

给用户分享的文本自动“贴标签”——

一个融合语义知识和阅读行为的深度学习模型

三人行语义沙龙, AI在图情—2019图书馆前沿技术论坛 (IT4L 2019),

8月13日, 上海图书馆

董行 (hang)

西交利物浦大学计算机科学系 博士生



Xi'an Jiaotong-Liverpool University

西交利物浦大学



UNIVERSITY OF
LIVERPOOL



研究目的: 自动社会化标注

- 在社会化网站中，自动给用户分享的文档标注，以方便在线文本的组织、检索和推荐。（下图: Bibsonomy 网站中的社会化论文标注）

输入: 标题
+ 摘要

The screenshot shows a Bibsonomy entry for a research paper. At the top left is a blue user icon. To its right is the title "Rules for Inducing Hierarchies from Social Tagging Data" in a dark blue box. Below the title is the author information "H. Dong, W. Wang, and F. Coenen. *Transforming Digital Worlds*, page 345–355. Cham, Springer International Publishing, (2018)". On the far right are four small icons: a pencil, an 'X', a folder, and a dropdown arrow. Below the title is a section titled "Abstract" with a grey icon. The abstract text reads: "Automatic generation of hierarchies from social tags is a challenging task. We identified three rules, set inclusion, graph centrality and information-theoretic condition from the literature and proposed two new rules, fuzzy set inclusion and probabilistic association to induce hierarchical relations. We proposed an hierarchy generation algorithm, ... (more)". To the right of the abstract is a vertical sidebar with the same title and author info, followed by a "Abstract" section with the same text. At the bottom of the sidebar is a note: "Automatic generation of hierarchies from social tags is a challenging task. We identified three rules, set inclusion, graph centrality".

输出: 标签



Images in website:
<https://www.bibsonomy.org/bibtex/2e41f059d3b72c2bbfdb0063eaf8772b0>

场景：社会化问答系统中问题的自动标注



问题(标题 & 内容) + 标签层级关系 -> 标签



克里斯蒂亚诺·罗纳尔多 (Cristiano Ronaldo)

修改

克里斯蒂亚诺·罗纳尔多 (Cristiano Ronaldo) » 话题组织 » 完整话题结构

克里斯蒂亚诺·罗纳尔多 (Cristiano Ronaldo)

2018 年俄罗斯世界杯

葡萄牙国家男子足球队

摩洛哥国家男子足球队

父级话题

父话题是一个完全包括该话题的更大的话题。

「根话题」

e. 朋友

· 体育从业者

② 运动员

• 早班

6. 前锋

• 前言

• 亮

如何评价 2010 埃多斯世界杯飞马再破门，葡萄牙 1.0 摩洛哥：

[世界杯-C罗4分钟破门葡萄牙胜 摩洛哥遭连败出局](#)

相关问题C罗622球成为进球最多葡萄牙人，是否说明了C罗的历史地位超过了尤西比奥？

Images in website:

- ```
[1] https://biendata.com/competition/zhihu/
[2] https://www.zhihu.com/question/281774641/answer/422318818
[3] https://www.zhihu.com/topic/19577949/organize/entire#anchor-children-topic
```

# 微博自动标注

自动给用户发布的微博添加标签，以方便微博的检索与管理。

CCTV5 V  
【最佳进球：C罗龙出浅海头槌破敌制胜】#2018世界杯#第7比赛日#最佳进球#，仍然来自众星之  
星，克里斯蒂亚诺-罗纳尔多！当世天骄封神将，龙出浅海战沙场。北非雄狮亡枪下，金靴金杯遥  
在望！摩洛哥求首胜心切，无奈开场就送出破绽。战术角球后莫雷拉精准制导，C罗低空腾跃头  
槌重砸皮球，攻陷球门！本届的金靴，C罗必会收入囊中！#微5世界杯# CCTV5的微...

收起全文^



今天08:00 来自 微博云剪

收藏 转发 2807 评论 214 点赞 1022

OptaJoe @OptaJoe · 14h  
85 - Cristiano Ronaldo has now scored more international goals than any other European player in the history of football (85 goals for Portugal). Historic.

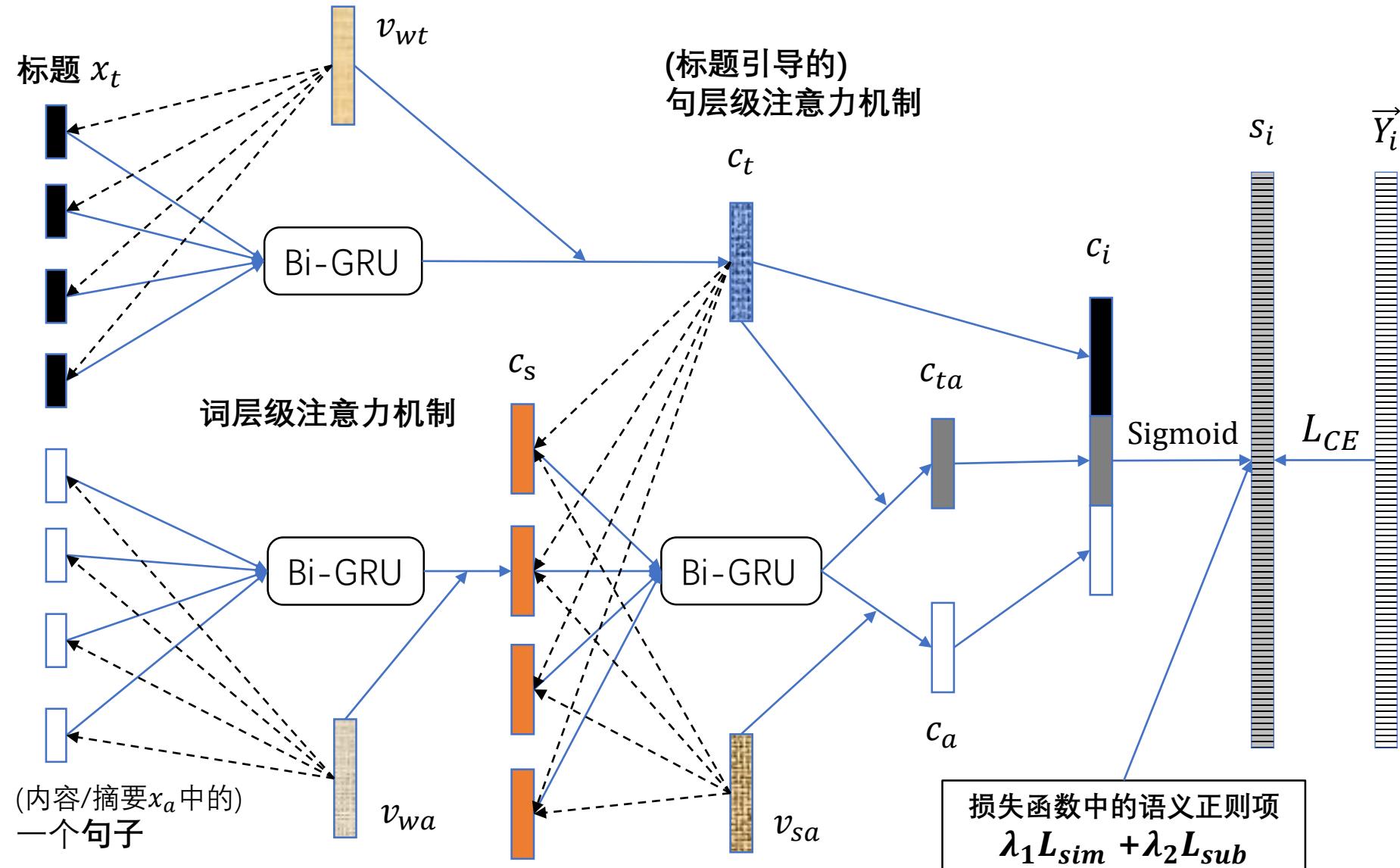
#POR 🇵🇹 #MAR 🇲🇷 #PORMAR #Ronaldo #WorldCup



91 4.8K 8.8K

# 当前的方法与问题

- 多标签分类 (Multi-label classification), 需要考虑标签间的关系。
- 方法: 注意力机制的神经网络 (Bi-GRU (Cho *et al.*, 2014), HAN (Yang *et al.*, 2016))
- 问题: (1) 在预测标签时没有融入标签关系 (知识组织系统 KOS)
  - 相似关系
  - 层级关系
- (2) 在“阅读”文档时, 忽视了元数据层次之间的关系
  - 题名对于阅读和标注的重要性: 具有很强的表达性和区分性 (Figueiredo *et al.*, 2013)
  - 题名与摘要(内容)之间的关系



Coloured version of the Figure in H. Dong, W. Wang, K. Huang, F. Coenen, Joint Multi-Label Attention Networks for Social Text Annotation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019)*, Volume 1 (Long and Short Papers), pp. 1348-1354.

# 损失函数中的语义正则项 (semantic-based loss regularisers)

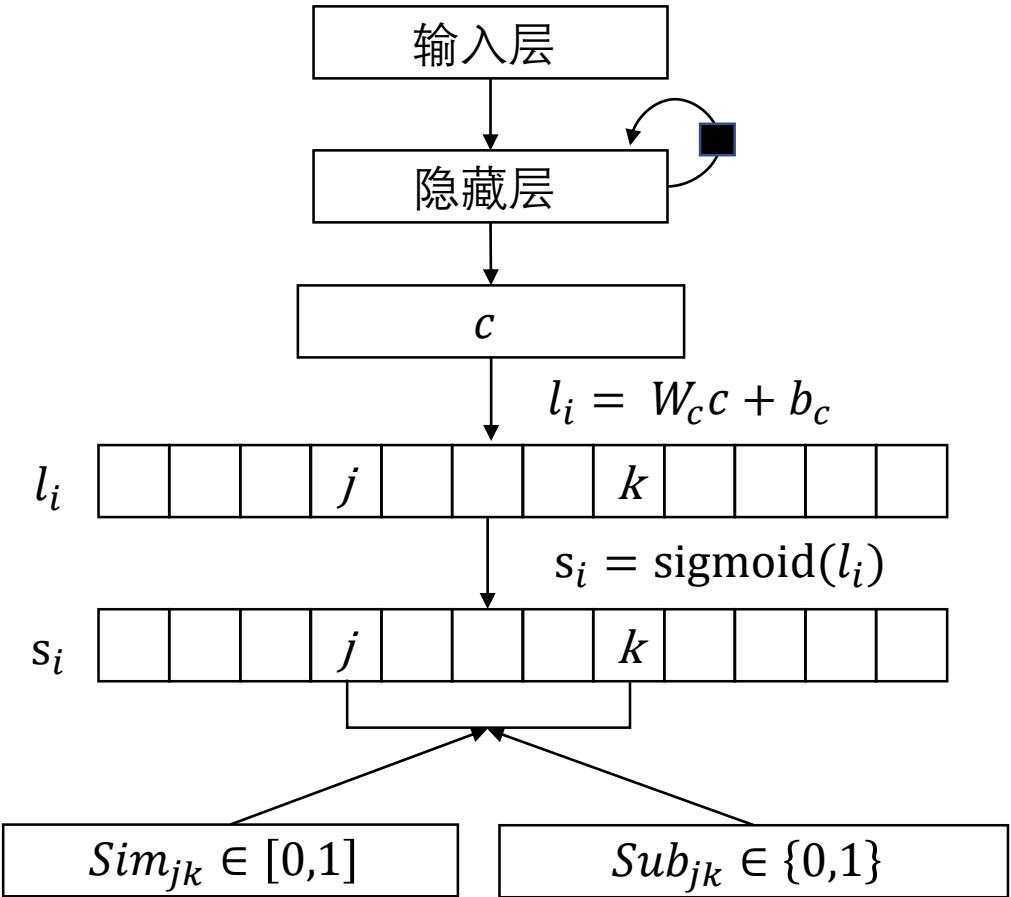
$$L = L_{CE} + \lambda_1 L_{sim} + \lambda_2 L_{sub}$$

$$L_{CE} = - \sum_i \sum_j (y_{ij} \log(s_{ij}) + (1 - y_{ij}) \log(1 - s_{ij}))$$

$$L_{sim} = \frac{1}{2} \sum_i \sum_{j,k | y_j, y_k \in Y_i} Sim_{jk} |s_{ij} - s_{ik}|^2$$

$$L_{sub} = \frac{1}{2} \sum_i \sum_{j,k | y_j, y_k \in Y_i} Sub_{jk} R(s_{ij})(1 - R(s_{ik}))$$

- $Sim_{jk} \in [0,1]$  表示 第j个标签 和 第k个标签 的相似程度.
- $Sub_{jk} \in \{0,1\}$  表示 第j个标签 是否是 第k个标签 的子标签.



# 标题引导的注意力机制

- 词层级注意力机制

$$v^{(i)} = \tanh(W_t h^{(i)} + b_t)$$

$$\alpha^{(i)} = \frac{\exp(v_{wt} \bullet v^{(i)})}{\sum_{i \in [1, n_t]} \exp(v_{wt} \bullet v^{(i)})}$$

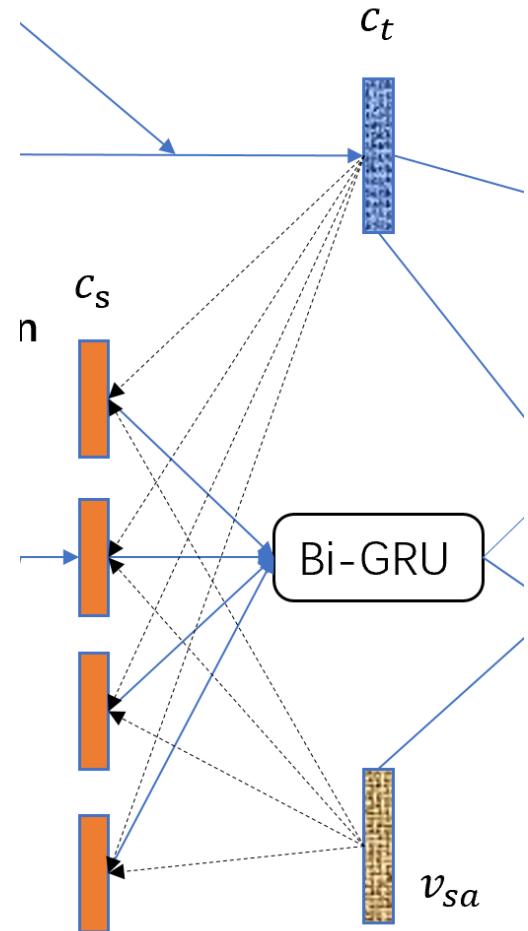
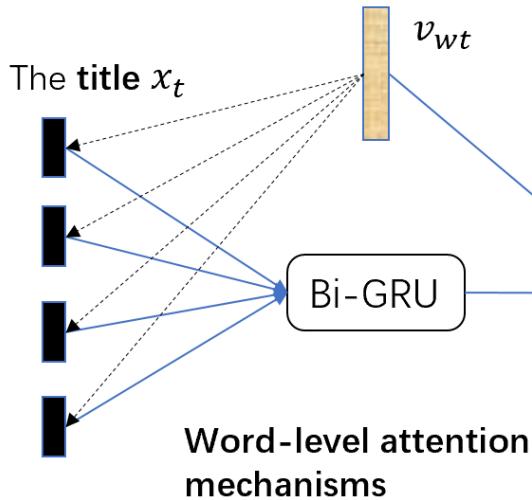
$$c_a = \sum_{i \in [1, n_t]} \alpha^{(i)} h^{(i)}$$

- 标题引导的句层级注意力机制

$$v_s^{(r)} = \tanh(W_s h_s^{(r)} + b_s)$$

$$\alpha_s^{(r)} = \frac{\exp(c_t \bullet v_s^{(r)})}{\sum_{k \in [1, n_s]} \exp(c_t \bullet v_s^{(k)})}$$

$$c_{ta} = \sum_{r \in [1, n_s]} \alpha_s^{(r)} h_s^{(r)}$$



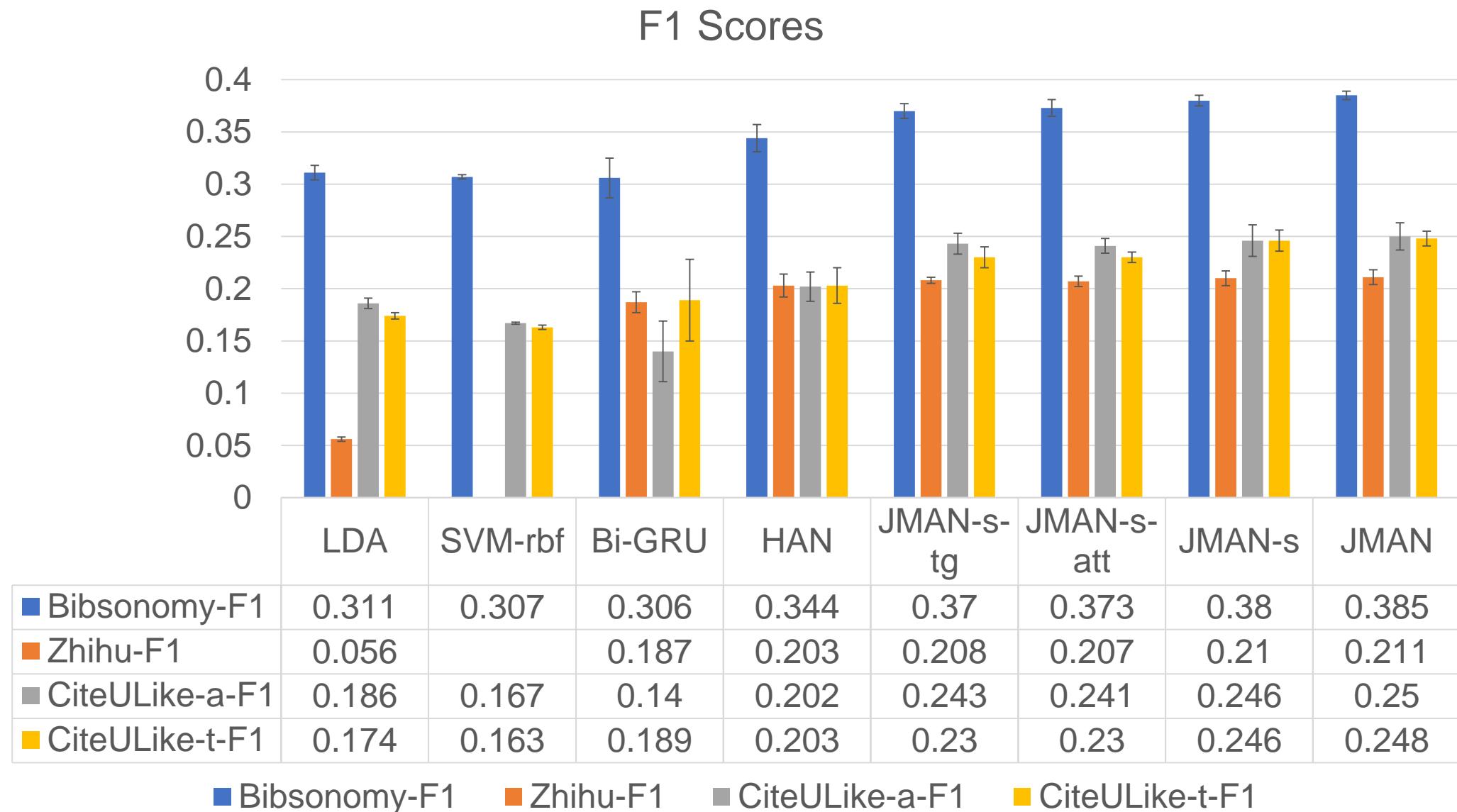
# 实验设置：数据集和实现

- 社会化论文标注 (Bibsonomy, CiteULike) 和问题标注 (知乎 Zhihu)
- 数据集统计:  $|X|$ , 文档数;  $|Y|$ , 标签数;  $|V|$ , 词汇数;  $Ave$ , 文档对应的平均标签数;  $\sum Sub$ , 标签层级关系对的数量

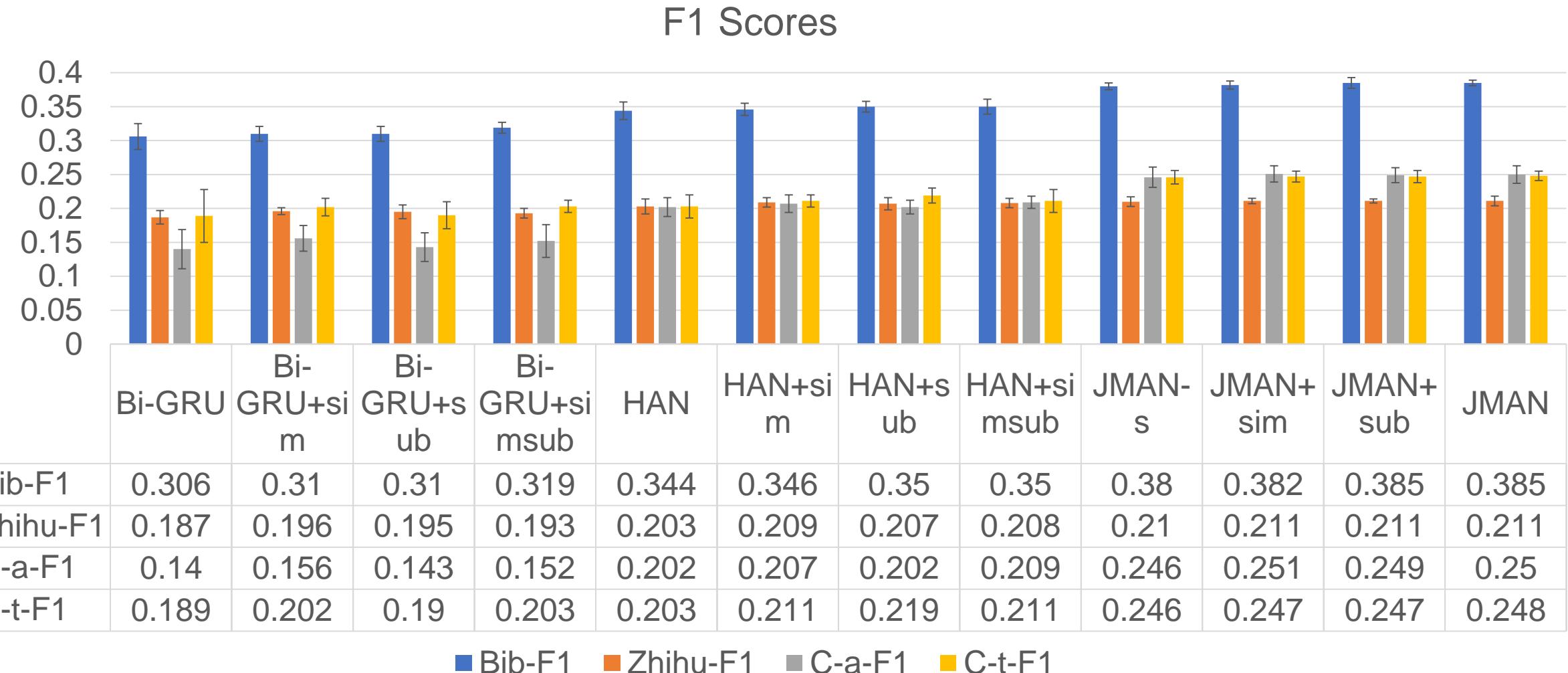
| Dataset             | $ X $   | $ Y $ | $ V $  | $Ave$ | $\Sigma Sub$ |
|---------------------|---------|-------|--------|-------|--------------|
| Bibsonomy (clean)   | 12,101  | 5,196 | 17,619 | 11.59 | 101,084      |
| CiteULike-a (clean) | 13,319  | 3,201 | 17,489 | 11.60 | 107,273      |
| CiteULike-t (clean) | 24,042  | 3,528 | 23,408 | 7.68  | 141,093      |
| Zhihu (sample)      | 108,168 | 1,999 | 62,519 | 2.45  | 2,655        |

- 20% 作为最终测试数据, 80% 进行10折交叉验证
- 标签相似矩阵  $Sim_{jk}$  通过计算标签集嵌入embedding的余弦距离获得;
- 标签层级矩阵  $Sub_{jk}$  通过匹配 Microsoft Concept Graph 获得 (知乎除外, 知乎提供用户编辑的标签层级关系)
- 参数: 词嵌入维度100; 隐藏层节点数100; 最大文本长度 对于 Bibsonomy, CiteULike 300个词 (句长30), 知乎 100个词 (句长25) ...

# 结果 (1) – 对比基准方法



# 结果 (2) – 关于语义正则项



# 结果 (3) – 模型运行和收敛速度

Comparison of training time for all models in seconds

|     | SVM           | LDA                               | Bi-GRU          | Bi-GRU+s       | HAN           | HAN+s          | JMAN-s-tg                          | JMAN-s-att                          | JMAN-s                             | JMAN           |
|-----|---------------|-----------------------------------|-----------------|----------------|---------------|----------------|------------------------------------|-------------------------------------|------------------------------------|----------------|
| Bib | $1107 \pm 12$ | <b><math>110 \pm 2(1)</math></b>  | $1480 \pm 92$   | $1683 \pm 78$  | $1164 \pm 52$ | $1434 \pm 74$  | $1075 \pm 87$                      | <b><math>1024 \pm 100(3)</math></b> | $894 \pm 55(2)$                    | $1138 \pm 86$  |
| C-a | $1660 \pm 31$ | <b><math>113 \pm 3(1)</math></b>  | $869 \pm 288$   | $877 \pm 57$   | $462 \pm 63$  | $554 \pm 45$   | $434 \pm 49$                       | <b><math>429 \pm 41(3)</math></b>   | $394 \pm 33(2)$                    | $468 \pm 38$   |
| C-t | $4796 \pm 50$ | <b><math>210 \pm 7(1)</math></b>  | $1635 \pm 1034$ | $1469 \pm 276$ | $858 \pm 100$ | $947 \pm 115$  | <b><math>752 \pm 52(3)</math></b>  | $780 \pm 69$                        | <b><math>744 \pm 62(2)</math></b>  | $839 \pm 49$   |
| Zhi | over 1 day    | <b><math>903 \pm 31(1)</math></b> | $1455 \pm 69$   | $2459 \pm 151$ | $1387 \pm 78$ | $2388 \pm 275$ | <b><math>1220 \pm 81(3)</math></b> | $1275 \pm 99$                       | <b><math>1147 \pm 44(2)</math></b> | $1712 \pm 105$ |

Training time of the three most efficient models are in **bold** and marked with a ranking index in brackets. BiGRU+s and HAN+s denote the models with semantic-based loss regularisers.

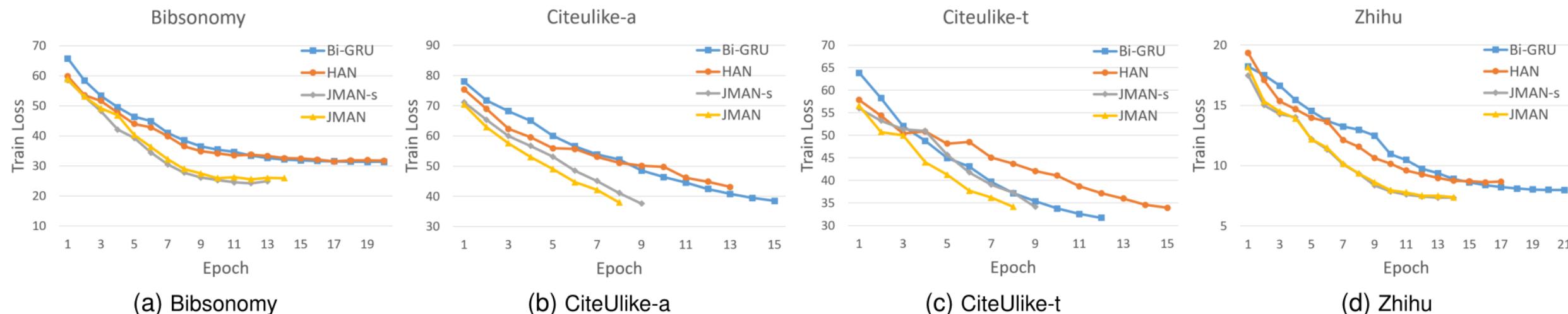


Fig. 3. Convergence plot: training loss with respect to the number of training epochs for the Bi-GRU, HAN, JMAN-s and JMAN models

# 可视化: 注意力机制和标注结果 (Bibsonomy的测试数据)

| ori                | tg | title:      |               | an            | information               | visualization | tool             | for           | personalized    |  |
|--------------------|----|-------------|---------------|---------------|---------------------------|---------------|------------------|---------------|-----------------|--|
|                    |    | knowledge   | work          | in            | many                      | fields        | requires         | examining     | several         |  |
|                    |    | attain      | meaningful    | understanding | that                      | is            | not              | explicitly    | available       |  |
|                    |    | despite     | recent        | advances      | in                        | document      | corpus           | visualization | research        |  |
|                    |    | principled  | approaches    | which         | enable                    | the           | users            | to            | personalize     |  |
|                    |    | in          | this          | paper         | ,                         | we            | present          | information   | for             |  |
|                    |    | innovative  | visualization | tool          | which                     | employs       | the              | personal      | model           |  |
|                    |    | not         | only          | does          | the                       | tool          | allow            | information   | to              |  |
|                    |    | exploration | and           | analysis      | of                        | a             | document         | users         | it              |  |
|                    |    | the         | usability     | of            | the                       | tool          | was              | ,             | the             |  |
|                    |    | it          | worthwhile    | to            | conduct                   | a             | usability        | and           |                 |  |
| <b>prediction:</b> |    | user        | information   | visualization | information_visualization |               |                  |               |                 |  |
| <b>labels:</b>     |    | user        | information   | interface     | user_interface            | semantic      | social           | management    |                 |  |
|                    |    | ontology    | visualization | personal      | information_visualization | exploratory   | semantic_desktop | desktop       | personal_inform |  |

紫色块显示句子（占标题title下的每两行）中每个词对于标注的重要性；

红色块显示“ori”（原始的）和“tg”（标题引导的）句层级注意力权重，表示文档中每个句子对于标注的重要性。颜色越深则权重越大。

文档下方展示了标注结果(prediction) 和真实的用户生成的标签(labels)。

# 可视化: 注意力机制和标注结果 (Bibsonomy的测试数据)

| ori                            | tg   | title: chinese culture and ecommerce an exploratory study                                   |                                                                                       |                                                   |                                         |                    |                    |                             |                    |                    |                    |                    |                          |
|--------------------------------|------|---------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------|---------------------------------------------------|-----------------------------------------|--------------------|--------------------|-----------------------------|--------------------|--------------------|--------------------|--------------------|--------------------------|
| 0.02                           | 0    | differing characteristics of local environments , both infrastructural and socioeconomic    | have created significant level of variation in the                                    | and the regions of the                            | and the                                 | the                | the                | the                         | the                | the                | the                | the                | socioeconomic            |
| 0.16                           | 0.01 | acceptance this paper focuses on the ecommerce development impact in china                  | and growth of ecommerce on the impact in china                                        | of different these                                | of these                                | china              | china              | china                       | china              | china              | china              | china              | china                    |
| 0.16                           | 0.11 | the ecommerce findings provide insights into the factors that may impact culture in broader | findings and development of ecommerce we present and discuss our findings             | the factors that may impact culture in broader    | that may impact culture in broader      | china              | china              | china                       | china              | china              | china              | china              | broader                  |
| 0.16                           | 0.12 | acceptance in this paper identify changes that will be required for findings broader        | and development , changes that will be required for findings broader                  | we present and be required for findings broader   | will be required for findings broader   | china              | china              | china                       | china              | china              | china              | china              | broader                  |
| 0.31                           | 0.38 | cultural issues institutional trust socializing attitudes toward debt were determined       | cultural issues institutional trust socializing attitudes toward debt were determined | socializing attitudes toward debt were determined | effect of commerce debt were determined | china              | china              | china                       | china              | china              | china              | china              | transactional determined |
| 0.16                           | 0.36 | however their means for payment research also shows that even most enlightened              | means for payment research also shows that even most enlightened                      | research also shows that even most enlightened    | shows that even most enlightened        | consumers in china | consumers in china | consumers in china          | consumers in china | consumers in china | consumers in china | consumers in china | enlightened              |
| <b>prediction:<br/>labels:</b> |      | culture study                                                                               | e_commerce culture                                                                    | china e_commerce                                  | chinese china                           | office             | commerce           | intercultural_communication | chinese            |                    |                    |                    |                          |

# 可视化: 注意力机制和标注结果 (CiteULike-t的测试数据)

| ori                    | tg              | title:                |                        | virtual       | machines                 | versatile                      | platforms       | for                  | systems       | an               |
|------------------------|-----------------|-----------------------|------------------------|---------------|--------------------------|--------------------------------|-----------------|----------------------|---------------|------------------|
| ori                    | virtual machine | machine compatibility | technology constraints | applies and   | the hardware             | concept resource               | of constraints  | virtualization       | to enable     | an a             |
|                        | virtual         | machines              | are                    | rapidly       | becoming                 | an                             | essential       | element              | in            | computer         |
|                        | they efficiency | provide               | system                 | security      | ,                        | flexibility                    | ,               | cross                | platform      | compatibility    |
|                        | designed        | to                    | solve                  | problems      | in role                  | combining                      | and             | using                | major         | computer         |
|                        | technologies    | play                  | a                      | key           | the                      | in process                     | many            | disciplines          | ,             | sy               |
|                        | for             | example               | ,                      | at            |                          | level                          | level           | ,                    | virtualizing  | including        |
|                        | platform        | independent           | network                | computing     |                          |                                |                 |                      |               | operating        |
|                        | at platform     | the and               | system in              | level servers | ,                        | they                           | support         | multiple             | operating     | system           |
|                        | br that         | br employ             | historically them      | ,             | individual               | virtual                        | machine         | techniques           | have referred | enviro           |
|                        | in a            | this unified          | text discipline        | ,             | some smith               | cases and                      | they take       | even a               | new           | been to approach |
| prediction:<br>labels: | pulling         | together              | cross cutting          | technologies  | allows                   | virtual                        | machine         | implementations      | to            |                  |
|                        | a topics        | well structured       | manner set             | emulation     |                          | dynamic system                 | program virtual | translation machines | and for       | opti b           |
| including              |                 | instruction           | java                   | and           |                          |                                |                 |                      |               |                  |
| prediction:<br>labels: |                 | machine book          | virtual machine        | comp virtual  | virtualization text_book | virtual_machine virtualization |                 |                      |               | 15               |

# 结论

- 现有的神经网络可以从融合知识和用户行为的角度提升性能。在自动社会化标注中，这主要体现在
  - 融合知识：用标签间相似和层级关系约束多标签分类的输出层。
  - 用户行为：在阅读和标注时，从题名字段出发，然后引导句子的理解，会产生更好的标注效果。
- 未来的研究：探索融合知识和用户行为的新方法，提升神经网络在相关应用上的性能与可解释性。
  - 不同的知识呈现形式
  - 用户行为的模式
  - 新的模型结构和融合方式

# 参考文献和推荐阅读

[1] Dong, H., Wang, W., Huang, K., & Coenen, F. (2019, June). Joint Multi-Label Attention Networks for Social Text Annotation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* (pp. 1348-1354). (自动社会标注的深度学习模型)

## 关于深度学习与注意力机制:

[2] Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural Machine Translation by Jointly Learning to Align and Translate. In *3rd International Conference on Learning Representations, ICLR 2015, Conference*. <http://arxiv.org/abs/1409.0473> (注意力机制最早在机器翻译中的提出和应用)

[3] Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., & Hovy, E. (2016, June). Hierarchical attention networks for document classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies* (pp. 1480-1489). (层级注意力网络 HAN 在文本分类中的提出和应用)

[4] Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014, October). Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 1724-1734). (双向门控循环单元 Bi-GRU 在机器翻译中的提出和应用)

## 关于自动社会化标注:

[4] Figueiredo, F., Pinto, H., Belém, F., Almeida, J., Gonçalves, M., Fernandes, D., & Moura, E. (2013). Assessing the quality of textual features in social media. *Information Processing & Management*, 49(1), 222-247. (不同元数据字段作为文本特征的质量研究：其中一个结论是文本的题名字段具有很强的表达性和区分性)

[5] Hassan, H. A. M., Sansonetti, G., Gasparetti, F., & Micarelli, A. (2018, September). Semantic-based tag recommendation in scientific bookmarking systems. In *Proceedings of the 12th ACM Conference on Recommender Systems* (pp. 465-469). ACM. (社会化科技文献标注)

[6] Nie, L., Zhao, Y. L., Wang, X., Shen, J., & Chua, T. S. (2014). Learning to recommend descriptive tags for questions in social forums. *ACM Transactions on Information Systems (TOIS)*, 32(1), 5. (社会化问答系统中的问题标注)

谢谢聆听

董行 | [hangdong@liverpool.ac.uk](mailto:hangdong@liverpool.ac.uk)

代码开源: <https://github.com/acadTags/Automated-Social-Annotation>

海报: <https://www.aclweb.org/anthology/attachments/N19-1136.Poster.pdf>

主页: <http://cgi.csc.liv.ac.uk/~hang/index.html>