

字符识别在图片文字提取中的应用研究：机器学习方法的探索

栗浩宇¹

东北大学 资源与土木工程学院, 辽宁省沈阳市

摘要: 在网络信息获取的过程中,经常会遇到一些网站上的文档或者文章中的字符无法复制的情况。这种情况下,如果需要获取这些信息,就需要寻找其他的方法。例如,可以通过截图的方式获取这些信息,但是这样就变成了图片,无法直接从图片中获取文字。这就需要使用光学字符识别(OCR)技术来从图片中提取文字。本实验的目的是探索如何使用OCR技术从无法复制的网络文档中提取文字信息。首先,收集了一系列不同字体的字母图片作为训练数据。然后,从这些图片中提取出有用的特征,并使用这些特征和对应的标签来训练字符识别分类器。在小型数据集上的实验验证了该方法的可行性,并确定了基于机器学习的图片文字识别的整体流程框架。在实验过程中,采用了K-近邻(KNN)算法进行模型训练。通过对训练集和测试集的特征提取,成功训练了KNN字符识别分类器。在测试集上的评估结果表明,该分类器能够准确地识别出图片中的字母,实现了从图片中提取字母的目标。此外,还开发了一个用于分割图像并获取边界框的函数。这个函数可以有效地将图像分割成多个小块,每个小块包含一个字符。然后,对每个小块提取特征并预测标签,从而实现了从图片中提取文字的目标。本研究表明,使用机器学习技术进行图片文字识别是可行的,并且有很大的应用前景。但是这只是一个初步的研究,还有许多不足之处。期待在未来的工作中进一步优化方法,使其能够处理更复杂的图片,提取更多类型的文字,以应用在更多领域。

关键词: 字符识别; KNN 算法; 机器学习

Machine Learning for Image Text Extraction via Character Recognition

Li Haoyu¹

1. Northeastern University, College of Resources and Civil Engineering, Shenyang, Liaoning Province, China.

Abstract: In the process of obtaining network information, we often encounter some documents or articles on the website that cannot be copied. In this case, if you need to obtain this information, you need to find other methods. For example, you can get this information by taking screenshots, but this becomes a picture, and you can't directly get the text from the picture. This requires the use of optical character recognition (OCR) technology to extract text from images. The purpose of this experiment is to explore how to use OCR technology to extract text information from network documents that cannot be copied. First, a series of letter pictures with different fonts were collected as training data. Then, useful features were extracted from these pictures, and these features and corresponding labels were used to train the character recognition classifier. The experiment on a small data set verified the feasibility of the method and determined the overall framework of image text recognition based on machine learning. In the experiment process, the K-nearest neighbor (KNN) algorithm was used for model training. By extracting features from the training set and test set, the KNN character recognition classifier was successfully trained. The evaluation results on the test set show that the classifier can accurately identify the letters in the picture, achieving the goal of extracting letters from the picture. In addition, a function for segmenting images and obtaining bounding boxes was developed. This function can effectively segment the image into multiple small blocks, each containing a character. Then,

features are extracted and labels are predicted for each small block, thus achieving the goal of extracting text from the image. This study shows that it is feasible to use machine learning technology for image text recognition, and has great application prospects. But this is only a preliminary study, and there are still many shortcomings. We look forward to further optimizing the method in future work, so that it can handle more complex images, extract more types of text, and apply to more fields.

Key words: OCR; K-Nearest Neighbor; machine learning

引言: 我们身处的信息时代, 大数据、算法、机器学习等等这些新兴事物, 已经渗透到每个人的日常生活之中。信息时代给我们的生活带来了极大的便利^[1], 然而, 我们在获取网络信息的过程中, 经常会遇到一些挑战。其中一个常见的问题是, 一些网站上的文档或者文章中的字符无法复制。这可能是由于版权保护、防止数据泄露或者其他的一些原因。这种情况下, 如果需要获取这些信息, 就需要寻找其他的方法。

光学字符识别 (OCR) 技术为我们提供了一种可能的解决方案。光学字符识别 (OCR), 是极其重要的人工智能技术之一, 符合国家四个面向重大需求^[1]。OCR 技术可以从图片中提取文字, 使得我们可以从无法复制的文档中获取信息。然而, OCR 技术的应用并不简单。它需要处理各种类型的图片, 包括但不限于截图、扫描件、照片等。此外, OCR 技术还需要能够准确地识别出图片中的文字, 这需要大量的训练数据和复杂的机器学习算法。

研究的目的是探索如何使用 OCR 技术从无法复制的网络文档中提取文字信息。我将介绍我的研究方法, 包括我如何提取特征, 以及我如何训练我的 KNN 字符识别模型。

我还将展示我的实验结果, 包括我的模型在测试集上的性能, 以及我的模型在实际图片上的表现。

1 实验过程分析

1.1 训练 KNN 模型的详细步骤

1.1.1 定义数据集路径

-代码: "dataFolder='p_dataset_26';"

-解释: 设置数据集存放的文件夹路径。这里

的"p_dataset_26"是数据集路径。

1.1.2 创建 imageDatastore 对象

-代码: "imaes=imageDatastore(dataFolder,'IncludeSubfolders',true,'LabelSource','foldernames');"

-解释: 创建一个 imageDatastore 对象来管理图片数据。这个对象会包含指定文件夹下的所有图片, 并根据文件夹名自动标记图片的标签。

1.1.3 分割数据集为训练集和测试集

-代码: "[trainingSet,testSet]=splitEachLabel(image s,0.7,'randomize');"

-解释: 将数据集分为 70% 的训练集和 30% 的测试集。"splitEachLabel"函数确保每个类别的图像都按照相同的比例分割。

1.1.4 提取特征

-代码: "trainingFeatures=featureExtractor(training Set);"和 "testFeatures=featureExtractor(testSet);"

-解释: 使用 "featureExtractor"函数从训练集和测试集中提取特征。这些特征将用于训练和测试 KNN 模型。

1.1.5 设置 KNN 分类器参数

-代码: "numNeighbors=5;"和 "distanceMetric='euclidean';"

-解释: 设定 KNN 分类器的参数。这里使用 5 个邻居 (k=5) 和欧几里得距离度量。

1.1.6 训练 KNN 分类器

-代码: "knnClassifier=fittknn(trainingFeatures,trainingLabels,'NumNeighbors',numNeighbors,'Distance',distanceMetric);"

-解释: 用训练集的特征和标签来训练 KNN 分类器。

1.1.7 在测试集上评估分类器

-代码: “predictedLabels=predict(knnClassifier,testFeatures);”

-解释: 使用训练好的 KNN 分类器对测试集的特征进行预测, 得到预测标签。

1.1.8 计算并显示准确率

-代码: “accuracy=sum(predictedLabels==testLabels)/numel(testLabels);”

-解释: 计算测试集上的预测准确率, 并打印出来。

1.2 预测和识别字符的详细步骤

1.2.1 加载训练好的 KNN 模型

-代码: “model=load('knnModel.mat');knnModel=model.knnClassifier;”

-解释: 加载之前保存的 KNN 模型。这里“.mat”文件中的变量名称为 “knnClassifier”

1.2.2 分割图像并获取边界框

-代码: “imagePath='hello_world.png';img=imread(imagePath);”和 “bboxes=splitAndDisplayBoundingBoxes(imagePath,widthThreshold,areaThreshold,heightThreshold);”

-解释: 读取要识别的图像文件, 并使用 “splitAndDisplayBoundingBoxes” 函数分割图像, 提取字符的边界框。阈值用于确定字符的大小和形状。

1.2.3 初始化识别文本和输出图像

-代码: “recognizedText='';outputImage=img;”

-解释: 初始化一个空字符串来存储识别的文本和一个输出图像。

1.2.4 对于每个边界框, 提取特征并预测标签

-代码: 包含在 “for” 循环中的一系列步骤, 如裁剪图像、调整大小、保存临时文件、创建 “ImageDatastore” 提取特征、使用 KNN 模型预测标签。

-解释: 对于每个边界框, 首先裁剪出字符图像, 调整其大小, 并将其保存为临时文件。然后创建 “ImageDatastore” 对象, 使用 “featureExtractor” 函

数提取特征, 最后使用 KNN 模型进行预测

1.2.5 将预测标签添加到识别文本和输出图像

-代码: “recognizedText=[recognizedText,label];”

和 “outputImage=insertText(outputImage,position,label,'AnchorPoint','Center');”

-解释: 将预测出的标签添加到识别文本字符串中, 并将标签写在输出图像的相应位置上。

1.2.6 删除临时文件

-代码: “delete(tempFileName);”

-解释: 处理完一个字符后, 删除其对应的临时文件

1.2.7 输出识别的文本

-代码: “disp(['识别的文本:',recognizedText]);”

-解释: 在命令窗口显示识别出的完整文本。

1.2.8 显示带有文本标签的图像

-代码: “imshow(outputImage);”

-解释: 显示最终的图像, 其中包含了识别出的字符和其在原图中的位置。

1.3 “splitAndDisplayBoundingBoxes” 函数详解

函数 “splitAndDisplayBoundingBoxes” 的目的是从给定的图像中提取字符的边界框并对这些边界框进行处理

1.3.1 读取图像

-代码: “img=imread(imagePath);”

-解释: 使用 MATLAB 的 “imread” 函数读取指定路径的图像。

1.3.2 转换为灰度图像

-代码: “grayImg=rgb2gray(img);”

-解释: 将读取的彩色图像转换为灰度图像。

1.3.3 应用阈值化以获得二值图像

-代码: “binaryImg=imbinarize(grayImg);”

-解释: 对灰度图像应用阈值化操作, 将其转换为二值图像, 以便于后续处理。

1.3.4 使用形态学操作以分离字符

-代码: “se=strel('rectangle',[1,5]);erodedImg=imero
de(binaryImg,se);dilatedImg=imdilate(erodedImg,s
e);”

-解释: 使用形态学操作(腐蚀和膨胀)来分离图像
中的字符。这些操作有助于减少字符间的相连。

1.3.5 查找连通组件

-代码: “[~,L]=bwboundaries(dilatedImg,'noholes');s
tats=regionprops(L,'BoundingBox','Area');”

-解释: 识别二值图像中的连通组件,并计算每个组
件的边界框和面积。

1.3.6 过滤掉小的连通组件

-代码: “stats=stats([stats.Area]>areaThreshold);sta
ts=stats(arrayfun(@(s)s.BoundingBox(4)>=heightTh
reshold,stats));stats=stats(arrayfun(@(s)s.Boundin
gBox(3)>=widthThreshold/2,stats));”

-解释: 移除面积、高度或宽度小于指定阈值的连通
组件,以排除非字符元素。

1.3.7 初始化边界框数组

-代码: “bboxes=[];”

-解释: 创建一个空数组用于存储符合条件的边界框。

1.3.8 处理每个连通组件的边界框

-代码: 在 “for” 循环中处理每个连通组件。

-解释: 对于每个连通组件,检查其边界框的宽度。
如果宽度大于阈值,则将该边界框分割为两个新的
边界框(左半部分和右半部分),否则保留原始边
界框。

1.3.9 将边界框添加到数组中

-代码: “bboxes=[bboxes;bb1;bb2];”或“bboxes=[bb
oxes;bb];”

-解释: 根据边界框的宽度,将分割后的或原始的边
界框添加到 “bboxes” 数组中。

这个函数最终返回一个包含所有符合条件的
边界框的数组。这些边界框随后可以用于字符识别

和其他处理。

1.4 “featureExtractor”函数的解释

featureExtractor 函数的主要目的是从图像数
据存储中提取特征。函数首先初始化一个空的特征
数组,然后遍历图像数据存储中的每一张图像。对
于每一张图像,函数使用 extractHOGFeatures 函数
提取其 HOG 特征,并将这些特征添加到特征数组
中。当所有图像都被处理后,函数使用 reset 函数
重置图像数据存储的读取状态,以便下次可以从头
开始读取。最后,函数返回包含所有图像特征的特
征数组。这个特征数组可以用来训练机器学习模型。

2. 实验结果

2.1 KNN 模型在测试集上的准确率

当我的模型被训练并在测试集上预测正确率
时,我得到了令人满意的结果。图 1 展示了使用训
练好的 KNN 分类器对测试集的特征进行预测所
得到的准确率,为 96.49%,显示出模型在实际应
用中的有效性和准确性。

```
>> trainAndSaveKNNModel
The accuracy of the KNN classifier is: 96.49%
```

图 1 KNN 模型在测试集上的准确率

Fig.1 KNN model accuracy on the test set

2.2 KNN 模型在小型数据集上的成功运用

在这个基础上,我选择了一张图片,并运行了
这个模型,成功识别出了图片中的字符,图 2 是模
型运行所识别出字符的结果。

```
>> KnnCharacterRecognition
识别的文本: HELLOWORWD
```

图 2 识别图片中字符的结果

Fig.2 Character recognition results

图 3 是所识别的字符在原图中的位置，可以清晰地看出它在图片上的内容。



图 3 字符在原图中的位置

Fig.3 Characters' positions in the original image.

这表明我的模型在经过充分的训练后，能够在实际应用中取得可靠的性能表现。

综合来看，尽管我选择的数据集比较有限，但通过对实验结果的观察和分析，我可以自信地说，我的 KNN 模型在正确识别任务中表现出了显著的能力，成功的完成了这个任务。这些结果为进一步深入研究和实际应用奠定了坚实的基础。

3 讨论

3.1 优点

数据处理：代码首先定义了数据集的路径，并创建了一个 imageDatastore 对象来管理图片数据。这使得程序可以处理大量的图片数据，而无需一次性加载所有的图片到内存中。

特征提取：HOG 特征作为图像的特征，这是一种在计算机视觉和图像处理中广泛使用的特征描述符，它通过计算和统计图像局部区域的梯度方向直方图来构建特征，这对于物体检测非常有用。

模型训练和评估：我使用了 K-近邻(KNN)分类器进行模型训练，并在测试集上评估了分类器的性能。这使得我的代码可以根据实际的性能反馈进行调整和优化。

3.2 缺点

特征提取：虽然 HOG 特征在许多情况下都表现得很好，但它可能不适用于所有类型的图像。对

于一些复杂的图像，可能需要更复杂的特征才能获得好的性能。

模型选择：虽然 KNN 分类器在许多情况下都表现得很好，但利用 KNN 算法对目标样本做预测时，需要对目标样本和所有训练样本之间的距离进行计算^[3]，当训练样本较大时，KNN 法会因为巨大的计算量而产生较大的误差^[4]，所以对于一些复杂的问题，可能需要使用更复杂的模型，如 bp 神经网络。

3.3 改进方案：

使用更先进的 OCR 技术：随着技术的发展，OCR 技术也在不断进步。例如，现在的 OCR 技术已经可以利用人工智能和机器学习来实现更高级的字符识别，如识别不同的语言或手写风格^[5]。此外，一些 OCR 技术还引入了图像质量评分、旋转校正等高级特性，以提高模型的准确性^[6]。

我需要在未来的研究中考虑如何改进我的方法，以克服这些缺点，提高 OCR 系统的性能。

4 结论

在本研究中，我通过使用一个小型数据集来验证了基于机器学习的图片文字识别的可行性，并确定了整体流程框架。通过实验结果和分析，我得出以下结论：

首先，我使用一个小型数据集进行了训练和测试，结果显示基于 KNN 算法的字符识别系统在识别准确率上表现出色。这表明，在较小的数据集上，机器学习算法能够有效地学习和识别不同字符，为后续扩展到更大数据集奠定了基础。

其次，我设计了一个整体的流程框架，包括数据集的加载、训练集和测试集的分割、特征提取、KNN 分类器的训练和预测等步骤。通过实验验证，我确认了该流程框架的合理性和有效性。这一流程

框架提供了一个清晰的指导,使得我能够系统地处理图片文字识别问题。

此外,我还提出了一种边界框分割的算法,用于定位和提取字符图像。该算法基于形态学操作能够有效地将字符分离出来。通过对分割后的字符图像进行特征提取和预测,我成功地实现了对图片文字的识别。

通过本次研究,我验证了基于机器学习的图片文字识别在小型数据集上的可行性,并确定了整体流程框架。我的实验结果表明,这一框架在字符识别准确率和处理效率方面表现出色。未来,我将进一步扩展数据集规模、尝试其他机器学习算法,并考虑应用于更广泛的领域,如自动化办公、图像检索、车牌识别等。

需要注意的是,尽管在小型数据集上验证了可行性,但在实际应用中,还需要考虑更大规模的数据集和更复杂的场景。同时,我也意识到在不同领域和应用中,可能会涉及到更多的挑战和需求,比如多语言文字识别、不同字体和风格的文字等。因此,进一步的研究和改进仍然是必要的。

5 结束语

通过这门课程近八周的学习,我的进步是巨大的,这门课开始的时候我从来没想过会学到这么多丰富的内容,更没想到自己可以独立完成这个结课作业。

结课作业是充满挑战的,有了字符识别这个想法后,最初也是毫无头绪、无从下手,但是从最基本的知识开始,查找一篇篇文献、了解一个个内容,并且多次向他人请教,期间多次怀疑自己是否能够

完成这个问题,但庆幸最后的结果是好的,当我看到“HELLO WORLLD”出现在屏幕上时,才发现自己是真的完成了这个任务,自己真的做到了。

我的进步离不开老师的帮助,在课堂上,当我听不懂不知道如何操作时,老师总是走到我身后给我耐心指导;当老师讲完一个知识点但发现大家都没听懂时,也会不厌其烦的再次讲解;从第一节课因怕被老师责备没好好听课而不敢向老师询问问题,到后来主动问老师自己的问题,赵老师的和蔼和专业让我对这门课程产生了浓厚的兴趣和信心。

通过这门课程的学习和结课作业的完成,我不仅学到了很多有价值的知识和技能,也锻炼了我的思维和创造力。我对自己的进步感到非常自豪。

再次衷心地感谢老师的辛勤付出!

参考文献:

- [1] 信息时代,文化的变与不变[J].江南,2023,(1): 149-163
- [2] 刘成林,金连文,殷绪成等.《中国图象图形学报》文档图像智能处理与识别专栏简介[J].中国图象图形学报,2023,28(08):2221-2222.
- [3] 向文平,马弢,梁瑜等.基于熵定权的 KNN 建筑内定位算法研究[J].现代雷达,2021,43(07):32-37.DOI:10.16592/j.cnki.1004-7859.2021.07.006.
- [4] 卢海钊,彭慧豪,唐滔等.基于 KNN 和 XGBoost 的室内指纹定位算法[J].电子测量技术,2023,46(02):81-86.DOI:10.19651/j.cnki.emt.2210501.
- [5] IBM Cloud Education. What is optical character recognition (OCR)?[EB]. (2022-01-05)[2023-11-16]
- [6] Goole Cloud Education. What is OCR[EB]. (2022-01-05)[2023-11-16].