

基于渐进式学习的神经网络端到端验证码识别

刘达荣,张远平,汤茂斌,李福芳

(广州大学 计算机科学与教育软件学院,广东 广州 510006)

摘要:针对验证码经过弯曲变形,无法采用传统的字符分割方法进行检测的问题,在模仿人类的渐进式学习过程的基础上,提出利用卷积神经网络,优化手写识别 MNIST 的三层网络结构,无需预先对验证码进行分割,直接对验证码进行端到端的识别。利用不确定度回收训练图片,减少训练集数据量,并利用可视化工具提高其网络识别性能。经过 55 万次训练后,生成了检测模型并对测试集验证码进行了检测,验证码识别速度达到 0.073 秒/张,准确率达到 86%。通过对比同一测试环境下的测试集,发现利用渐进式学习方法具有更高的建模效率和更好的识别准确率,并对识别错误的验证码进行分析。

关键词:渐进式学习;卷积神经网络;验证码;无分割;端到端

中图分类号:TP311

文献标识码:A

文章编号:1673-629X(2018)09-0016-04

doi:10.3969/j.issn.1673-629X.2018.09.004

End-to-end Verification Code Identification of Neural Network Based on Progressive Learning

LIU Da-rong, ZHANG Yuan-ping, TANG Mao-bin, LI Fu-fang

(School of Computer and Education Software, Guangzhou University, Guangzhou 510006, China)

Abstract: For authentication code can't be measured by traditional character segmentation method after bending deformation, on the basis of the progressive learning of imitating human, we put forward convolution neural network to optimize three layers network structure of handwritten recognition MNIST, without prior segmentation for authentication code, directly identifying in the end-to-end to authentication code. The training images are recovered according to uncertainty to reduce the amount of training data, and the network recognition performance is improved with visual tools. After training by 550 000 times, the detection model is generated and the verification code from the test set is detected. The recognition speed of the verification code reaches 0.073 seconds per piece, and the accuracy is 86%. By contrast with the test set in the same test environment, it is found that the progressive learning method has higher modeling efficiency and better recognition accuracy. The wrong verification code is analyzed.

Key words: progressive learning; convolution neural network; verification code; no segmentation; end to end

1 概述

验证码^[1]作为区分人类和机器的测试工具,在图灵测试上占有一席之地。如铁道部在高铁票购买过程中,进行验证码识别以防止黄牛党抢票;论坛注册时,提供的验证码测试防止脚本自动注册。目前,比较常用的验证码都会包含一些扭曲变形的字符进行测试,测试者必须能看懂并输入正确的字符串才能通过。对于人类来说,扭曲变形的字符并不难识别,但对于机器就相当困难,一个好的验证码可以有 80% 的概率被人类识别^[2],但对于机器只有 0.01% 的概率被识别。一

个典型的验证码样例^[3]如图 1 所示。

To continue, please type the characters below:



图 1 谷歌提供的验证码样例

(包含一串扭曲变形的字符)

传统的光学字符识别^[4]的方法是:首先定位字符

收稿日期:2017-08-24

修回日期:2017-12-27

网络出版时间:2018-05-16

基金项目:国家自然科学基金(61472092)

作者简介:刘达荣(1988-),男,硕士研究生,研究方向为图像处理、深度学习;张远平,教授,研究方向为计算机算法设计与分析;汤茂斌,副教授,研究方向为人工智能、数据挖掘、系统分析与设计;李福芳,副教授,研究方向为计算机网络、网络(分布式)计算。

网络出版地址: <http://cnki.net/kcms/detail/61.1450.TP.20180515.1645.006.html>

bdrvmw	cgfatt	ddjycc	ejqzyr
bdrvmw	cgfatt	ddjycc	ejqzyr
fcbjwz	ftsflm	gnqbcm	gzyanm
fcbjwz	ftsflm	gnqbcm	gzyanm

(a)没有扭曲变形的简单验证码

aawjlg	bggzgy	cdhlnj	dasric
aawjlg	bggzgy	cdhlnj	dasric
ddimfc	dpshje	eaansn	lzhwz
ddimfc	dpshje	eaansn	lzhwz

(b)经过扭曲变形的困难验证码

图 3 验证码

2.2 卷积层网络结构参数和特征可视化

如图 2(b)所示,图片以像素形式转换成矩阵输入网络。在卷积层 1 中,对数据层以窗口大小 5×5、步长

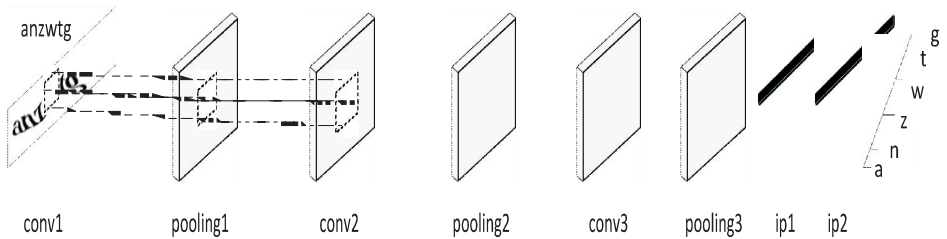


图 4 网络结构

2.3 输出层

对于输出层,为预测每个验证码图片的字符概率,采用 SoftmaxWithLoss 多分类函数:

$$P(y_i | v) = \frac{e^{v^T \omega_i}}{\sum_{k=1}^K e^{v^T \omega_k}} \tag{1}$$

即对于隐层的输出 v ,第 i 个输出单元 y_i 的概率为 $P(y_i | v)$, ω 为隐层与输出层的权重值。在输出层,由于不能直接输出字符,采用如下映射函数:

$$\Theta(y_i) = \begin{cases} '0' \cdots '9', y_i = 1 + 62i, 2 + 62i, \cdots, 10 + 62i \\ 'A' \cdots 'Z', y_i = 11 + 62i, \cdots, 36 + 62i \\ 'a' \cdots 'z', y_i = 37 + 62i, \cdots, 62 + 62i \end{cases} \tag{2}$$

其中, y_i 为输出层第 i ($i=0,1,\cdots,5$) 个单元,对应从数字 0 到字母 z 的其中一个,其余如此类推。如第一个输出单元 y_0 ,输出最大概率对应的横坐标为 37,代表输出字符 a;第二个输出单元 y_1 ,输出最大概率对应的横坐标为 112,代表输出字符 n,如此类推,如图 5 所示。

以上文的验证码图片“anzwtg”为例,图中为六位字符经过识别后每一位对应的最大概率,即 x 轴的第 0 号到 61 号为第一个字符的分布区间,第 62 到 123 为第二个字符的分布区间,依此类推。

为 2 进行特征提取,把原图片 50×180 压缩编码以 25×90 的矩阵向向量输出到下一层,直到最后的分类输出层。在卷积层 1、2、3 的输出中,第一个输出图学习到了验证码的边缘,第二个输出图学习到了验证码的局部特征,第三个输出图学习到了验证码的全局特征。

基于 Caffe 框架^[18],硬件采用 GTX960 GPU 进行训练,基础学习率 (base_lr) 初始值设置为 0.01,动量为 0.9,权值衰减为 0.000 5,学习率策略采用多项式,权重 (power) 为 0.75,最大迭代次数 (max_iter) 为 20 000,并使用 GPU 进行加速训练,迭代次数 (iter) 为两千次。由于选取学习率策略为多项式模式 (Poly),故在训练过程中的有效学习率为 $base_lr * (1 - iter / max_iter)^{(power)}$,训练两万张图片。

对超参数进行设置和建模,以及对每一层的输入与输出进行特征可视化,如图 4 所示。

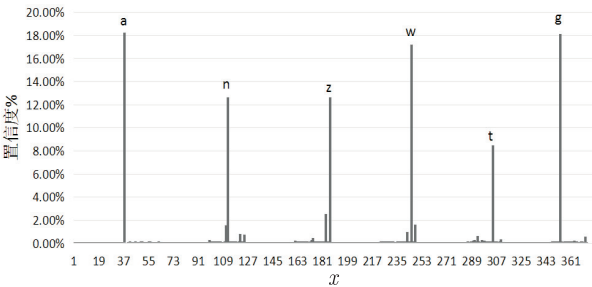


图 5 提取输出层的概率分布图

2.4 利用不确定度回收训练图片

文中采用“best-vs-second-best”来统计已经分类正确但具有较大不确定性的图片,重新放回下一次迭代训练,以此减少训练样本大小,公式如下:

$$\eta = \frac{1}{d} \cdot \sum_{i=1}^d \frac{\argmax\{P(y_i) \setminus \argmax P(y_i)\}}{\argmax P(y_i)} \tag{3}$$

其中, d 为迭代次数; $P(y_i)$ 为每个字符的概率大小。分母是每次迭代的字符最大概率,分子是每次迭代的字符次大概率,即选取总不确定度 η 值最大的图片,重新放入下一次迭代训练,以此来减少训练图片规模。

3 实验结果与分析

文中提出渐进式学习,其准确率和学习过程如图 6 所示。其中,实线为直接采用困难验证码图片训练

的效果,虚线为先进行简单图片预训练,当准确率达到 98% 时,替换成困难验证码图片进行训练的效果。

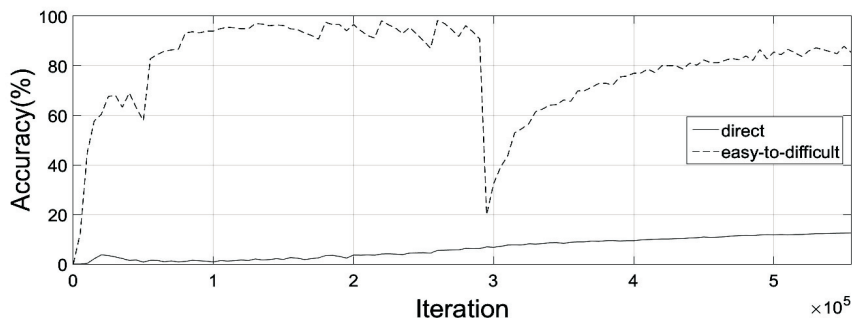


图6 迭代次数与准确率

从图6可以明显看出,采用渐进式学习在同一训练时间内准确率更高,达到86%。本次实验中,训练的总时间约为48 h,若继续增加训练时间,准确率还可以更高。与传统方法相比较,在单张验证码识别时间上,基于SVM识别需要0.76 s^[19],在本测试中需要0.073 s,识别效率约为传统方法的10倍。

4 结束语

结合传统神经网络的结构进行优化,采用模仿人类学习过程的渐进式学习模式,提出一种无分割、端到端验证码识别算法。该算法充分利用卷积神经网络的自学习特性,简化传统验证码识别中字符分割等人工干预手段,使得网络具有抗旋转的优良特性,并利用概率统计方法减少训练样本,最终使得神经网络具有收敛速度快、检测效果好的特性。在后续工作中,将在背景加入更多干扰来进行测试,以期提高网络的鲁棒性。

参考文献:

- [1] VON AHN L, BLUM M, HOPPER N J, et al. CAPTCHA: using hard AI problems for security[M]//Advances in cryptography. Berlin: Springer, 2003: 294–311.
- [2] CHELLAPILLA K, SIMARD P Y. Using machine learning to break visual human interaction proofs (HIPs) [C]//Advances in neural information processing systems. [s. l.]: [s. n.], 2004: 265–272.
- [3] AHN L V, MAURER B, MCMILLEN C, et al. reCAPTCHA: human-based character recognition via web security measures[J]. Science, 2008, 321(5895): 1465–1468.
- [4] JADERBERG M, VEDALDI A, ZISSERMAN A. Deep features for text spotting[C]//European conference on computer vision. Berlin: Springer, 2014: 512–528.
- [5] LU Yi. Machine printed character segmentation: an overview [J]. Pattern Recognition, 1995, 28(1): 67–80.
- [6] 吕刚, 郝平. 基于神经网络的数字验证码识别研究[J]. 浙江工业大学学报, 2010, 38(4): 433–436.
- [7] 左保河, 石晓爱, 谢芳勇, 等. 基于神经网络的网络验证码

- 识别研究[J]. 计算机工程与科学, 2009, 31(12): 20–22.
- [8] 邓介一, 刘黎志, 谭培祥. 基于神经网络的数字字符识别系统设计与实现[J]. 软件导刊, 2017, 16(5): 47–50.
- [9] GOODFELLOW I J, BULATOV Y, IBARZ J, et al. Multi-digit number recognition from street view imagery using deep convolutional neural networks [C]//Computer vision and pattern recognition. [s. l.]: [s. n.], 2013.
- [10] 付先琨. 基于RPROP人工神经网络对验证码识别的研究与实现[D]. 重庆: 重庆大学, 2011.
- [11] KRIZHEVSKY A, SUTSKEVER I, HINTON G E. ImageNet classification with deep convolutional neural networks[C]//Proceedings of the 25th international conference on neural information processing systems. Lake Tahoe, Nevada: Curran Associates Inc., 2012: 1097–1105.
- [12] LECUN Y, BOTTOU L, BENGIO Y, et al. Gradient-based learning applied to document recognition[J]. Proceedings of the IEEE, 1998, 86(11): 2278–2324.
- [13] 高原. 基于BP神经网络的文本验证码破解[J]. 电子科技, 2012, 25(7): 37–42.
- [14] 吕霖. 基于神经网络的验证码识别技术研究[D]. 泉州: 华侨大学, 2015.
- [15] NAIR V, HINTON G E. Rectified linear units improve restricted boltzmann machines [C]//Proceedings of the 27th international conference on machine learning. Haifa, Israel: Omnipress, 2010: 807–814.
- [16] 陈超, 毛坚恒, 刘寅. 基于卷积神经网络的铁路货运网站验证码识别[J]. 指挥信息系统与技术, 2016, 7(4): 91–96.
- [17] HINTON G E, SRIVASTAVA N, KRIZHEVSKY A, et al. Improving neural networks by preventing co-adaptation of feature detectors [J]. Computer Science, 2012, 3(4): 212–223.
- [18] JIA Yangqing, SHELHAMER E, DONAHUE J, et al. Caffe: convolutional architecture for fast feature embedding [C]//International conference on multimedia. [s. l.]: ACM, 2014: 675–678.
- [19] 唐海涛. 自组织增量神经网络的验证码识别模型与算法[D]. 广州: 广东工业大学, 2016.