Well Begun is Half Done: An Implicitly Augmented Generative Framework with Distribution Modification for Hierarchical Text Classification

Huawen Feng, Jingsong Yan, Junlong Liu, Junhao Zheng, School of Computer Science and Engineering South China University of Technology (SCUT)

Outlines



Long-tailed Distributions in HTC





Implicit Augmentation



Experimental Results



Conclusion



Hierarchical Text Classification (HTC) is a challenging task which aims to extract the labels in a tree structure corresponding to a given text.



07 Long-tailed Distributions in HTC



Figure 1: The diagram of the imbalance of HTC. Only a few labels (head classes) frequently appear, while others (medium and tail classes) rarely appear. The global frequency of labels shows a long-tailed distribution.





• Due to long-tailed distribution, the head classes are always allocated with large weights, which distorts the label space.

Current methods for long-tailed distribution are inapplicable to multi-label classification because they have to handle each class separately.















$$\Delta h_{(q,t)} = \sum_{i \in S_k^{agn}} w_{ik} h_i$$

$$w_{ik} = Norm \left(\frac{1}{\log|r_i - r_k| + 1}\right)$$







$$\triangle \sigma_k = \sum_{i \in S_k^{co}} w_{ik} \sigma_k$$
$$L_{\sigma} = \sum_{k=1}^{N_c} \frac{1}{|S_k^{co}|} \sum_{i \in S_k^{co}} (w_{ik} - \log X_{ik})^2$$

$$\widetilde{h}_{(q,t)}^p \sim \mathcal{N}(h_{(q,t)} + \alpha \triangle h_{(q,t)}, \beta(\sigma_k + \triangle \sigma_k))$$



03 Implicit Augmentation

(1)
$$L_{aug} = -\frac{1}{P} \frac{1}{Q} \frac{1}{T} \sum_{p=1}^{P} \sum_{q=1}^{Q} \sum_{t=1}^{T} \log g_{(q,t)}^{p}$$
$$= -\frac{1}{P} \frac{1}{Q} \frac{1}{T} \sum_{p=1}^{P} \sum_{q=1}^{Q} \sum_{t=1}^{T} \log \frac{e^{w_{k} \tilde{h}_{(q,t)}^{p}}}{\sum_{j=1}^{N_{C}} e^{w_{j} \tilde{h}_{(q,t)}^{p}}}$$

$$L_{aug} = \frac{1}{Q} \frac{1}{T} \sum_{q=1}^{Q} \sum_{t=1}^{T} G$$

$$(2) \qquad G = \mathbb{E}_{\widetilde{h}_{(q,t)}^{p}} [\log \sum_{j=1}^{N_{C}} e^{U}]$$

$$U = \psi_{k}^{j} \widetilde{h}_{(q,t)}^{p}, \quad \psi_{k}^{j} = w_{j} - w_{k}$$

$$(3) \qquad G \leq \log \mathbb{E}_{\widetilde{h}_{(q,t)}^{p}} [\sum_{j=1}^{N_{C}} e^{U}] = \log \sum_{j=1}^{N_{C}} \mathbb{E}_{\widetilde{h}_{(q,t)}^{p}} [e^{U}]$$

(4)
$$U \sim \mathcal{N}(\psi_k^j(h_{(q,t)} + \alpha \bigtriangleup h_{(q,t)}), \Sigma_{jk})$$
$$\Sigma_{jk} = \beta \psi_k^j(\sigma_k + \bigtriangleup \sigma_k) \psi_k^j^T$$

(5)
$$\operatorname{E}\left[e^{tX}\right] = e^{t\mu + \frac{1}{2}\sigma^{2}t^{2}}, \quad X \sim \mathcal{N}\left(\mu, \sigma^{2}\right)$$

(6)
$$\log \sum_{j=1}^{N_C} \mathbb{E}_{\widetilde{h}_{(q,t)}^p}[e^U] = \log \sum_{j=1}^{N_C} e^V$$

(7)
$$V = \psi_k^j(h_{(q,t)} + \alpha \triangle h_{(q,t)}) + \frac{1}{2}\beta \psi_k^j(\sigma_k + \triangle \sigma_k) \psi_k^{j^T}$$

8)
$$\overline{L_{aug}} = \frac{1}{Q} \frac{1}{T} \sum_{q=1}^{Q} \sum_{t=1}^{T} \log \sum_{j=1}^{N_C} e^V = -\frac{1}{Q} \frac{1}{T} \sum_{q=1}^{Q} \sum_{t=1}^{T} e^{w_k (h_{(q,t)} + \alpha \bigtriangleup h_{(q,t)})} \log(\frac{e^{w_k (h_{(q,t)} + \alpha \bigtriangleup h_{(q,t)})}}{\sum_{j=1}^{N_C} e^{w_j (h_{(q,t)} + \alpha \bigtriangleup h_{(q,t)}) + \frac{1}{2} \beta \psi_k^j (\sigma_k + \bigtriangleup \sigma_k) \psi_k^{j^T}}})$$

LREC-COLING 2024

03 Implicit Augmentation + Random Shuffling



BFS: 1=>2=>5=>6 DFS: 1=>5=>2=>6



Implicit Augmentation + Random Shuffling







IMplicitly Augmented GenerativE framework with distribution modification (IMAGE)

Model	WOS		RCV1-V2		BGC	
	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1
TextRNN	77.94	69.65	-	-	-	-
TextCNN	82.00	76.18	76.60	43.00	-	-
TextRCNN	83.55	76.99	81.57	59.25		-
HiAGM	85.82	80.28	83.96	63.35	77.22	57.91
BERT+HiAGM	86.04	80.19	85.58	67.93	-	-
HTCInfoMax	85.58	80.05	83.51	62.71		-
BERT+HTCInfoMax	86.30	79.97	85.53	67.09	-	-
HiMatch	86.20	80.53	84.73	64.11	76.57	58.34
BERT+HiMatch	86.70	81.06	86.33	68.66	78.89	63.19
HGCLR	87.11	81.20	86.49	68.31	-	-
HPT	87.16	81.93	87.26	69.53		-
SGM-T5	85.83	80.79	84.39	65.09	77.84	60.91
Seq2Tree	87.20	82.50	87.20	70.01	79.72	63.96
PAAM-HiA-T5	90.36	81.64	87.22	70.02	-	-
IMAGE	87.76	82.38	87.69	70.77	81.16	67.37

Table 2: Experimental results of on three benchmarks. The best result is in red, the second is in green, and the third is in blue. The table is divided into three parts, and from top to bottom they are: discriminative methods, generative methods, and our method.





Model	Size
BERT+GCN	148M
BERT+Graphormer	156M
T5	220M
BART	140M

Table 3: The size of different backbones.





Model	Micro-F1	Macro-F1
IMAGE	87.69	70.77
-w/o Distribution Modification	87.43	69.26
 -w/o Implicit Augmentation 	87.06	69.22
-w/o Random Shuffling	86.67	69.35
BART (Backbone)	86.51	68.15

Table 4: The results of the ablation study on the benchmark dataset - RCV1-V2.



04 Experimental Results



The output of BART (the left graph) compared with IMAGE (the right graph). Class 0, 1 and 2 (green bars) corresponds to <s>, <pad>, and </s>, which are ignorable in the global predictions. LREC-COLING 2024

Thanks for your listening.

Huawen Feng

