

传送门

Image Caption 任务综述 | PaperWeekly

2017-01-31 科研圈

本文经授权转载自公众号 Paper Weekly（ID：paperweekly）。Paperweekly 每周分享自然语言处理领域好玩的paper。

引言

Image Caption是一个融合计算机视觉、自然语言处理和机器学习的综合问题，它类似于翻译一副图片为一段描述文字。该任务对于人类来说非常容易，但是对于机器却非常具有挑战性，它不仅需要利用模型去理解图片的内容并且还需要用自然语言去表达它们之间的关系。除此之外，模型还需要能够抓住图像的语义信息，并且生成人类可读的句子。

随着机器翻译和大数据的兴起，出现了Image Caption的研究浪潮。当前大多数的Image Caption方法基于encoder-decoder模型。其中encoder一般为卷积神经网络，利用最后全连接层或者卷积层的特征作为图像的特征，decoder一般为递归神经网络，主要用于图像描述的生成。由于普通RNN存在梯度下降的问题，RNN只能记忆之前有限的时间单元的内容，而LSTM是一种特殊的RNN架构，能够解决梯度消失等问题，并且其具有长期记忆，所以一般在decoder阶段采用LSTM。

问题描述

Image Caption问题可以定义为二元组(I,S)的形式，其中I表示图，S为目标单词序列，其中S={S1,S2,...}，其中St为来自于数据集提取的单词。训练的目标是使最大似然 $p(S|I)$ 取得最大值，即使生成的语句和目标语句更加匹配，也可以表达为用尽可能准确的用语句去描述图像。

数据集

论文中常用数据集为Flickr8k,Flick30k,MSCOCO,其中各个数据集的图片数量如下表所示。



分享这篇文章

数据

阅读 1009

点赞 2

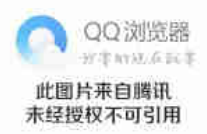
更新 2月2日 0:47

科研圈 最新头条文章

她35岁进大学，39岁读博，破解性别之谜，却因性别被遗忘 2017-07-20

成功才是成功之母，“胜利者效应”的神经机制首次被发现 2017-07-19

掌握这些技巧，生活科研两不误 2017-07-18



数据集图片和描述示例如图

其中每张图像都至少有5张参考描述。为了使每张图像具有多种互相独立的描述，数据集使用了不同的语法去描述同一张图像。如示例图所示，相同图像的不同描述侧重场景的不同方面或者使用不同的语法构成。

模型

本文主要介绍基于神经网络的方法

1 NIC[1]

Show and Tell: A Neural Image Caption Generator

本文提出了一种encoder-decoder框架，其中通过CNN提取图像特征，然后经过LSTM生成目标语言，其目标函数为最大化目标描述的最大似然估计。



该模型主要包括encoder-decoder两个部分。encoder部分为一个用于提取图像特征的卷积神经网络，可以采用VGG16，VGG19，GoogleNet等模型，decoder为经典的LSTM递归神经网络，其中第一步的输入为经过卷积神经网络提取的图像特征，其后时刻输入为每个单词的词向量表达。对于每个单词首先通过one-hot向量进行表示，然后经过词嵌入模型，变成与图像特征相同的维度。

2 MS Captivator[2]

From captions to visual concepts and back

本文首先利用多实例学习，去训练视觉检测器来提取一副图像中所包含的单词，然后学习一个统计模型用于生成描述。对于视觉检测器部分，由于数据集对图像并没有准确的边框标注，并且一些形容词、动词也不能通过图像直接表达，所以本文采用Multiple Instance Learning(MIL)的弱监督方法，用于训练检测器。



科研圈 热门头条文章

- 2016年中国高被引学者榜单发布 10万+阅读
- 《自然》重磅：肺竟然是造血器官？ 98436阅读
- “玩具”登上 Nature 子刊：纸片离心机，只要一块钱！ 84693阅读
- 这位博士的五年总结，或许能解答你许多问题 69032阅读
- 大面积重复失败，韩春雨的基因编辑技术恐遭弃用 43032阅读
- Cell 重磅 | 长生不老药从未如此接近 41701阅读
- 这位导师给博士生的五年建议，或许能帮你避免不少弯路 34076阅读
- 一个博后的自述：5年零3个月后，我选择离开学术圈 25723阅读
- 张锋登上《时代周刊》封面：编辑基因组的下一代领袖 16980阅读
- 水蒸发，能发电？Nature 子刊报道中国科学家的突破 13806阅读



3 Hard-Attention Soft-Attention[3]

Show, atten and tell: Neural image caption generation with visual attention

受最近注意机制在机器翻译中发展的启发，作者提出了在图像的卷积特征中结合空间注意机制的方法，然后将上下文信息输入到encoder-decoder框架中。在encoder阶段，与之前直接通过全连接层提取特征不同，作者使用较低层的卷积层作为图像特征，其中卷积层保留了图像空间信息，然后结合注意机制，能够动态的选择图像的空间特征用于decoder阶段。在decoder阶段，输入增加了图像上下文向量，该向量是当前时刻图像的显著区域的特征表达。



4 gLSTM[4]

Guiding long-short term memory for image caption generation

使用语义信息来指导LSTM在各个时刻生成描述。由于经典的NIC[1]模型，只是在LSTM模型开始时候输入图像，但是LSTM随着时间的增长，会慢慢缺少图像特征的指导，所以本文采取了三种不同的语义信息，用于指导每个时刻单词的生成，其中guidance分别为Retrieval-based guidance (ret-gLSTM), Semantic embedding guidance(emb-gLSTM) ,Image as guidance (img-gLSTM)。



5 sentence-condition[5]

Image Caption Generation with Text-Conditional Semantic Attention





该模型首先利用卷积神经网络提取图像特征，然后结合图像特征和词嵌入的文本特征作为gLSTM的输入。由于之前gLSTM的guidance都采用了时间不变的信息，忽略了不同时刻guidance信息的不同，而作者采用了text-conditional的方法，并且和图像特征相结合，最终能够根据图像的特定部分用于当前单词的生成。

6 Att-CNN+LSTM [6]

What value do explicit high level concepts have in vision to language problems?

如图，作者首先利用VggNet模型在ImageNet数据库进行预训练，然后进行多标签数训练。给一张图片，首先产生多个候选区域，将更多候选区域输入CNN产生多标签预测结果，然后将结果经过max pooling作为图像的高层语义信息，最后输入到LSTM用于描述的生成。该方法相当于保留了图像的高层语义信息，不仅在Image Caption上取得了不错的结果，在VQA问题上，也取得很好的成绩。



未经授权不可引用



外贸网站建设



深度强化学习

7 MSM[7]

BOOSTING IMAGE CAPTIONING WITH ATTRIBUTES



该文研究了图像属性特征对于描述结果的影响，其中图像属性特征通过多实例学习[2]的方法进行提取。作者采用了五种不同的组合形式进行对比。其中第3种、第5种，在五种中的表现出了比较好的效果。由于提取属性的模型，之前用于描述图像的单词的生成，所以属性特征能够更加抓住图像的重要特征。而该文中的第3种形式，相当于在NIC模型的基础上，在之前加上了属性作为LSTM的初始输入，增强了模型对于图像属性的理解。第5种，在每个时间节点将属性和文本信息进行结合作为输入，使每一步单词的生成都能够利用图像属性的信息。

8 When to Look[8]

Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning



该文主要提出了何时利用何种特征的概念。由于有些描述单词可能并不直接和图像相关，而是可以从当前生成的描述中推测出来，所以当前单词的生成可能依赖图像，也可能依赖于语言模型。基于以上思想，作者提出了“视觉哨兵”的概念，能够以自适应的方法决定当前生成单词，是利用图像特征还是文本特征。

结果

本文列出的模型的在COCO测试集上的结果如下：





以下为online MSCOCO testing server的结果：



总结

最近的Image Caption的方法，大多基于encoder-decoder框架，而且随着flickr30,mscoco等大型数据集的出现，为基于深度学习的方法提供了数据的支撑，并且为论文实验结果的比较提供了统一的标准。模型利用之前在机器翻译等任务中流行的Attention方法，来加强对图像有效区域的利用，使在decoder阶段，能够更有效地利用图像特定区域的特征[3]。模型利用图像的语义信息在decoder阶段指导单词序列的生成，避免了之前只在decoder开始阶段利用图像信息，从而导致了图像信息随着时间的增长逐渐丢失的问题[4][5]。模型为了更好的得到图像的高层语义信息，对原有的卷积神经网络进行改进，包括利用多分类和多实例学习的方法，更好的提取图像的高层语义信息，加强encoder阶段图像特征的提取[6][7]。随着增强学习，GAN等模型已经在文本生成等任务中取得了不错的效果，相信也能为Image Caption效果带来提升。

参考文献

1. Vinyals O, Toshev A, Bengio S, et al. Show and tell: A neural image caption generator[J]. Computer Science, 2015:3156-3164.
2. Fang H, Gupta S, Iandola F, et al. From captions to visual concepts and back[C]// IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2015:1473-1482.
3. Xu K, Ba J, Kiros R, et al. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention[J]. Computer Science, 2016:2048-2057.
4. Jia X, Gavves E, Fernando B, et al. Guiding Long-Short Term Memory for Image Caption Generation[J]. 2015.
5. Zhou L, Xu C, Koch P, et al. Image Caption Generation with Text-Conditional Semantic Attention[J]. 2016.
6. Wu Q, Shen C, Liu L, et al. What Value Do Explicit High Level Concepts Have in Vision to Language Problems?[J]. Computer Science, 2016.
7. Yao T, Pan Y, Li Y, et al. Boosting Image Captioning with Attributes[J]. 2016.

8.Lu J, Xiong C, Parikh D, et al. Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning[J]. 2016.

作者

朱欣鑫，北京邮电大学在读博士，研究方向为视觉语义理解

邮箱：zhuxinxin@bupt.edu.cn

阅读更多

▽ 故事

- 过年讲究这么多，竟因汉语擅长玩谐音？
- 摩尔定律的终结与量子计算机的诞生
- 21部科学纪录片 伴你过一个不一样的春节假期
- 跨年脑洞：鸡不会飞，为什么能遍布全球？

▽ 论文推荐

- NgAgo 或只能切割 RNA | BioRxiv 论文推荐
- 大脑也负责抗击细菌？ | Cell Immunity 论文推荐
- 噬菌体之间的信息传递和群体决策 | Nature 论文推荐

▽ 论文导读

- 【新年快乐】Science 等一周论文导读（上）| 2017年第3期
- 【来点干货】Science 等一周论文导读（下）| 2017年第3期

内容合作请联系

keyanquan@huanqiukexue.com

这里是“科学美国人”中文网《环球科学》服务科研人的微信号“科研圈”。我们：

查看标识获取更多信息

公众号推广

外贸网站建设

· 关注科学进展与科研生态

Hadoop是什么

· 推荐重要前沿研究

废钢破碎机

· 发布科研招聘

· 推送学术讲座与会议预告。

欢迎[长按二维码](#)关注。





觉得不错，分享给更多人看到 [微信](#) [QQ空间](#) [新浪微博](#) [腾讯微博](#) [人人](#) [Twitter](#) [豆瓣](#) [百度贴吧](#)

科研圈 最新文章:

- 2017-07-20 她35岁进大学，39岁读博，破解性别之谜，却因性别被遗忘
- 2017-07-20 追忆首位菲尔兹奖女性得主：坚持不懈的曲面探险家【附视频】
- 2017-07-20 复旦大学高分子科学系优秀教师招聘启事 | 长期有效
- 2017-07-20 复旦大学附属华山医院感染科招聘博士后3-4名
- 2017-07-20 来聊 | 科研人的待遇究竟怎么样？
- 2017-07-19 成功才是成功之母，“胜利者效应”的神经机制首次被发现
- 2017-07-19 酵母进化帮助人们理解癌症机制 | PNAS 论文推荐
- 2017-07-19 北京理工大学物理学院凝聚态方向招生和招聘计划 | 长期有效
- 2017-07-19 来聊 | 科研人的待遇究竟怎么样？
- 2017-07-18 掌握这些技巧，生活科研两不误

怎样锻炼阴茎	美国移民	女用充气娃娃	
	处女膜多少钱		
为什么性欲低下	不够硬进不去	老年斑怎么去掉	性生活持久方法
老年斑如何消除			
怎么让阴茎勃起	打胎需要多少钱	阴茎多大合适	
	移民美国		



科研圈 热门文章:

2016年中国高被引学者榜单发布 阅读/点赞 : 10万+/161

《自然》重磅：肺竟然是造血器官？ 阅读/点赞 : 98436/299

"玩具"登上 Nature 子刊：纸片离心机，只要一块钱！ 阅读/点赞 : 84693/218

这位博士的五年总结，或许能解答你许多问题 阅读/点赞 : 69032/132

大面积重复失败，韩春雨的基因编辑技术恐遭弃用 阅读/点赞 : 43032/144

Cell 重磅 | 长生不老药从未如此接近 阅读/点赞 : 41701/126

这位导师给博士生的五年建议，或许能帮你避免不少弯路 阅读/点赞 : 34076/124

一个博后的自述：5年零3个月后，我选择离开学术圈 阅读/点赞 : 25723/121

张锋登上《时代周刊》封面：编辑基因组的下一代领袖 阅读/点赞 : 16980/116

水蒸发，能发电？Nature 子刊报道中国科学家的突破 阅读/点赞 : 13806/119

