

Joint Learning of Hyperbolic Label Embeddings for Hierarchical Multi-label Classification - HiddenN

Soumya Chatterjee¹ Ayush Maheshwari¹ Ganesh Ramakrishnan¹ Saketha Nath Jagarlapudi²
¹{soumya, ayusham, ganesh}@cse.iitb.ac.in, ²saketha@iith.ac.in

¹ Indian Institute of Technology Bombay ² Indian Institute of Technology Hyderabad

Problem Statement

Given a set of documents and labels, classify the documents into multiple labels respecting the hierarchy. For eg., *Voice Recognition Is Improving, but Don't Stop the Elocution Lessons* - Labels are *Top/News/Technology*.

Assumption : Label hierarchy is not available.

Key Contributions

- Our approach, HIDDEN learns label embeddings using the joint optimisation approach
- HIDDEN sometimes generalizes even better than state-of-the-art hierarchical multi-label classifiers that have complete access to the true label hierarchy
- We show significant improvement over classical multi-label classification methods as well as baselines that employ hyperbolic label embeddings.

Background: Poincaré Embeddings

- Let $\mathcal{B}^n = \{x \in \mathbb{R}^n \mid \|x\| < 1\}$ be the open n -dimensional unit ball, where $\|\cdot\|$ is the Euclidean 2 norm.
- The Poincaré ball model is a Riemannian Manifold (\mathcal{B}^n, g_x) , the open unit ball equipped with the Riemannian metric tensor $g_x = \left(\frac{2}{1-\|x\|^2}\right)^2 g^E$, where $x \in \mathcal{B}^d$ and g^E is the Euclidean metric tensor.
- The geodesic distance between two points $u, v \in \mathcal{B}^d$ is given as

$$d(u, v) = \operatorname{arcosh} \left(1 + 2 \frac{\|u - v\|^2}{(1 - \|u\|^2)(1 - \|v\|^2)} \right)$$

Source Code

<https://github.com/soumyac1999/hyperbolic-label-emb-for-hmc/>

Our Model : HiddenN

- L are nodes of a fixed hierarchy but hierarchy is unknown to our model.
- Document Model** - $\mathcal{F}_w(D) \in \mathbb{R}^n$
- Label Embedding Model** - $\mathcal{G}_\Theta(l) \equiv \Theta * y^l = \Theta_l$, where $\Theta \in \mathbb{R}^{n \times L}$
- Projection of $\Theta(l)$ into Poincaré manifold to get $\Pi(\Theta_l)$
 $\Pi(x) = \frac{x}{1 + \sqrt{1 + \|x\|^2}}$
- Alignment Model**: $\hat{y}_D^l(w, \Theta) \equiv \sigma(\mathcal{F}_w(D)^\top \Theta_l)$

Joint Learning

- First Term (Cross Entropy Loss for Classification)** -

$$\mathcal{L}_1(w, \Theta) = \sum_{i=1}^m \sum_{l=1}^L [y_i^l \log(\hat{y}_i^l(w, \Theta)) + (1 - y_i^l) \log(1 - \hat{y}_i^l(w, \Theta))]$$

- Second Term (Geodesic Distance Loss for Label Embeddings)** -

$$\mathcal{L}_2(\Theta) = \sum_{\substack{l, l' \in L \\ l \neq l'}} \log \left(\frac{e^{-d(\Pi(\Theta_l), \Pi(\Theta_{l'}))}}{\sum_{z \in (l, l')} e^{-d(\Pi(\Theta_l), \Pi(\Theta_z))}} \right)$$

- Overall objective function**

$$\mathcal{L}(w, \Theta) = \mathcal{L}_1(w, \Theta) + \lambda \mathcal{L}_2(\Theta) \quad (1)$$

- Inference**: Labels with $\hat{y}_D^l(\hat{w}, \hat{\Theta}) > 0.5$

Variants of HiddenN

- HIDDEN_{jnt} - $(w_{\text{jnt}}, \Theta_{\text{jnt}}) \in \arg \min_{w, \Theta} \mathcal{L}(w, \Theta)$
- HIDDEN_{cas}
 - \mathcal{L}_2 is minimized to obtain label embeddings $\hat{\Theta}_{\text{cas}} \in \arg \min_{\Theta} \mathcal{L}_2(\Theta)$.
 - These are then used in \mathcal{L}_1 to obtain document parameters: $\hat{w}_{\text{cas}} \in \arg \min_w \mathcal{L}_1(w, \hat{\Theta}_{\text{cas}})$.
- HIDDEN_{flt} - Θ_{flt} is fixed to the identity matrix
- HIDDEN_{auc} $\mathcal{L}_{2\text{Euc}}(\Theta) = \sum_{\substack{l, l' \in L \\ l \neq l'}} \log \left(\frac{e^{-\|\Theta_l - \Theta_{l'}\|_2}}{\sum_{z \in (l, l')} e^{-\|\Theta_l - \Theta_z\|_2}} \right)$

Synthetic Experiments

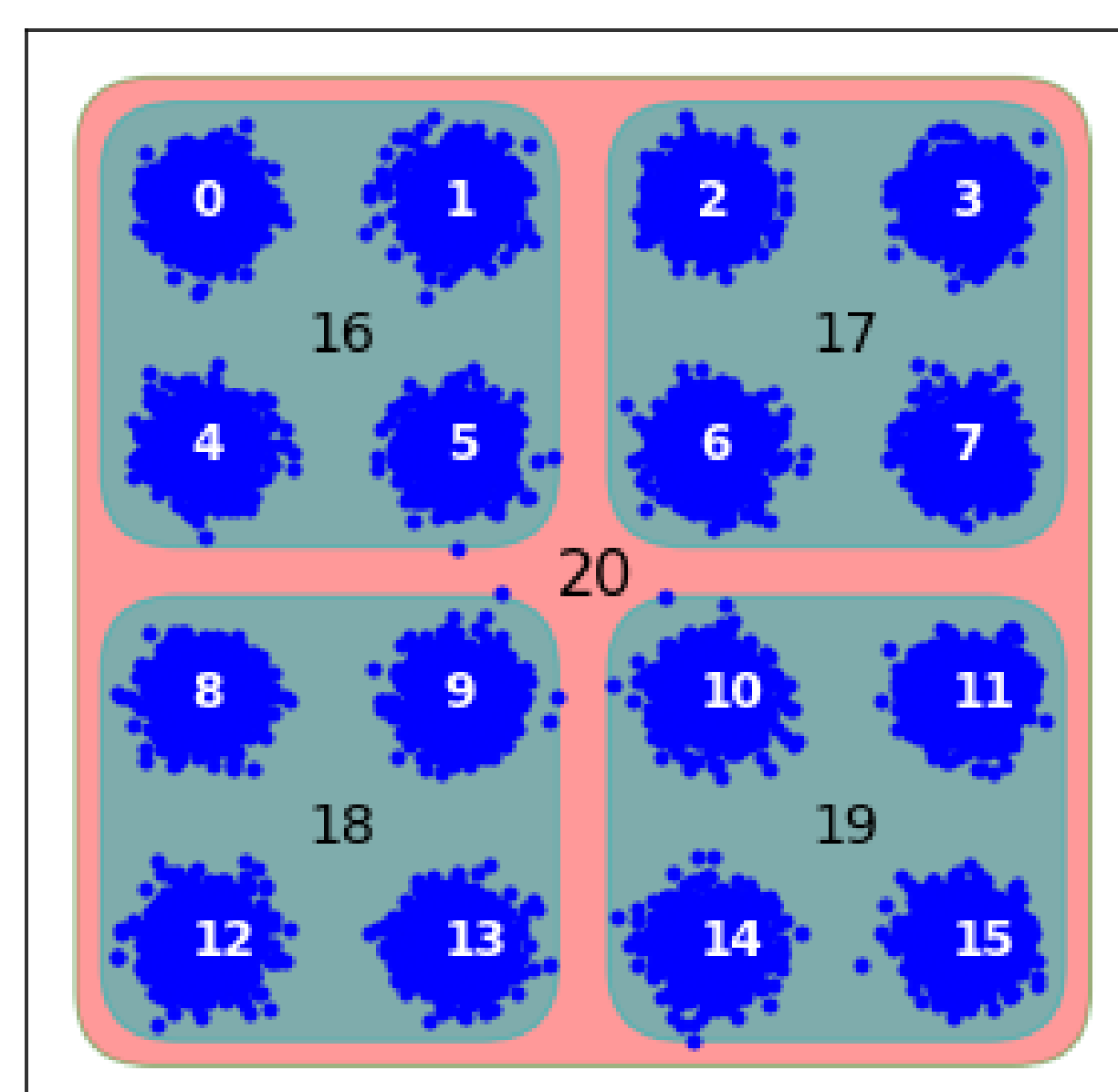


Figure: Gaussian used for the synthetic experiment

- 16 gaussians corresponds to a single label l_1, l_2, \dots, l_{16} .
- 3 layered tree hierarchy of labels

Prob	0.00		0.20		0.40	
	Micro F1	Macro F1	Micro F1	Macro F1	Micro F1	Macro F1
HIDDEN _{flt}	96.8	89.1	93.2	87.8	90.4	87.7
HIDDEN _{cas}	98.0	93.4	94.4	88.9	91.9	91.0
HIDDEN _{jnt}	98.1	94.0	94.8	91.6	92.3	91.7

Table: Synthetic data here has 12000 training and 8000 test samples.

Experiments

Dataset	Hierarchy	Hyperbolicity	L	Avg(L)	Max(L)	Train	Val	Test
RCV1	Tree	0	104	3.24	17	20833	2314	781265
NYT	Tree	1	120	6.58	24	86461	9606	9903
Yelp	DAG	1	539	4.07	32	98460	10939	46884

Table: Statistics of the datasets.

Dataset	Method	Micro-F1	Macro-F1
RCV1	TextCNN-Flat*	76.6	43.0
	HIDDEN _{flt}	77.9	44.5
	HIDDEN _{cas}	78.0	45.5
	HIDDEN _{jnt}	79.3	47.3
NYTimes	TextCNN-Flat*	69.5	39.5
	HIDDEN _{flt}	76.4	37.1
	HIDDEN _{cas}	74.6	33.2
	HIDDEN _{jnt}	77.0	43.6
Yelp	TextCNN-Flat*	62.8	27.3
	HIDDEN _{flt}	62.5	37.9
	HIDDEN _{cas}	60.5	33.9
	HIDDEN _{jnt}	60.8	35.6

Table: Performance comparison on all three datasets with TextCNN as the base classification model.

Dataset	Method	Micro-F1	Macro-F1
RCV1	HIDDEN _{auc}	78.4	47.6
	HIDDEN _{jnt}	79.3	47.3
NYTimes	HIDDEN _{auc}	76.4	40.4
	HIDDEN _{jnt}	77.0	43.6
Yelp	HIDDEN _{auc}	61.1	34.2
	HIDDEN _{jnt}	60.8	35.6

Table: Performance comparison for HIDDEN_{jnt} with HIDDEN_{auc}.

Dataset	HIDDEN _{jnt}		HiLAP	
	Micro	Macro	Micro	Macro
RCV1	79.3	47.3	78.6	50.5
NYTimes	77.0	43.6	69.9	43.2
Yelp	60.8	35.6	65.5	37.3

Table: Performance comparison of HIDDEN_{jnt} with HiLAP

	HIDDEN _{flt}	HIDDEN _{jnt}	HIDDEN _{cas}
RCV1	21.2	53.9	44.1
NYTimes	11.4	39.5	36.1
Yelp	16.3	31.9	28.8

Table: Spearman rank correlation test for the generated embeddings for all the datasets. Each method is compared against the ground truth hierarchy.

References

- [1] Yuning Mao and J et al. Tian. Hierarchical text classification with reinforced label assignment (hilap). In *Proceedings of the 2019 Conference of EMNLP-IJCNLP*, pages 445–455, 2019.

Acknowledgements

We thank Bamdev Mishra and Pratik Jawanpuria (Microsoft India, Hyderabad) for valuable discussions that gave us impetus to work towards this problem. Ayush Maheshwari is supported by a Fellowship from Ekal Foundation (www.ekal.org).