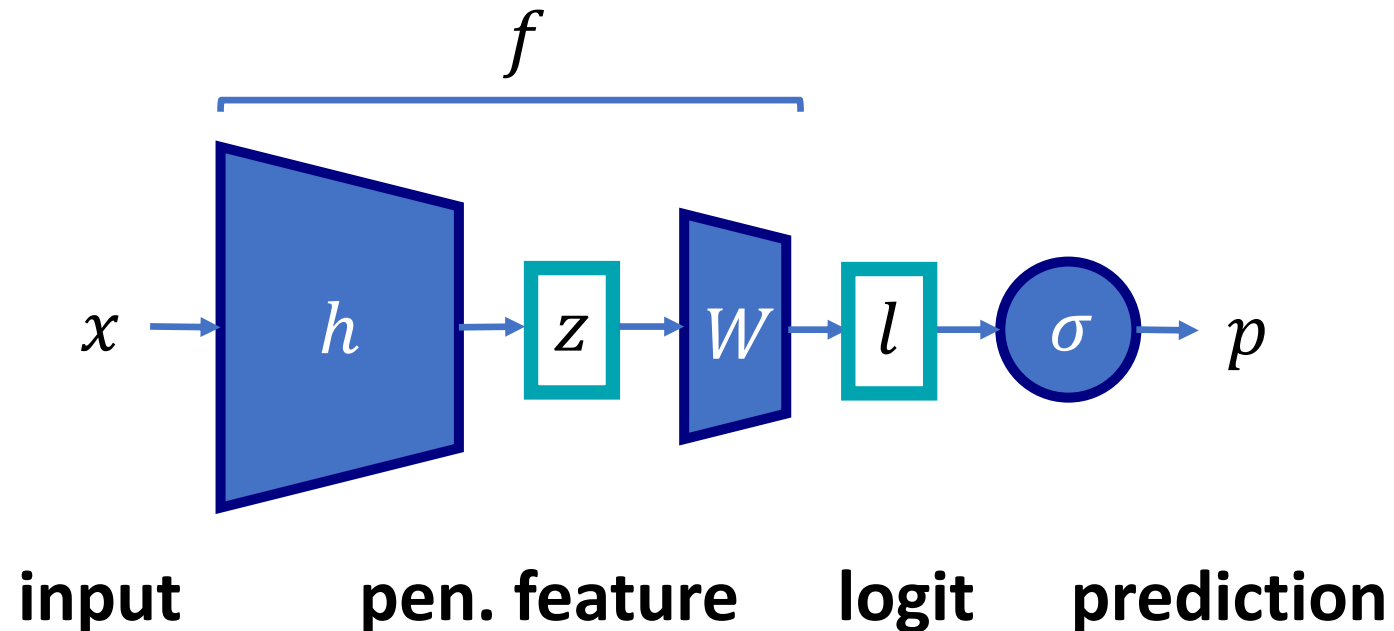
Out of Distribution Detection

Let σ be the softmax function and f be an ANN classifier



Out of Distribution Detection

Let σ be the softmax function and f be an ANN classifier



Maximum Softmax Probability (MSP)

One example for a **score function** is MSP

$$\text{MSP}(x) = \max_{i=1,\dots,C} \sigma(l)_i$$

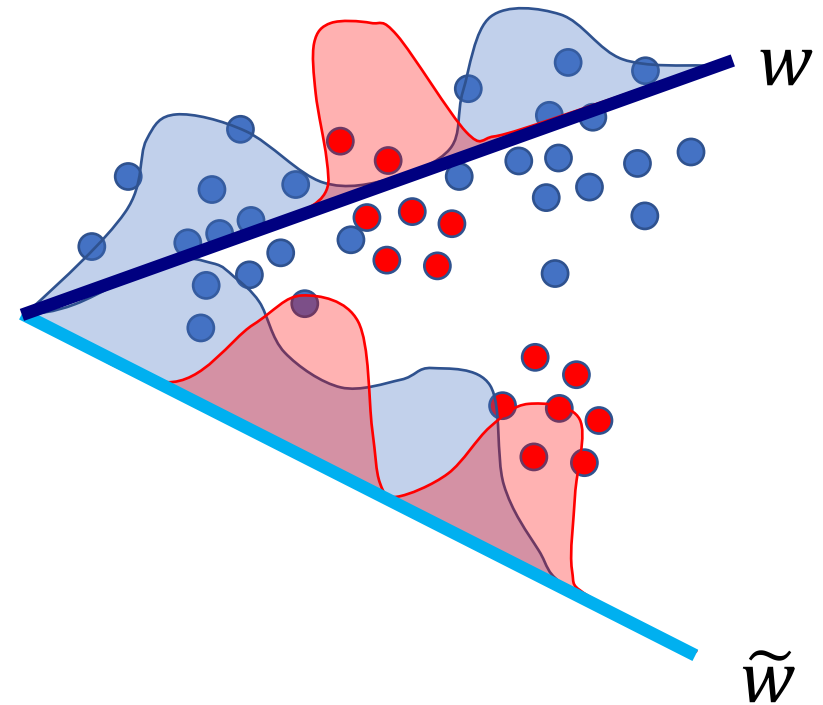
How can we extract more information from
the penultimate feature space?

$$\text{MSP}(x) = \max_{i=1,\dots,C} \sigma(Wz)_i$$

Creating a richer Representation

Cramer-Wold Thm.

If we sample all projections $\tilde{w} \in \mathbb{R}^k$, the penultimate distribution Z is fully determined by its one-dimensional projections $\tilde{w}(Z)$.

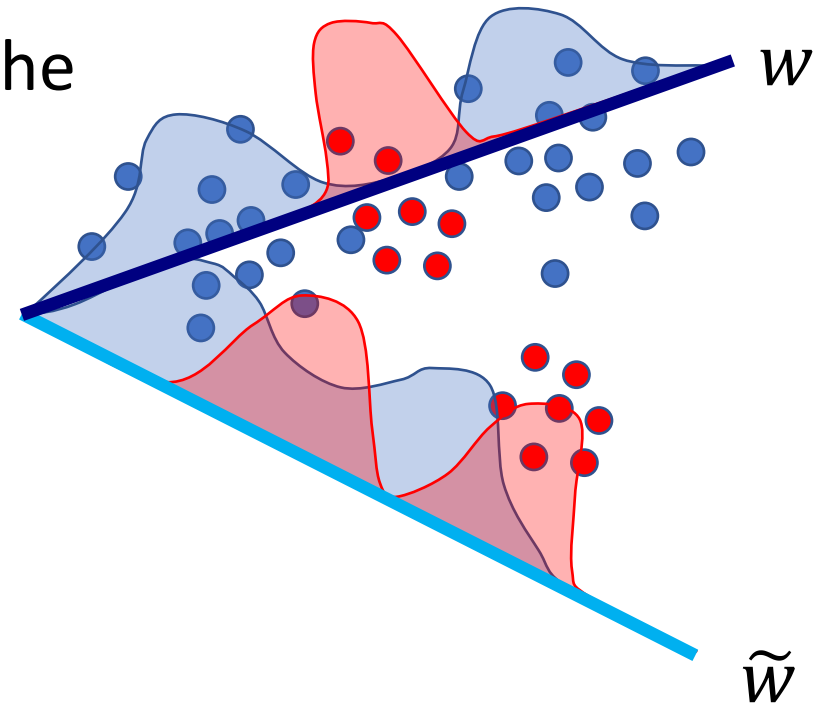


Creating a richer Representation

Cramer-Wold Thm.

If we sample all projections $\tilde{W} \in \mathbb{R}^{C \times k}$, the penultimate distribution Z is fully determined by its **C-dimensional projections** $\tilde{W}(Z)$.

$$S(x) := \frac{1}{r} \sum_{j=1}^r \max_i \sigma(\tilde{W}_j z)$$



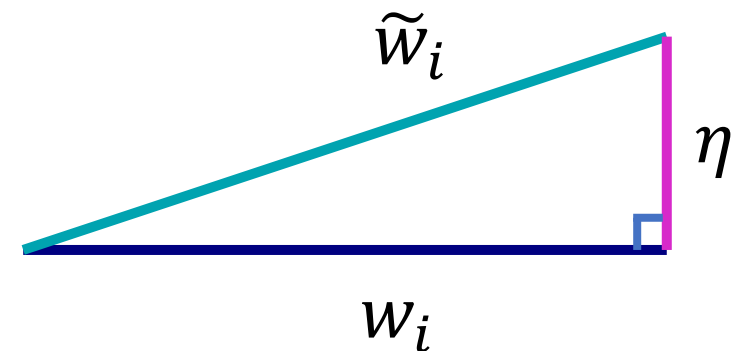
Weight Perturbations (WeiPer)

- Define projections that **correlate** with the class projections W .
- We control the **similarity** with a hyperparameter δ .

$\widetilde{w}_i \sim w_i$ row vector of W

$$\widetilde{w}_i = w_i + \delta \cdot \eta, \quad \eta \perp w, \quad \eta \in S_K$$

$$\angle(\widetilde{w}_i, w_i) = \arctan \delta$$

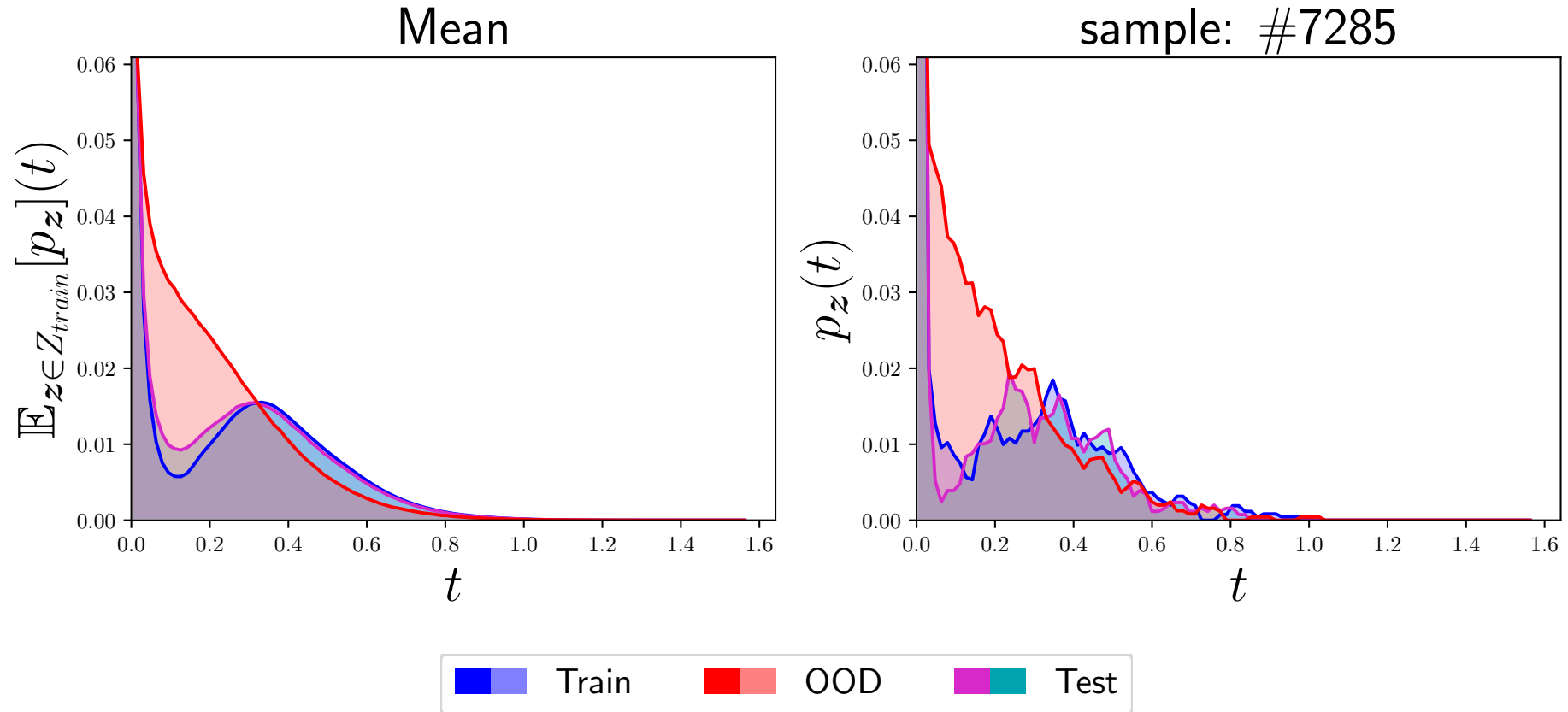


MSP on WeiPer space

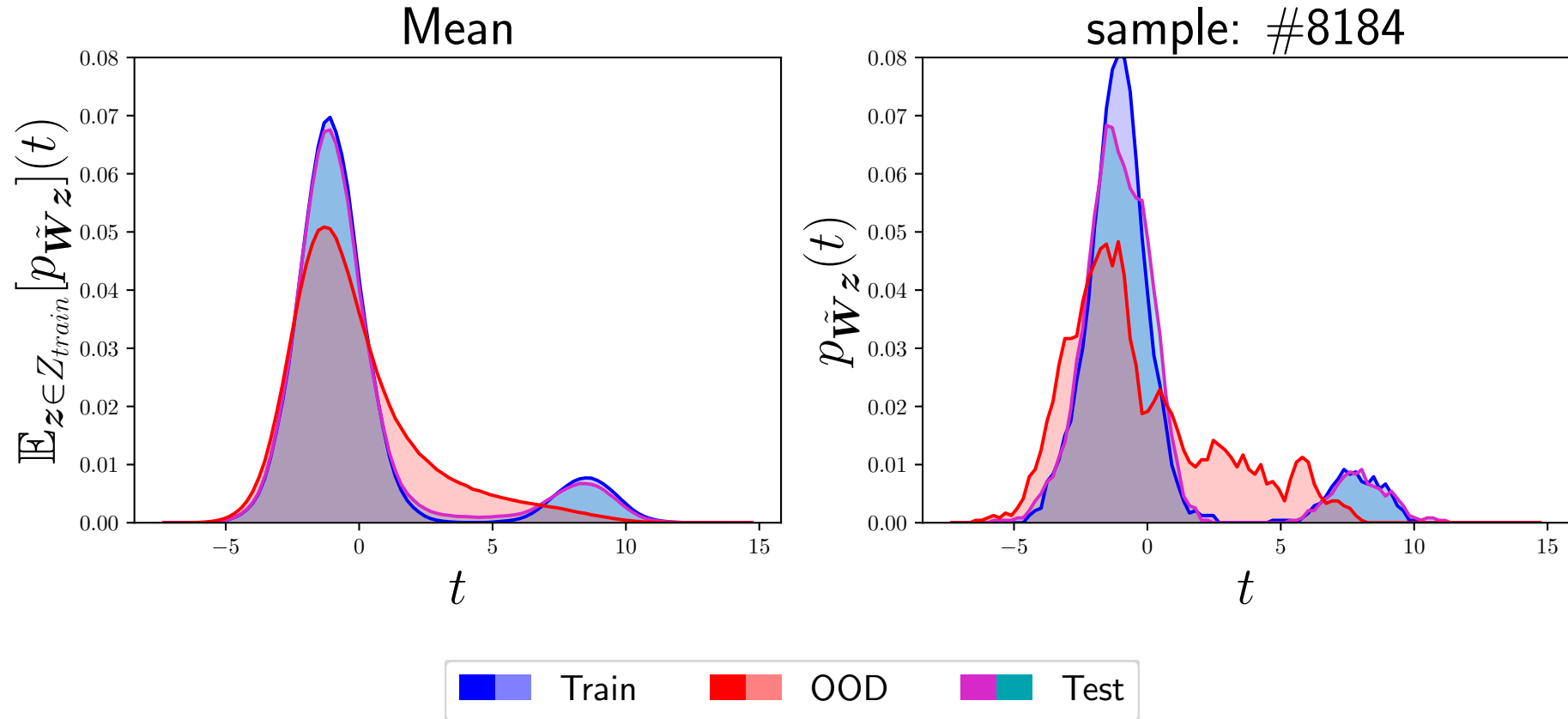
Repeat perturbing each row vector r times to get $\tilde{W}_1, \dots, \tilde{W}_r$

$$\text{MSP}_{\text{WeiPer}}(z) := \frac{1}{r} \sum_{j=1}^r \text{MSP}(\tilde{W}_j z)$$

Penultimate Layer Distribution

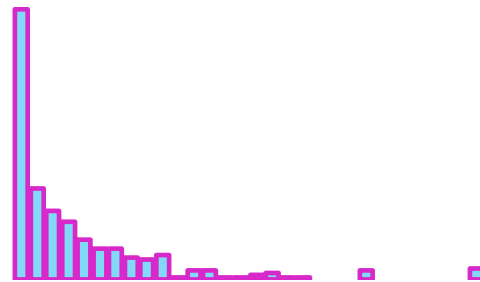


WeiPer space Distribution

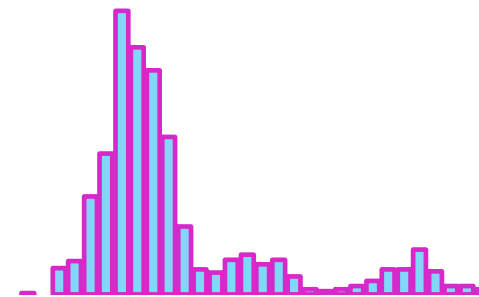


KL Divergence Score Function

1. Calculate the **binned density** p_z of the **penultimate activations** and $p_{\tilde{W}_z}$ of the **WeiPer outputs** of z .



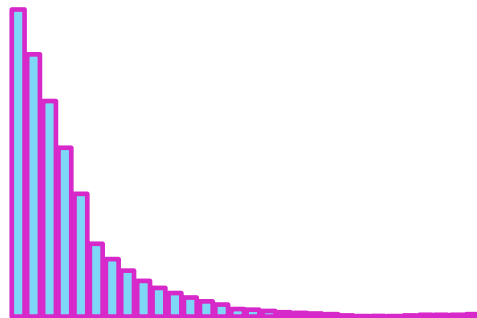
p_z



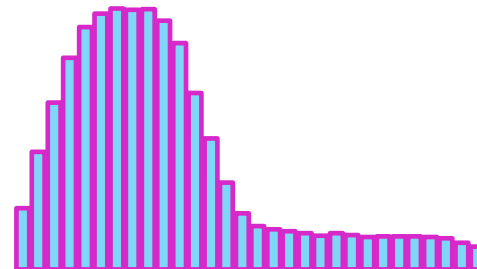
$p_{\tilde{W}_z}$

KL Divergence Score Function

2. Apply **normalization and smoothing** T to p_z and $p_{\tilde{w}_z}$.



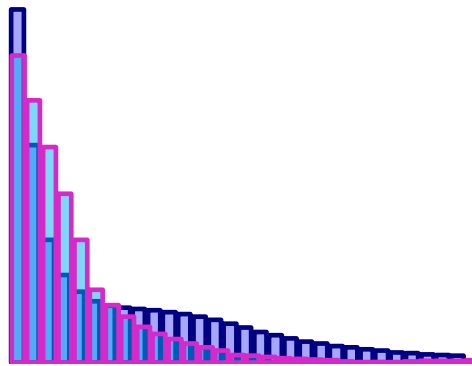
$T(p_z)$



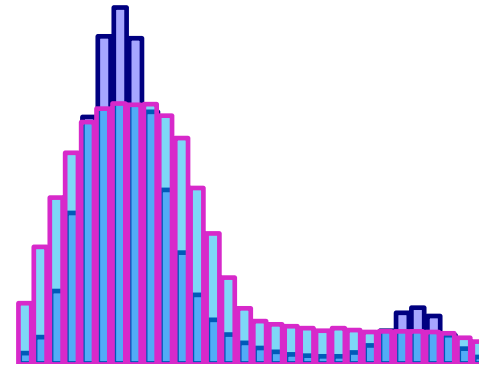
$T(p_{\tilde{w}_z})$

KL Divergence Score Function

3. Calculate a **KL-Divergence based distance** D_{KL} from each sample distribution to the mean distribution.



$$D_{KL}(T(p_z), \mathbb{E}(p_z))$$



$$D_{KL}(T(p_{\tilde{w}_z}), \mathbb{E}(p_{\tilde{w}_z}))$$

KL Divergence Score Function

3. Combine all **scores** into one:

$$\text{WeiPer+KLD}(x) := D_{KL}(T(p_z), \mathbb{E}(p_z)) + \lambda_1 D_{KL}(T(p_{\tilde{W}_z}), \mathbb{E}(p_{\tilde{W}_z})) - \lambda_2 \text{MSP}_{\text{WeiPer}}(z)$$

Results on OpenOOD

CIFAR10		CIFAR100		ImageNet (ResNet)		ImageNet(ViT)	
Near	Far	Near	Far	Near	Far	Near	Far
<u>90.54</u>	93.12	81.37	79.01	80.05	<u>95.54</u>	75.0	90.32
<u>(2nd)</u>	(3rd)	(1st)	(13th)	(1st)	<u>(2nd)</u>	(6th)	(8th)

Thank You

Code: <https://github.com/mgranz/weiper>
Paper #96381