Hyperspherical Prototype Networks Pascal Mettes, Elise van der Pol, and Cees Snoek 👰 🏠

Hyperspheres as output spaces

In hyperspherical output spaces, classes can be described by prototypes that are defined a priori with large margin separation and prior knowledge.





How would you position two classes? For cat, airplane tiger, what should go where?

Benefits over softmax cross-entropy

- \checkmark Enables an embedding of inductive biases prior to learning.
- \checkmark Freedom to choose any output dimensionality.
- \checkmark Unification with regression using the same loss.

Benefits over other prototype approaches

 \checkmark No chicken-egg problem (prototypes and inputs depend on each other). \checkmark No need for expensive prototype updating.



Regression:

- Retain two prototypes on extreme ends.
- Prototypes state lower and upper bound.
- Same loss as for classification.
- Multi-task learning possible in same space.





Match examples with fixed class prototypes based on their angle:

 $\mathcal{L}_{c} = \sum_{i=1}^{N} (1 - \cos \theta_{\mathbf{z}_{i}, \mathbf{p}_{y_{i}}})^{2} = \sum_{i=1}^{N} (1 - \frac{|\mathbf{z}_{i} \cdot \mathbf{p}_{y_{i}}|}{||\mathbf{z}_{i}||})^{2}$

Prototypes not updated, backpropagation through examples.

 $\frac{d}{\mathbf{z}_i}(1 - \cos\theta_i)^2 = 2(1 - \cos\theta_i) \left(\frac{\cos\theta_i \cdot \mathbf{z}_i}{||\mathbf{z}_i||^2} - \frac{\mathbf{p}_{y_i}}{||\mathbf{z}_i||||\mathbf{p}_{y_i}||}\right)$

During inference, simply select class of nearest prototype: $c^* = \arg \max_{c \in C} \left(\cos \theta_{f_{\phi}(\mathbf{\tilde{x}}), \mathbf{p}_c} \right)$

Prototype construction

Optimally separting arbitrary numbers of points for any dimensionality is an open mathematical problem. We resort to an approximation through optimization with the following objective:

 $\mathbf{P}^* = \arg \min_{\mathbf{P}' \in \mathbb{P}}$

Separating nearest pair per step is inefficient. We opt to separate each prototype from its neighbour:

 $\mathcal{L}_{\mathrm{HP}} = \frac{1}{K} \sum_{i=1}^{K} \max_{j \in C} \mathbf{M}_{ij}, \quad \mathbf{M} = \hat{\mathbf{P}} \hat{\mathbf{P}}^{T} - 2\mathbf{I}, \quad \text{s.t. } \forall_i ||\hat{\mathbf{P}}_i|| = 1$

We also propose an extra ranking loss to preserve privileged class similarities:

Number of triplets

 $\left[\!\left[\cos\theta_{\mathbf{w}_{i},\mathbf{w}_{i}}\geq\cos\theta_{\mathbf{w}_{i},\mathbf{w}_{k}}\right]\!\right]$ Word embedding similarities

Regression and multi-task learning



Maintain two prototypes pointing in opposite directions. Prototypes denote lower and upper regression bound:

Joint regression and classification possible in same space. One axis for regression, other axes for classification. No need to balance tasks as they have the same loss.

Classification



$$\left(\max_{(k,l,k\neq l)\in C}\cos\theta_{(\mathbf{p}'_k,\mathbf{p}'_l)}\right)$$

 $\mathcal{L}_{\rm PI} = \frac{1}{|T|} \sum_{(i,j,k)\in T} -\bar{S}_{ijk} \log S_{ijk} - (1 - \bar{S}_{ijk}) \log(1 - S_{ijk})$

Sigmoid over prototype similarities

 $\mathcal{L}_{\rm r} = \sum_{i=1}^{N} (r_i - \cos\theta_{\mathbf{z}_i, \mathbf{p}_u})^2, \qquad r_i = 2 \cdot \frac{y_i - v_l}{v_u - v_l} - 1$

Regression value as interpolation between bounds







This paper with privileged info.



Predict creation year for 20th century paintings. Our approach is stable across learning rates and output spaces.



Evaluation: OmniArt Classify style, regress creation year. We do both, without task weighting.





Classification results

	CIFAR-100				ImageNet-200			
imensions	10	25	50	100	25	50	100	200
ne-hot	-	-	-	62.1 ± 0.1	-	-	-	33.1 ± 0.6
ord2vec	29.0 ± 0.0	44.5 ± 0.5	54.3 ± 0.1	57.6 ± 0.6	20.7 ± 0.4	27.6 ± 0.3	29.8 ± 0.3	30.0 ± 0.4
nis paper	51.1 ± 0.7	$\textbf{63.0} \pm 0.1$	$\textbf{64.7} \pm 0.2$	$\textbf{65.0} \pm 0.3$	$\textbf{38.6} \pm 0.2$	$\textbf{44.7} \pm 0.2$	$\textbf{44.6} \pm 0.0$	44.7 ± 0.3

Our approach is preferred over one-hot and word2vec embeddings. Separation matters especially when output spaces are small. Above results for ResNet-32, results hold for DenseNet-121.







This paper vs cross-entropy.

Regression results



	1890
MAL	
	de la

	Omniart				
Output space	\mathbf{S}^1		\mathbb{S}^2		
Learning rate	1e-2	1e-3	1e-2	1e-3	
MSE	$210.7 \pm \text{140.1}$	110.3 ± 0.8	339.9 ± 0.0	109.9 ± 0.5	
This paper	84.4 ± 10.7	$\textbf{76.3} \pm 5.6$	$\textbf{82.9} \pm 1.9$	$\textbf{73.2} \pm 0.6$	

Multi-task results

Visualization: Rotated MNIST

Classify on xy-plane, regress on z-axis. Tasks are separated on same sphere.

Task weight	0.01	0.10	0.25	0.50	0.90
Creation year (MAE ↓) MTL baseline This paper	262.7 65.2	344.5 64.6	348.5 64.1	354.7 68.3	352.3 83.6
Art style (acc ↑) MTL baseline This paper	44.6 46.6	47.9 51.2	49.5 54.5	47.2 52.6	47.1 51.4

