软件学报 Software Journal of Warre



主办

🤰 中国科学院软件研究所

(中国计算机学会

出版

3 科学出版社



主 编 赵 琛,中国科学院软件研究所

执行主编 金 芝,北京大学

副 主 编 Yao, Andrew C., 清华大学

编委会成员(以汉语拼音为序)

外籍和港澳地区编委

Bjφrner, Dines, Technical University of Denmark, Denmark Chin, Francis Yuk-Lun, University of Hong Kong, China

Clarke, Edmund M., Carnegie Mellon University, USA

Graham, Ronald L., University of California at San Diego, USA

Krieg-Bruckner, Bernd, Universitat Bremen, Germany

Li, Ming, University of Waterloo, Canada

Li, Kai, Princeton University, USA

Motiwalla, Juzar, National University of Singapore, Singapore

Zhu, Hong, Oxford Brookes University, UK

资深编委

戴国忠, 中国科学院软件研究所

董韫美,中国科学院软件研究所

冯玉琳, 中国科学院软件研究所

何新贵, 北京大学

李德毅, 电子系统工程研究所

李 未,北京航空航天大学

林惠民,中国科学院软件研究所

陆汝钤,中科院数学与系统科学研究院 陆维明,中科院数学与系统科学研究院

钱华林,中国科学院计算机网络信息中心

孙家广,清华大学

王 珊,中国人民大学

吴建平,清华大学

杨芙清, 北京大学

张 钹,清华大学

赵沁平, 北京航空航天大学

郑南宁, 西安交通大学

周龙骧, 中科院数学与系统科学研究院

领域编委

陈克非, 杭州师范大学

杜小勇, 中国人民大学

冯登国, 中国科学院软件研究所

黄 涛,中国科学院软件研究所

蒋昌俊, 东华大学

李 华,中国科学院计算技术研究所

李宣东,南京大学

罗军舟,东南大学

吕 建,南京大学

梅 宏,北京理工大学

沈一栋,中国科学院软件研究所

田 捷,中国科学院自动化研究所

王 戟,国防科技大学

巫英才, 浙江大学

肖 侬,中山大学

徐 恪,清华大学

于 剑,北京交通大学

张 民, 苏州大学

郑纬民,清华大学

周傲英, 华东师范大学

周志华,南京大学

责任编委

陈文光,清华大学

陈翌佳,复旦大学

崔 斌,北京大学

丁佐华, 浙江理工大学

冯晓兵, 中国科学院计算技术研究所

冯新宇,南京大学

高 宏,哈尔滨工业大学

高 阳,南京大学

葛敬国,中国科学院信息工程研究所

金 海, 华中科技大学

阚海斌, 复旦大学

刘 璘,清华大学

刘云浩,清华大学

马晓星,南京大学

毛文吉, 中国科学院自动化研究所

毛晓光,国防科技大学

欧阳丹彤, 吉林大学

彭 鑫,复旦大学

任丰原,清华大学

申恒涛, 电子科技大学

舒继武,清华大学

苏金树, 国防科技大学

苏开乐, 暨南大学

眭跃飞, 中国科学院计算技术研究所

孙晓明, 中国科学院计算技术研究所

田 聪,西安电子科技大学

汪 芸,东南大学

王国仁, 北京理工大学

王建勇,清华大学

王文成,中国科学院软件研究所

王兴伟, 东北大学

魏 峻,中国科学院软件研究所

翁 健,暨南大学

徐 杨,电子科技大学

薛 锐,中国科学院信息工程研究所

杨 珉,复旦大学

尹一通,南京大学

雍俊海,清华大学

于 戈, 东北大学

曾庆凯,南京大学

詹乃军, 中国科学院软件研究所

张军平,复旦大学

张 康,香港科技大学(广州)

张 路,北京大学

张敏灵,东南大学

张自力, 西南大学

周国栋, 苏州大学

软 件 学 报

(Ruanjian Xuebao)

第 33 卷第 3 期 2022 年 3 月

目 次

数据库系统新型技术专题

数据库系统新型技术专题前言		李国良	于 戈	杨俊	范 举	(771)
	陈小强	周士俊	卞福升	吴 非	陈兵	(774)
基于树型门控循环单元的基数和代价估计器		元昌安	Louis A	lborto GLII	riedde7	(707)
AlphaQO: 鲁棒的学习型查询优化器 ···············	七				HERREZ	(191)
	柴成亮	张辛宁 刘睿诚	汤 南 张俊晨	孙 信 罗永平	李国良 金培权	
基于 NVM 和 HTM 的低时延事务处理	魏星达	陆放明	旅 核	タポー 陈海波	並培权 臧斌宇	
内存数据库并发控制算法的实验研究			卢卫	木冶知	山 1. 高	(967)
	及依冶	杨皖晴	P L	李海翔	杜小勇	(867)
屠要峰	陈河堆	王涵毅	闫宗帅	孔鲁	陈 兵	(891)
数据库管理系统中数据异常体系化定义与分类	李晓燕	刘畅	杜小勇	卢卫	潘安群	(909)
	V- h =17		عاد مال دا	76 N. 1h		
	郑志明	童咏昕	张瑞升	魏淑越	李卫华	(931)
	李瑞远	郭阳	蒋忠元	鲍捷	郑宇	
时间序列对称模式挖掘	岳晓飞	史 岚	李盼盼 赵宇海	宋韶旭 季航旭	王建民王国仁	
新型分布式计算系统中的异构任务调度框架						
		赵恒泰	金福生	李荣华	王国仁	(1005)
	崔鹏杰	袁 野	李岑浩	张 灿	王国仁	(1018)
面向大规模二部图的分布式 Tip 分解算法 周 旭		杨志邦	李博仁	张 吉	李肯立	(1043)
联邦学习中的隐私保护技术		刘艺璇	陈 红	刘宇涵	李翠平	
混洗差分隐私下的多维类别数据的收集与分析 	 王志刚	谷 峪	魏志强	 张啸剑	 于 戈	(1093)
面向多方安全的数据联邦系统				7K /// 21		
• • • • • • • • • • • • • • • • • • • •	史鼎元	廖旺冬	张利鹏	童咏昕	许 可	(1111)
模式识别与人工智能						
融合信息增益比和遗传算法的混合式特征选择算基于多覆盖模型的神经机器翻译		许召召 黄锴宇	申德荣 李玖一	聂铁铮 宋鼎新	寇 月 黄德根	
坐 1 夕 後 皿 佚 至 的 秤 左 机 顧 附 年	八	夹珀丁	丁八	/_ \text{\tau} \tau \tau \	虫1心化	(1141)
《软件学报》投稿指南						(封三)

期刊基本参数: CN11-2560/TP*1990*m*16*384*zh+en*P*\footnote{70.2022*21*2022-03

771	Preface
	LI Guo-Liang, YU Ge, YANG Jun, FAN Ju
774	Geno: Cost-based Heterogeneous Fusion Query Optimizer
	TU Yao-Feng, CHEN Xiao-Qiang, ZHOU Shi-Jun, BIAN Fu-Sheng, WU Fei, CHEN Bing
797	Cardinality and Cost Estimator Based on Tree Gated Recurrent Unit
	QIAO Shao-Jie, YANG Guo-Ping, HAN Nan, QU Lu-Lu, CHEN Hao, MAO Rui, YUAN Chang-An,
	Louis Alberto GUTIERREZ
814	AlphaQO: Robust Learned Query Optimizer
	YU Xiang, CHAI Cheng-Liang, ZHANG Xin-Ning, TANG Nan, SUN Ji, LI Guo-Liang
832	Heterogeneous Index for Non-volatile Memory
	LIU Rui-Cheng, ZHANG Jun-Chen, LUO Yong-Ping, JIN Pei-Quan
849	Low-latency Transaction Processing Using NVM and HTM
	WEI Xing-Da, LU Fang-Ming, CHEN Rong, CHEN Hai-Bo, ZANG Bin-Yu
867	Experimental Study on Concurrency Control Algorithms in In-Memory Databases
	ZHAO Hong-Yao, ZHAO Zhan-Hao, YANG Wan-Qing, LU Wei, LI Hai-Xiang, DU Xiao-Yong
891	Optimal Design of NUMA-aware Persistent Memory Storage Engine
	TU Yao-Feng, CHEN He-Dui, WANG Han-Yi, YAN Zong-Shuai, KONG Lu, CHEN Bing
909	Systematic Definition and Classification of Data Anomalies in Data Base Management Systems
	LI Hai-Xiang, LI Xiao-Yan, LIU Chang, DU Xiao-Yong, LU Wei, PAN An-Qun
931	Intelligent System for Distributed Social Governance Based on Big Data
	LÜ Wei-Feng, ZHENG Zhi-Ming, TONG Yong-Xin, ZHANG Rui-Sheng, WEI Shu-Yue, LI Wei-Hua
950	Distributed Time Series Similarity Search Method Based on Key-value Data Stores
	YU Zi-Sheng, LI Rui-Yuan, GUO Yang, JIANG Zhong-Yuan, BAO Jie, ZHENG Yu
968	Time Series Symmetric Pattern Mining
	LI Pan-Pan, SONG Shao-Xu, WANG Jian-Min
985	Dynamic Resource Allocation Strategy for Flink Iterative Jobs
	YUE Xiao-Fei, SHI Lan, ZHAO Yu-Hai, JI Hang-Xu, WANG Guo-Ren
1005	Heterogeneous Task Scheduling Framework in Emerging Distributed Computing Systems
	LIU Rui-Qi, LI Bo-Yang, GAO Yu-Jin, LI Chang-Sheng, ZHAO Heng-Tai, JIN Fu-Sheng, LI Rong-Hua,
	WANG Guo-Ren
1018	RGraph: Effective Distributed Graph Data Processing System Based on RDMA
	CUI Peng-Jie, YUAN Ye, LI Cen-Hao, ZHANG Can, WANG Guo-Ren
1043	Distributed Algorithm for Tip Decomposition on Large Bipartite Graphs
	ZHOU Xu, WENG Tong-Feng, YANG Zhi-Bang, LI Bo-Ren, ZHANG Ji, LI Ken-Li
1057	Privacy-preserving Techniques in Federated Learning
	LIU Yi-Xuan, CHEN Hong, LIU Yu-Han, LI Cui-Ping
1093	Collecting and Analyzing Multidimensional Categorical Data Under Shuffled Differential Privacy
	LIU Yi-Fei, WANG Ning, WANG Zhi-Gang, GU Yu, WEI Zhi-Qiang, ZHANG Xiao-Jian, YU Ge
1111	Data Federation System for Multi-party Security
	LI Shu-Yuan, JI Yu-Dian, SHI Ding-Yuan, LIAO Wang-Dong, ZHANG Li-Peng, TONG Yong-Xin, XU Ke
TERN	RECOGNITION AND ARTIFICIAL INTELLIGENCE

PAT

- 1128 Hybrid Feature Selection Algorithm Combining Information Gain Ratio and Genetic Algorithm XU Zhao-Zhao, SHEN De-Rong, NIE Tie-Zheng, KOU Yue
- 1141 Multi-coverage Model for Neural Machine Translation LIU Jun-Peng, HUANG Kai-Yu, LI Jiu-Yi, SONG Ding-Xin, HUANG De-Gen

©Copyright 2022, Institute of Software, the Chinese Academy of Sciences. All rights reserved. No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form without the prior written permission of the Institute of Software, the Chinese Academy of Sciences.

基于多覆盖模型的神经机器翻译

刘俊鹏, 黄锴宇, 李玖一, 宋鼎新, 黄德根

(大连理工大学 计算机科学与技术学院, 辽宁 大连 116024)

通信作者: 黄德根, E-mail: huangdg@dlut.edu.cn



E-mail: jos@iscas.ac.cn

http://www.jos.org.cn

Tel: +86-10-62562563

摘 要:覆盖模型可以缓解神经机器翻译中的过度翻译和漏翻译问题.现有方法通常依靠覆盖向量或覆盖分数等单一方式存储覆盖信息,而未考虑不同覆盖信息之间的关联性,因此对信息的利用并不完善.针对该问题,基于翻译历史信息的一致性和模型之间的互补性,提出了多覆盖融合模型.首先定义词级覆盖分数概念;然后利用覆盖向量和覆盖分数存储的信息同时指导注意力机制,降低信息存储损失对注意力权重计算的影响.根据两种覆盖信息融合方式的不同,提出了两种多覆盖融合方法.利用序列到序列模型在中英翻译任务上进行了实验,结果表明,所提方法能够显著提升翻译性能,并改善源语言和目标语言的对齐质量.与只使用覆盖向量的模型相比,过度翻译和漏翻译问题的数量得到进一步减少.

关键词: 神经机器翻译; 注意力机制; 序列到序列模型; 多覆盖模型; 过度翻译; 漏翻译中图法分类号: TP183

中文引用格式: 刘俊鹏, 黄锴宇, 李玖一, 宋鼎新, 黄德根. 基于多覆盖模型的神经机器翻译. 软件学报, 2022, 33(3): 1141-1152. http://www.jos.org.cn/1000-9825/6201.htm

英文引用格式: Liu JP, Huang KY, Li JY, Song DX, Huang DG. Multi-coverage Model for Neural Machine Translation. Ruan Jian Xue Bao/Journal of Software, 2022, 33(3): 1141–1152 (in Chinese). http://www.jos.org.cn/1000-9825/6201.htm

Multi-coverage Model for Neural Machine Translation

LIU Jun-Peng, HUANG Kai-Yu, LI Jiu-Yi, SONG Ding-Xin, HUANG De-Gen

(School of Computer Science and Technology, Dalian University of Technology, Dalian 116024, China)

Abstract: The over-translation and under-translation problem in neural machine translation could be alleviated by coverage model. The existing methods usually store the coverage information in a single way, such as coverage vector or coverage score, but do not take the relationship among different coverage methods into consideration, leading to insufficient use of the information. This study proposes a multi-coverage mechanism based on the consistency of translation history and complementarity between different models. The concept of word-level coverage score is defined first, and then the coverage information stored in both coverage vector and coverage score are incorporated into the attention mechanism simultaneously, aiming to reduce the influence of information loss. According to different fusion methods, two models are introduced. Experiments are carried out on Chinese-to-English translation task based on sequence-to-sequence model. Results show that the proposed method could significantly enhance translation performance and improve alignment quality between source and target words. Compared with the model with coverage vector only, the number of over-translation and under-translation problem is further reduced.

Key words: neural machine translation; attention mechanism; sequence-to-sequence model; multi-coverage model; over-translation; under-translation

机器翻译是使用计算机将一种语言(源语言)自动翻译成另一种语言(目标语言)的技术. 近年来, 神经机器翻译(neural machine translation, NMT)经过不断发展和完善, 显示出强大的翻译性能. 现有的神经机器翻译模型通常采用"编码器-解码器"架构, 其中, 编码器和解码器可以采用不同的网络结构, 依据神经网络拓扑结构

* 基金项目: 国家重点研发计划(2020AAA0108004); 国家自然科学基金(61672127, U1936109)

收稿时间: 2020-04-01; 修改时间: 2020-06-17, 2020-09-30; 采用时间: 2020-10-28

的特点, 分为循环神经网络(recurrent neural network)^[1-3]、卷积神经网络(convolutional neural network)^[4]和自注 意力网络(self-attention network)^[5]等. 目前, 神经机器翻译的性能已经显著超越传统的统计机器翻译, 成为当 前机器翻译领域的前沿热点^[6].

虽然大规模训练数据和强大的计算能力使得神经机器翻译的性能大幅提升,但是神经机器翻译还存在着过度翻译(over-translation)和漏翻译(under-translation)问题.在基于短语的统计机器翻译模型中,文献[7]通过将加入输出序列的翻译结果所对应的源语言短语标记为"已翻译"的方式,确保所有的源语言短语均被翻译覆盖且只被翻译一次.但是在神经机器翻译模型中没有显式存储历史翻译信息的结构,并且每一次解码过程都需要所有源语言词语的参与,因此极易出现过度翻译和漏翻译问题.

为了向神经机器翻译模型中增加覆盖信息,文献[8]借鉴统计机器翻译中的覆盖思想提出了覆盖模型 (coverage model). 设计一个覆盖向量(coverage vector, CV)记录解码过程中的翻译历史信息,并指导注意力权重计算,降低已翻译词语的权重,使注意力机制更多地关注未翻译词语. 类似地,文献[9]为源语言词语设置全覆盖编码向量(full coverage embedding vector),在解码过程中不断削减已翻译词语的编码向量,以此降低其在未来解码中的作用. 文献[10]利用循环注意力机制(recurrent attention mechanism)给注意力单元提供更多的重调序信息,并通过条件解码器(conditioned decoder)减少重复翻译. 文献[11]增加了两个额外的循环神经网络分别记录翻译过程中的历史(past)和未来(future)信息,并利用该信息指导注意力机制和解码状态.

覆盖度还可以与解码算法相结合,用于在已生成的译文中筛选对原文忠实度最高的结果.文献[12]通过覆盖惩罚(coverage penalty)和长度归一化(length normalization)改进束搜索(beam search)算法,使模型在选择翻译结果时考虑该翻译结果对源语言信息的覆盖程度,避免偏向句子长度更短的翻译结果.文献[13]提出将覆盖分数(coverage score, CS)引入每一次束搜索过程,以此减少搜索错误,并且使覆盖分数的计算适用于源语言词汇和目标语言词汇的多种映射关系.

在翻译过程中,源语言上下文通常影响翻译的忠实度,而目标语言上下文则与译文的流利度有关.基于这种思路,文献[14]提出了上下文门(context gate)方法,根据源语言和目标语言上下文的重要程度,动态控制两种上下文信息在生成目标词语时的影响比重.该方法可以与覆盖机制相结合,能够在改善翻译结果对源语言的覆盖度的同时,提升译文的流利度.文献[15]通过解码历史增强注意力机制(decoding-history enhanced attention mechanism)建立所有源语言词和目标语言词的结构关系,使神经机器翻译模型更好地选择源语言端和目标语言端的信息.

此外,文献[16]针对源语言中被漏翻译的词语特点展开研究,发现在源语言中,翻译熵(translation entropy)越大的词语越容易被漏翻译,并据此设计了一种粗粒度到细粒度(coarse-to-fine)框架,分别解决句子级和词级的训练和翻译问题,从而减少熵高词语的漏翻译情况.

虽然上述方法均能在一定程度上缓解神经机器翻译中的过度翻译和漏翻译问题,但该问题依然不能被完全避免.如表 1 所示,在基线系统的翻译结果中,源语言句子中的"事关全局(related to the overall situation)"被遗漏翻译,且"体制(system)"被重复翻译一次.而引入覆盖模型后,上述问题并未得到明显改正,"事关全局"虽然被未被遗漏但却被错误地翻译为"in the world",而"体制"则被漏翻译.

文本类别	例句	
源语言句子	深化国家监察体制改革是 事关全局 的重大政治 体制 改革.	
参考译文	Deepening the reform of the state oversight system is a major political structural reform related to the overall situation.	
基线系统	Deepening the reform of the national oversight system is a major political system system reform.	
覆盖模型	模型 Deepening the reform of the state supervisory system is a major political reform in the world .	
多覆盖融合模型	Deepening the reform of the country's supervision system is a major political system reform to the overall situation.	

表 1 过度翻译和漏翻译示例

注:基线系统采用基于循环神经网络的 Seq2Seq 模型,覆盖模型是根据文献[8]在基线系统上的复现

由此可见, 现有的覆盖模型对覆盖信息的记录和使用并不完善. 可能的原因是, 覆盖向量在使用 GRU 网

络更新时存在信息损失,导致注意力权重分配不准确,进而产生了重复翻译和漏翻译现象.针对上述问题,本 文提出一种多覆盖融合的神经机器翻译模型,首先定义一种词级覆盖分数,用于记录源语言词语在解码过程 中的注意力累积情况;而后在解码阶段,利用覆盖向量和覆盖分数所记录的覆盖信息同时指导注意力权重计 算,使两种覆盖信息互相补充,从而最大限度减少信息损失.

本文的主要创新点包括两个方面: (1) 定义了词级覆盖分数概念, 并利用覆盖分数指导注意力机制的权重 计算、拓展了覆盖分数在神经机器翻译模型中的使用方式: (2) 提出了多种覆盖信息的融合模型及两种实现方 式,并利用实验验证了方法的可行性和有效性. 在 CWMT2018 中英数据集上的实验结果表明, 多覆盖融合方 法能显著提升翻译质量,并在覆盖模型基础上,进一步减少过度翻译和漏翻译现象.

1 研究背景

1.1 基于循环神经网络的神经机器翻译模型

基于注意力机制的神经机器翻译模型通常采用"编码器-解码器"结构, 编码器将源语言句子编码成隐层向 量表示, 解码器根据编码器的输出逐字预测目标端句子的单词序列. 给定源端输入句子 X={x1.x2,....xm}, 神经 机器翻译模型对目标端句子 $Y=\{y_1,y_2,...,y_n\}$ 的条件概率 P(Y|X)进行建模.

$$P(Y \mid X) = \prod_{j=1}^{n} P(y_j \mid y_{< j}, X)$$
 (1)

具体来说, 若当前已生成的目标端输出为 $\{y_1,y_2,...,y_{i-1}\}$, 则生成下一个目标词语 y_i 的概率计算公式为

$$P(y_i|y_{< i},X) = g(y_{i-1},t_i,s_i)$$

$$\tag{2}$$

其中, $g(\cdot)$ 是非线性函数, t_i 是 j 时刻解码器的隐层状态, s_i 是将编码器的所有隐层状态加权得到的源语言上下文 向量. t_i 和 s_i 的计算公式如下.

$$t_{j} = f(y_{j-1}, t_{j-1}, s_{j})$$
 (3)

$$s_j = \sum_{i=1}^m a_{ij} \cdot \boldsymbol{h}_i \tag{4}$$

$$s_{j} = \sum_{i=1}^{m} a_{ij} \cdot \mathbf{h}_{i}$$

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{m} \exp(e_{kj})}$$

$$(5)$$

$$e_{ij} = ATT(\boldsymbol{t}_{j-1}, \boldsymbol{h}_i) = v_a^T \tanh(W_a \boldsymbol{t}_{j-1} + U_a \boldsymbol{h}_i)$$
(6)

其中,

- h_i 是源语言词语 x_i 经编码器编码后得到的隐层向量,由于在实际中常使用双向循环神经网络作为编 码器,因此 $h_i = [\vec{h}_i; \vec{h}_i]$,其中, \vec{h}_i 和 \vec{h}_i 分别是前向和后向循环神经网络的隐层状态,[:::]表示拼接操作;
- *a*_{ii}为注意力权重,表示目标端词汇 *y*_i与源语言词语 *x*_i之间的关联度;
- *ATT*(·)是计算注意力权重的匹配函数;
- v_a^T , W_a 和 U_a 为权重矩阵.

由于注意力机制的引入、源语言上下文向量 s_i 不再是一个单一固定的向量、而是随着不同时刻 a_{ii} 的不同 而变化、从而使解码器在不同时刻对源语言句子中各个词语的信息考量有所侧重. f(·)通常使用长短期记忆网 络(long-short term memory, LSTM)或门控循环单元(gated recurrent unit, GRU), 其网络结构可采用单层或多层, 并且多层网络结构的性能比单层网络有更加显著的提升[12]. 多层堆叠长短期记忆网络(stacked long-short term memory)中,第 i 层和 i+1 层状态的计算公式如下所示。

$$c_{t}^{i}, h_{t}^{i} = LSTM_{i}(c_{t-1}^{i}, h_{t-1}^{i}, x_{t}^{i-1})$$

$$(7)$$

$$x_t^i = h_t^i \tag{8}$$

$$x_{i}^{i} = h_{i}^{i}$$

$$c_{t}^{i+1}, h_{i}^{i+1} = LSTM_{i+1}(c_{t-1}^{i+1}, h_{t-1}^{i+1}, x_{t}^{i})$$

$$(8)$$

其中, $LSTM_i$ 和 $LSTM_{i+1}$ 分别表示第 i 和 i+1 层 $LSTM_i$ x_i^i, c_i^i 和 h_i^i 分别表示 t 时刻 $LSTM_i$ 的输入、记忆单元状态 和隐层状态.

最后,在训练集 $\{(x^p, y^p)\}_{p=1}^J$ 上,利用极大似然估计对目标函数的参数 θ 进行迭代训练,目标函数如公式 (10)所示.

$$L(\theta) = -\frac{1}{J} \sum_{t=1}^{J} \sum_{t=1}^{n} \log P(y_t^p \mid y_{< t}^p, x^p)$$
 (10)

1.2 覆盖向量

与统计机器翻译模型中双语词语的硬对齐(hard-alignment)关系不同,神经机器翻译模型的注意力机制提供的是一种软对齐(soft-alignment)方法,因此难以对覆盖机制进行建模.文献[8]提出通过覆盖向量显式存储解码过程中历史信息的覆盖模型,利用覆盖向量所存储的信息指导注意力评分过程,从而使注意力机制更多地关注未翻译的词语,其模型结构如图 1 所示.

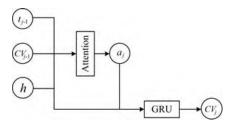


图 1 基于覆盖的注意力模型结构

加入覆盖向量指导后, 注意力权重的计算方法由公式(6)修改为公式(11).

$$e_{ij} = ATT(\mathbf{t}_{j-1}, \mathbf{h}_i, CV_{i,j-1}) = v_a^T \tanh(W_a \mathbf{t}_{j-1} + U_a \mathbf{h}_i + V_a CV_{i,j-1})$$
(11)

其中, V_a 为权重矩阵, $CV_{i,i-1}$ 表示 j 时刻前源语言词语 x_i 对应的覆盖向量.

在每一步解码完成后,利用 GRU 网络对每个源语言词语覆盖向量中存储的覆盖信息进行更新,更新方法如公式(12)所示.

$$CV_{i,j} = GRU(CV_{i,j-1}, a_{ij}, \boldsymbol{h}_{i}, \boldsymbol{t}_{j-1})$$

$$(12)$$

1.3 覆盖分数

文献[13]提出了表示源语言句子翻译程度的覆盖分数概念. 给定一个翻译句对(X,Y), 将源语言词语 x_i 的覆盖度定义为所有目标词语对该源语言词语的注意力权重之和, 如公式(13)所示.

$$coverage_{x_i} = \sum_{j=1}^{|Y|} a_{ij}$$
 (13)

在此基础上,利用所有源语言词语的覆盖度表示源语言句子的覆盖分数,计算公式如公式(14)所示.

$$cs(X,Y) = \sum_{i=1}^{|X|} \log \max(coverage_{x_i}, \beta)$$
 (14)

其中, β 为可调参数. 最后,将模型预测的条件概率与译文的覆盖分数线性组合,使模型在选择译文时兼顾对源语言句子翻译覆盖程度. 改进后的评价函数如公式(15)所示.

$$score(X,Y) = a \cdot \log P(Y|X) + b \cdot cs(X,Y)$$
 (15)

其中, logP(Y|X)表示模型预测的条件概率值, a 和 b 是用于平衡条件概率和覆盖分数作用的参数.

2 多覆盖融合模型

虽然覆盖向量和覆盖分数均能显式记录翻译过程中的覆盖信息,但二者在信息的存储方式和使用方式上都有所不同:前者以向量的形式对信息进行存储和更新,并通过指导注意力权重计算的方式传递翻译历史;

而后者以常量的形式进行累加,并作为评价指标用于翻译结果的选择.两种方法各有优缺点:覆盖向量存储的信息抽象程度更高,但在利用 GRU 网络进行更新时,由于重置门(reset gate)会自动丢弃一定比例的历史信息,并且更新门(update gate)也会丢掉一部分新的覆盖信息,因此可能造成覆盖信息损失;而覆盖分数虽然表达直观,但难以确定取值界限,因而无法直接根据数值大小衡量不同词语的覆盖程度.

由于在任意时刻覆盖向量和覆盖分数中存储的信息具有一致性,且由上述分析可知二者具备一定的互补性,因此提出一种将两种方法优点相结合的多覆盖融合模型.利用覆盖向量和覆盖分数存储的信息同时指导注意力机制,从而降低信息损失对注意力权重分配的影响.为了将覆盖分数引入注意力机制,首先定义了词级覆盖分数概念;然后,根据覆盖向量和覆盖分数融合方式的不同提出了层次型和平行型两种多覆盖融合模型.总体框架如图 2 所示.

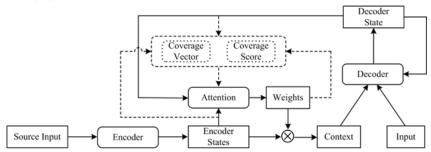


图 2 多覆盖融合的神经机器翻译模型结构

2.1 词级覆盖分数

词级覆盖分数用于表示任意时刻译文结果对每个源语言词语的覆盖程度. 给定源语言句子 X 和 j 时刻目标端的预测序列 $\{y_1,y_2,...,y_j\}$,那么 j 时刻源语言 x_i 的覆盖分数定义为当前已生成的所有目标端词语 $\{y_1,y_2,...,y_j\}$ 对源语言 x_i 的注意力权重之和的截断取值,如公式(16)所示.

$$CS_{ij} = \max\left(\sum_{k=1}^{j} a_{ik}, 1\right) \tag{16}$$

这里, 截断函数 max(·)的作用主要有以下两个方面.

(1) 改善源语言词语覆盖程度的可比较性

在翻译过程中,存在某些源语言词语对应多个目标词语的情况.因此,这些源语言的注意力分数累加和存在大于1的可能,这使得不同源语言词语之间的覆盖程度难以比较.经过截断取值之后,覆盖分数 CS_{ij} 的取值范围设定在[0,1]之间,那么对于所有覆盖分数取值为1源语言词语,认为该词已被翻译覆盖;而覆盖分数小于1的词语,在后续解码过程中,仍有继续被翻译覆盖的可能.

(2) 借鉴硬对齐思想并融入注意力机制

将源语言词语的覆盖分数取值控制在[0,1]的区间内, 更容易直观地表示当前翻译结果对源语言的覆盖情况和未来翻译过程中仍需被注意力机制关注的, 因此便于融入注意力机制并指导注意力权重的计算.

2.2 层次多覆盖模型

层次多覆盖模型(hierarchical multi-coverage, HMC)主要基于覆盖向量和覆盖分数存储信息的一致性原理,模型结构如图 3 所示,首先,利用覆盖向量指导的注意力模型计算当前时刻分配给所有源语言词语的注意力权重;然后,利用覆盖分数存储的信息对注意力权重进行校验和重新分配,从而减少由于覆盖信息丢失而造成注意力权重分配错误的现象.层次多覆盖模型注意力权重的计算方法如公式(17)—公式(19)所示.

$$\tilde{e}_{ij} = ATT(\boldsymbol{t}_{j-1}, \boldsymbol{h}_{i}, CV_{i,j-1}) \tag{17}$$

$$\tilde{a}_{ij} = \frac{\exp(\tilde{e}_{ij})}{\sum_{k=1}^{m} \exp(\tilde{e}_{kj})}$$
(18)

2022 年第	第 3 期	数据库系统新型技术	李国良,于戈,杨俊,范举
2022 年第	第4期	面向开放场景的鲁棒机器学习研究	陈恩红,李宇峰,邹权
2022 年第	第 5 期	领域软件工程	汤恩义,江贺,陈俊洁,李必信,唐滨
2022 年第	第6期	系统软件安全	杨珉,张超,宋富,张源
2022 年第	第6期	定理证明理论与应用	曹钦翔,詹博华,赵永望
2022 年第	第7期	智能系统的分析和验证	明仲,张立军,秦胜潮
2022 年第	第8期	形式化方法与应用	陈立前,孙猛
2022 年第	第9期	融合媒体环境下的媒体内容分析与信息服务技术	汪萌,张勇东,俞俊,张伟
2022 年第	等10期	智慧信息系统新技术	邢春晓,王鑫,张勇,于戈

登录软件学报网站: http://www.jos.org.cn 免费下载专刊/专题全文.

软 件 学 报

Ruanjian Xuebao

(月刊, 1990年创刊)

第33卷第3期 2022年3月

Journal of Software

(monthly)

(Started in 1990)

Vol.33 No.3 Mar. 2022

主管单位 中国科学院

主办单位 中国科学院软件研究所 中国计算机学会

主 编 赵琛

编 辑 《软件学报》编辑部

(北京 8718 信箱 邮编 100190)

电话: 010-62562563, E-mail: jos@iscas.ac.cn

http://www.jos.org.cn

编辑部主任 方 梅

出 版

(北京东黄城根北街 16号 邮编 100717) EП 刷 北京宝昌彩色印刷有限公司 总发行处 中国邮政集团公司北京市报刊发行局

订购处 全国各地邮局

国外总发行 中国国际图书贸易总公司 (北京 399 信箱 邮编 100044)

Sponsored by the Chinese Academy of Sciences

Published by Institute of Software, The Chinese Academy of

Sciences (ISCAS) and China Computer Federation

Editor-in-Chief: ZHAO Chen

Edited by Editorial Board of Journal of Software (P.O.Box 8718, Beijing 100190, P.R.China) Tel: 8610-62562563, E-mail: jos@iscas.ac.cn http://www.jos.org.cn

Distributed by Science Press (16 Donghuangchenggen North Street, Beijing 100717, China)

Printed by Beijing Baochang Color Printing Co., Ltd

Generally Distributed by Beijing Bureau for Distribution of Newspapers and Journals

Domestically Distributed by All Local Post Offices in China Overseas Distributed by China International Book Trading Corporation (P.O.Box 399, Beijing 100044, China)

ISSN 1000-9825

国内邮发代号: 82-367 CN 11-2560/TP

国外发行代号: M4628

©2022 ISCAS (版权所有)

定价: 70.00 元

