

软件定义世界(SDX)(/author/index?authorId=626)

团队号 (/channel/in/dex?type=263&name=team)

11月21日 1新/0 阅读: 23971 ℃: 0

简介:❤

移动开发 (/channel/hdex?type=272&name=mobile)

# 用TensorFlow为图片添加字幕

Java (/channel/index?type=257&name=java)

. 2017-05-04 06:51 软件定义世界(SDX)

■ 160 阅读4

Python (/channel/index?type=268&name=python)

个人号 (/channel/index?type=2下载name=按照即可阅读)

前端 (/channel/index?type=264%) p2045中国数据分析师行业峰会精彩PPT下载(共计21个文件)(/html

?a /626/201509/216409056/4.html) 后端 (/channel/index?type=269&name=server)

ut 测试 (/channel/index?tymndefferenters)来构建和训练一个图片字幕生成器

ho 更多 (/channel/index<u>**图/der寫生3成</u>模型結准还**近年来计算机视觉和机器翻译方面的进步,通过使用神</u>

- 经网络来生成现实图片的字幕。对于一个给定的输入图片,神经图像字幕模型被 训练来最大化生成一个字幕的可能性。可以被用来产生新颖的图像描述。例如,
- 下面是用<u>MS COCO数据集</u>训练的一个神经图像字幕生成器所产生的字幕。

62

6)

The man in grey swings a bat while the man in black looks on.



A big bus sitting next to a person.

COReillyData

图1. 来源:Paul Puri。图片来自MS COCO数据集

在这篇文章里,我们会介绍一个中级程度的教程,教大家如何使用谷歌的"Show and Tell"模型的一种变形和Flickr30k数据集来训练一个图片字幕生成器。我们使 用<u>TensorFlow</u>的框架来构建、训练和测试我们的模型,因为它相对容易使用而 且也有一个日益庞大的在线社区。

下载《开发者大全》



为什么要生成字幕?

团队号 (/channel/index?type=263&name=team)

移动开发 (/channel/index.type=272&name=motifie) 理任务上应用深度神经网络的成功激励着AI研究人 员去探索新的研究机会,交叉连接这些之前互相独立的领域。字幕生成模型就必须去 Java (/channel/index 对视觉 经素和自然语声的 理解进行平衡。

Python (/channel/index. 两个传统的表面的知识的 Application (/channel/index. 两个传统的表面的知识的 Application (/channel/index. 西方特别的一种特别的 Application (/channel/index. 西方特别的一种特别的 Application (/channel/index. 西方特别的 Application (/channel/index. 西方特别 Application (/channel/index. 西方特别的 Application (/channel/index. 西方特别 Appli 有一些很直接的应用。比如,为YouTube视频自动生成摘要或是标注未标记的图片。 个人号(/channel/index主要的子创造分的应用阶级火幅度提高一个更广泛的人群的生活质量。与传统的计 前端 (/channel/index ftype=264& name=web) 让这个世界对人类更加可达与可理解的潜力。它可以是一个导游,甚至可以成为日常 后端 (/channel/index建分的26分例说解助服务er比如意大利的AI公司Evra所开发的Horus可穿戴设备所展 示的这个场景。 测试 (/channel/index?type=266&name=test)

在我们正式开始前,需要先做一些整理工作。

首先,你需要安装TensorFlow。如果这是你第一次使用TensorFlow,我们推荐你先 看看这篇文章《你好, TensorFlow!从零开始构建和训练你的第一个TensorFlow图》

你需要安装pandas、OpenCV2和Jupyter库来保证相关的代码可以运行。不过为了简 化安装的过程,我们强烈推荐你使用与本文关联的GitHub库里的这个Docker安装指 南。

你还需要下载Flickr30k图片文件和图片字幕数据集。我们的GitHub库里有也提供了 下载链接。

现在,让我们开始吧!

图片字幕生成模型

图2. 来源:Shannon Shih获取自加州大学伯克利分校机器学习组织。马的图片来 自MS COCO

概括来看,这就是我们将要训练的模型。每张图片都可以被一个深度卷积神经网络编 码成一个4096维的向量表示。一个语言生成RNN(循环神经网络)将会对这个表示 按顺序解码,并生成自然语言的描述。

下载《开发者大全》



字幕生成是图像分类的一种扩展

团队号 (/channel/index?type=263&name=team)

移动开发 (/channel/index:type=272%name=mobile) 分类有着很长的历史,且有非常多好的模型。分类需要模型能把图像里与形状和物体相关的视觉信息拼接在一起,然后把这个图片分入 Java (/channel/index:type类别。78如他的计算机视觉的机器学习模型,诸如物体检查和图片分割等,则不仅 对呈现的信息进行识别,还通过学习如何解读二维的空间,并把这两种理解融合起来 Python (/channel/index:type=268&name=python) ,从而能判断出图片里分布的物体信息。对于字幕生成,有两个主要的问题:

个人号 (/channel/index?type=2718年里的重要管息时,我们如何基于图像分类模型的成功结果?

前端 (/channel/index?typ 發化的模型如何能学) 习去融合对于语言的理解和对于图片的理解?

后端 (/channel/index?type=269&name=server) 利用迁移学习

测试 (/channel/index?type=266&name=test)

我们可以利用已有的模型来帮助实现图片生成字幕。<u>迁移学习</u>让我们可以把从其他任更多 (/channel/index?type=270&name=other) 务上训练出来的神经网络的数据变换应用到我们自己的数据上。在我们的这个场景里 ,<u>VGG-16</u>图片分类模型是把224 x224像素的图片作为输入,产生4096维度的特征向量表示,用于分类图片。

我们可以利用这些VGG-16模型生成的表示(也叫做图向量)来训练我们的模型。限于本文的篇幅,我们省略了VGG-16的架构,而是直接用已经计算出来的4096维的特征来加快训练的过程。

导入VGG图片特征和图片字幕是相当得简单直接:

def get\_data(annotation\_path, feature\_path):

annotations = pd.read\_table(annotation\_path, sep='\t', header=None, names=[
'image', 'caption'])

return np.load(feature\_path,'r'), annotations['caption'].values

#### 理解字幕

现在我们已经有了图片的表示了,还需要我们的模型能学会把这些表示解码成能被理解的字幕。因为文字天生的序列特性,我们会利用一个RNN/LSTM网络里的循环特点(想了解更多,请参考这篇"理解LSTM网络")。这些网络被训练来预测在给定的一系列前置词汇和图片表示的情况下的下一个词。

长短期记忆(LSTM)单元能让这个模型能更好地选择什么样的信息可以用于字幕词汇序列,什么需要记忆,什么信息需要忘掉。TensorFlow提供了一个封装的功能可以下载。从去表面输入和输出组成,从downdrage/dev.apk) ★

最新 (/) 为了把词汇能变化成适合LSTM的固定长度表示的输入,我们使用一个向量层来把词汇映射成256维的特征(也叫词向量)。词向量帮助我们把词汇表示成向量,同时相团队号 (/channel/index?type=263&name=team)似的词向量在语义上也相似。想了解更多关于词向量如何获取不同词汇之间的关系,移动开发 (/channel/index.

前端 (/channel/index 物理和過程學內模型web)

后端 (/channel/index?type=269&name=server) 全合在一起 , Show and Tell模型就大概像这个样子:

测试 (/channel/index?type=266&name=test)

更多 (/channel/ind图3typ来源70%mamendtffer)1:2015 MSCOCO图片字幕大赛所获得的经验教训》

图3中, $\{s0, s1, ..., sN\}$ 表示我们试着去预测的字幕词汇, $\{wes0, wes1, ..., wesN-1\}$ 是每个词的词向量。LSTM的输出 $\{p1, p2, ..., pN\}$ 是这个模型产生的句子里下一个词的概率分布。模型的训练目标是最小化对所有的词概率取对数后求和的负值。

def build\_model(self):

# declaring the placeholders for our extracted image feature vectors, our caption, and our mask

# (describes how long our caption is with an array of 0/1 values of length `maxlen .

img = tf.placeholder(tf.float32, [self.batch\_size, self.dim\_in])

caption\_placeholder = tf.placeholder(tf.int32, [self.batch\_size, self.n\_lstm\_steps])

mask = tf.placeholder(tf.float32, [self.batch\_size, self.n\_lstm\_steps])

# getting an initial LSTM embedding from our image\_imbedding

image\_embedding = tf.matmul(img, self.img\_embedding) + self.img\_embedding\_bi
as

# setting initial state of our LSTM

下载a《干发者大全》。statese(f.download(dev-apk)at32)



```
最新 (/)
                      total_loss = 0.0
团队号 (/channel/index块块 基础 2008 cms med (tell )):
移动开发 (/channel/index?type=272&name=mobile)
Java (/channel/index?type=257&name=java)
                      #if this isn't the first iteration of our LSTM we need to get the word embedding co
Python (/channel/index?type=1268&name=python)
个人号 (/channel/ind#x2tybe+27)tl&wondenpensom)tion
前端 (/channel/index ** tybe = 2648 f a file = 60)
current_embedding = tf.nn.embedding_lookup(self.word_embedding, caption_plac
后端 (/channel/index?type=269&name=server)
eholder[:,i-1]) + self.embedding_bias
测试 (/channel/index?type=266&name=test)
  更多 (/channel/index?type=?704psqme=?theb) our LSTM we utilize the embedded image as our inp
                      ut
                      current_embedding = image_embedding
                      if i > 0:
                      # allows us to reuse the LSTM tensor variable on each iteration
                      tf.get_variable_scope().reuse_variables()
                      out, state = self.lstm(current_embedding, state)
                      print (out, self.word encoding, self.word encoding bias)
                      if i > 0:
                      #get the one-hot representation of the next word in our caption
                      labels = tf.expand_dims(caption_placeholder[:, i], 1)
                      ix_range=tf.range(0, self.batch_size, 1)
                      ixs = tf.expand_dims(ix_range, 1)
                      concat = tf.concat([ixs, labels],1)
                      onehot = tf.sparse_to_dense(
                      concat, tf.stack([self.batch_size, self.n_words]), 1.0, 0.0)
                   下载《开发者大全》
                                             下载 (/download/dev.apk)
                                            laccification to generate the next word in the caption
```

第5页 共11页

```
最新 (/)
                  logit = tf.matmul(out, self.word_encoding) + self.word_encoding_bias
移动开发 (/channel/index?type=xentropy * mask[; i]
loss = tf.reduce_sum(xentropy)
Java (/channel/index?type=257&name=java)
                  total loss += loss
Python (/channel/index?type=268&name=python)
                  total loss = total loss / tf.reduce sum(mask[:,1:])
前端 (/channel/index?type=264&name=web)
使用推断来生成字幕
后端 (/channel/index?type=269&name=server)
测试 (/channel/index完成训练后&r我们就获得了一个在给定图片和所有的前置词汇的前提下,可以给出字
                  幕里下一个词概率的模型。那么我们怎么用这个模型来生成字幕?
  更多 (/channel/index?type=270&name=other)
                  最简单的方法就是把一张图片作为输入,循环输出下一个概率最大的词,由此生成一
                  个字幕。
                  def build generator(self, maxlen, batchsize=1):
                  #same setup as `build_model` function
                  img = tf.placeholder(tf.float32, [self.batch_size, self.dim_in])
                  image_embedding = tf.matmul(img, self.img_embedding) + self.img_embedding_bi
                  as
                  state = self.lstm.zero state(batchsize,dtype=tf.float32)
                  #declare list to hold the words of our generated captions
                  all_words = []
                  print (state,image embedding,img)
                  with tf.variable scope("RNN"):
                  # in the first iteration we have no previous word, so we directly pass in the image
                  embedding
                  # and set the `previous_word` to the embedding of the start token ([0]) for the fut
                  ure iterations
                  output, state = self.lstm(image_embedding, state)
                下载el/所发者cce/ltf.nn.em时dc/lines/lanks/lines/conference/le/mbedding, [0]) + self.embed 🗶
```

 最新 (/) for i in range(maxlen):

移动开发 (/channel/index?type=2/2&name=mobile)

# get a one-hot word encoding from the output of the LSTM Java (/channel/index?type=257&name=java)

logit = tf.matmul(out, self.word encoding) + self.word encoding bias

Python (/channel/index?type=268&name=python) best\_word = tf.argmax(logit, 1)

个人号 (/channel/index?type=271&name=person) with tf.device("/cpu:0")!

后端 (/channel/index?type=269&name=server)

previous\_word = tf.nn.embedding\_lookup(self.word\_embedding, best\_word)

测试 (/channel/index?type=266&name=test)

previous word += self.embedding bias

更多 (/channel/index?type=270&name=other) all\_words.append(best\_word)

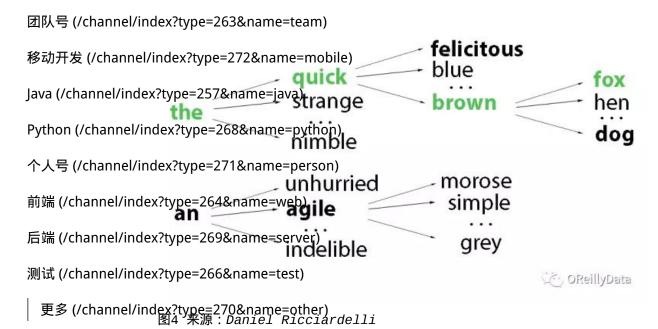
return img, all\_words

很多情况下,这个方法都能用。但是"贪婪"地使用下一个概率最大的词可能并不能产 生总体上最合适的字幕。

规避这个问题的一个可行的方法就是"<u>束搜索(Beam Search)</u>"。这个算法通过递归 的方法在最多长度为t的句子里寻找k个最好的候选者来生成长度为t+1的句子,且每 次循环都仅仅保留最好的那k个结果。这样就可以去探索一个更大的最佳的字幕空间 ,同时还能让推断在可控的计算规模内。在下图的例子里,算法维护了在每个垂直时 间步骤里的一系列k=2的候选句子(用粗体字表示)。

下载《开发者大全》





# 局限和讨论

神经图片字幕生成器给出了一个有用的框架,能学习从图片到人能理解的图片字幕间的映射。通过训练大量的图片-字幕对,这个模型学会提取从视觉特征到相关语义信息的关系。

但是,对于一个静态图片,我们的字幕生成器是关注图片里有利于分类的特征,而这并不一定是有利于字幕生成的特征。为了改进每个特征里与字幕相关的信息量,我们可以把这个图片向量模型(这个用来编码特征的VGG-16模型)作为整个字幕生成模型的一部分。这就可以让我们能更精细地调优图片编码器来更好地承担字幕生成的角色。

而且,如果我们去仔细地观察生成的字幕,就会发现它们其实相当的模糊与普通化。 用下面这个图片-字幕对为例:

下载《开发者大全》





团队号 (/channel/index?type=263&name 移动开发 (/channel/index?type=272&name= Java (/channel/index?type= name=java) Python (/channel/ind 个人号 (/channel/ind 前端 (/channel/index?type=264&name=web) 后端 (/channel/index?type=265%name=server)t



测试 (/channel/index2type=266&name=test), 图片来自MS COCO数据集

更多 (/channel/ind 这个图片当然是个铁玻璃站立在树旁边"。但是如果看看其他的图片,我们就可能注意 到它会对于任何有长颈鹿的图片都生成"长颈鹿站立在树旁边",因为在训练集里,长 颈鹿通常都出现在树的附近。

#### 下一步工作

首先,如果你想改进这里介绍的模型,请阅读以下谷歌的开源"Show and Tell网络"。 它可以用Inception-v3图片向量和MS COCO数据集来训练。

目前最前沿的图片字幕模型包含了一个视觉注意力机制。可以让模型在生产字幕时, 发现图片里的引起兴趣的区域来有选择地关注图片内容。

同时,如果你有对最前沿的字幕生成器的实现感兴趣,请阅读这个论文《展示、关注 和说出:使用视觉注意力的神经图片字幕生成》

注意:别忘了访问GitHub上与这篇文章对应的Python代码和iPython notebook。

这篇博文是O'Reilly和TensorFlow的合作产物。请阅读我们的编辑独立声明。

# 作者介绍:

#### Raul Puri

Paul Puri是加州大学伯克利分校CO 2017届毕业的本科生研究人员。Raul已经对多个 领域的研究项目做出了贡献,包括但不限于:机器人和自动化、计算机视觉、医疗成 图、生物记忆设备等。不过所有这些研究工作都专注于机器学习和机器学习系统在安 全、自主驾驶、自然语言处理、计算机视觉和机器人里面的应用。Raul也非常热情于 载《开友者大全》 下载 (/download/dev.apk) 通过数授应用机器学习概念课程来回馈社会。他还是多门伯克利分校机器学习课程的

第9页 共11页 2017/11/30 下午6:01 最新 (/) 助教和讲师。

团队号 (/channel/index nip & Eizis & Alime=team)

移动开发 (/channel/index?type=2/2&name=mobile)
移动开发 (/channel/index?type=2/2&name=mobile) 金融和工业的自然语言处理、计算机视觉、深度主动学习和自动知识发现。Dan在伯 Java (/channel/index**海利机器学**须组织里非常热心于让机器学习能为技术及非技术学生和专业人员所接触

Python (/channel/index?type=268&name=python)

This article originally appeared in English: "Caption this, with TensorFlow"

个人号 (/channel/index?type=271&name=person)

前端 (/channel/index?typ

=264&name=web)

推荐文章

后端 (/channel/index?type=269&name=server)

测试 (/channel/index?type=266&n<del>熏面蓝色</del>驗 (/html/626/201607/2649980106/1.html)即可阅读全文

10万读者睿选:【<u>2016年TOP1</u>0 更多 (/channel/index?type=270&name=other) 【2016年TOP100 (/html/626/201701/2649981531/1.html)】【2015年TOP10

CCTV大数据名人讲堂PPT&视频:【<u>万亿元产业 (/html/626/201607/2649980106/1.html)</u>】【 安全 (/html/626/201607/2649980284/1.html)】【城市 (/html/626/201607/2649980223/1.ht ml)】【农业 (/html/626/201607/2649980331/1.html)】【航运 (/html/626/201608/26499806 20/1.html)】【数据资产变现 (/html/626/201609/2649980796/1.html)】

DTiii: 【1203家大数据产业地图PPT及下载 (/html/626/201611/2649981121/1.ht

上下滑动查看更多精选专题

微信公众号:软件定义世界(

微信ID: SDx-SoftwareDefinedx

SDX)

- ●软件定义世界, 数据驱动未来;
- ❷ 大数据思想的策源地、产业变革的 指南针、创业者和VC的桥梁、政府和 企业家的智库、从业者的加油站;
- ❸个人微信号: sdxtime,

邮箱:sdxtime@126.com:

=>> 长按右侧二维码关注。



底部新增导航菜单,下载200多个精彩PPT,持续更新中!

下载《开发者大全》



#### 最新 (/) 阅读4 ₺0

团队号 (/channel/index?type=263&name=team)

移动开发 (/channel/index?type=272&name=mobile)

Java (/channel/index?type=257&name=java) 软件定义世界(SDX)更多文章

Python (/channel/index?type=268&name=python) 十货:揭秘信息可视化图表的设计方法 (/html/626/201705/2649982379/1.html)

个人号 (/channel/index?type=28718-pape=106500用实践 (/html/626/201705/2649982376/1.html)

2374/1.html)

后端 (/channel/index?type=269&name=server)

麦肯锡万字报告对比中美AI竞争力:学术、产业生态、算法、数据、计算力 (/html/626/20

测试 (/channel/index?秋始至260&h37~de=test)

#### 猜您喜欢

基于用户的协同过滤算法 (/html/178/201604/2650106997/1.html)

大型网站架构体系的演变 (/html/331/201704/2651688432/1.html)

15张PPT看清:为什么同样很忙,但你始终成不了大牛? (/html/202/201701/2649715459 /1.html)

CDH大数据平台搭建 (/html/344/201606/2652453276/1.html)

震惊小伙伴的单行Python代码 (/html/160/201607/2649639202/1.html)

Copyright © 十条网 (http://www.10tiao.com/) | 京ICP备13010217号 (http://www.miibeian.gov.cn/) | 关于十条 (/html/aboutu s/aboutus.html) | 开发者大全 (/download/index.html)

下载《开发者大全》

