

Character Recognition in Natural Images

Arshiya Nain¹, Sukhwinder Singh²

¹Dept. of Electronics Engg., Punjab Engineering College University of Technology, Chandigarh, Punjab, India

²Prof., Punjab Engineering College, Chandigarh Punjab, India

Abstract: This paper tackles the problem of recognizing characters in images of natural scenes. In particular, we focus on recognizing characters in situations that would traditionally not be handled well by OCR techniques. We present an annotated database of images containing English and Kannada characters. The database comprises of images of street scenes taken in Bangalore, India using a standard camera. The problem is addressed in an object categorization framework based on a bag-of-visual-words representation. We assess the performance of various features based on nearest neighbor and SVM classification. It is demonstrated that the performance of the proposed method, using as few as 15 training images, can be far superior to that of commercial OCR systems. Furthermore, the method can benefit from synthetically generated training data obviating the need for expensive data collection and annotation.

Keywords: Object recognition, camera-based character recognition, Latin characters, digits, Kannada characters, off-line handwritten character recognition.

1. INTRODUCTION

This paper presents work towards automatic reading of text in natural scenes. In particular, our focus is on the *recognition* of individual characters in such scenes. Figures 1, 2 and 3 highlight why this can be a hard task. Even if the problems of clutter and text segmentation were to be ignored for the moment, the following sources of variability still need to be accounted for: (a) font style and thickness; (b) background as well as foreground color and texture; (c) camera position which can introduce geometric distortions; (d) illumination and (e) image resolution. All these factors combine to give the problem a flavor of object recognition rather than optical character recognition or handwriting recognition. In fact, OCR techniques can not be applied out of the box precisely due to these factors. Furthermore, viable OCR systems have been developed for only a few languages and most Indic languages are still beyond the pale of current OCR techniques. Many problems need to be solved in order to read text in natural images including text localization, character and word segmentation, recognition, integration of language models and context, *etc.* Our focus, in this paper, is on the basic character recognition aspect of the problem (see Figures 2, 3, 5 and 6). We introduce a database of images containing English and Kannada text¹. In order to assess the feasibility of posing the problem as an object recognition task, we benchmark the performance of various features based on a bag-of-visual-words representation. The results indicate that even the isolated character recognition task is challenging. The number of classes can be moderate (62 for English) to large (657 for Kannada) with very little inter-class variation as highlighted by Figures 2 and 3. This problem is particularly acute for Kannada where two characters in the alphabet can differ just by the placement of a single dot like structure. Furthermore, while training data is readily available for some characters others might occur very infrequently in natural scenes. We therefore investigate whether surrogate training data, either in the form of font generated characters or hand-printed characters, can be used to bolster recognition in such a scenario. We also present baseline recognition results on front hand printed character database to contrast the difference in performance when reading text in natural images.

2. RELATED WORK

The task of character recognition in natural scenes is related to problems considered in camera based document analysis and recognition. [1] Most of the work in this field is based on locating and rectifying the text areas (e.g. (Kumar et al., 2007), (Krempp et al., 2002), (Clark and Mirmehdi, 2002) and (Brown et al., 2007)), followed by the application of OCR techniques (Kise and Doermann, 2007). Such approaches are therefore limited to scenarios where OCR works well. Furthermore, even the rectification step is not directly applicable to our problem, as it is based on the detection of printed document edges or assumes that the image is dominated by text. [2][3][4] Methods for off-line recognition of hand printed characters (Plamondon and Srihari, 2000), (Pal et al., 2007) have successfully tackled the problem of intra-class variation due to differing writing styles. However, such approaches typically consider only a limited number of appearance classes, not dealing with variations in foreground/background color and texture.

For natural scenes, some researchers have designed systems that integrate text detection, segmentation and recognition in a single framework to accommodate contextual relationships. For instance, (Tu et al., 2005) used insights from natural language processing and present a Markov chain framework for parsing images. (Jin and Geman, 2006) introduced composition machines for constructing probabilistic hierarchical image models which accommodate contextual relationships.[5][6] This approach allows re-usability of parts among multiple entities and non-Markovian distributions. (Weinman and Learned Miller, 2006) proposed a method that fuses image features and language information (such as bi-grams and letter case) in a single model and integrates dissimilarity information between character images. [7]

Simpler recognition pipelines based on classifying raw images have been widely explored for digits recognition (see (le Cun et al., 1998), (Zhang et al., 2006) and other works on the MNIST and USPS datasets). Another approach is based on modeling this as a shape matching problem (e.g. (Belongie et al., 2002)): several shape descriptors are detected and extracted and point-by-point matching is computed between pairs of images.

3. DATA SETS

Our focus is on recognizing characters in images of natural scenes.[8] Towards this end, we compiled a database of English and Kannada characters taken from images of street scenes in Bangalore, India. However, gathering and annotating a large number of images for training can be expensive and time consuming. Therefore, in order to provide complementary training data, we also acquired a database of hand-printed characters and another of characters generated by computer fonts.

For English, we treat upper and lower case characters separately and include digits to get a total of 62 classes. Kannada does not differentiate between upper and lower case characters. It has 49 basic characters in its alpha-syllabary, but consonants and vowels can combine to give more than 600 visually distinct classes.

3.1 Natural Images Data Set

We photographed a set of 1922 images, mostly of sign boards, hoardings and advertisements but we also included a few images of products in supermarkets and shops. [9] We experimented with two types of segmentations: rectangular bounding boxes and finer polygonal segments as shown in Figure 4. For the types of features investigated in this paper, it turned out that polygonal segmentation masks presented almost no advantage over bounding boxes. Therefore, all the results presented in Section 5 are using the bounding box segmentations. Our English dataset has 12503 characters, of which 4798 were labeled as bad images due to excessive occlusion, low resolution or noise. For our experiments, we used the remaining 7705 character images. [10] Similarly, for Kannada, a total of 4194 characters were extracted out of which only 3345 were used. Figures 5 and 6 show examples of the extracted characters. These datasets will be referred to as the *Img* datasets.

3.2 Font and Hand-printed Datasets

The hand-printed data set (*Hnd*) was captured using a tablet PC with the pen thickness set to match the average thickness found in hand [11][12] painted information generated by 55 volunteers. For Kannada, a total of 16425 characters were generated by 25 volunteers. Some sample images are shown in Figure 7.

The font dataset was synthesized only for English characters. We tried 254 different fonts in 4 styles (normal, bold, italic and bold+italic) to generate a total of 62992 characters. This dataset will be referred to as the *Fnt* dataset.

4. FEATURE EXTRACTION AND REPRESENTATION

Bag-of-visual-words is a popular technique for representing image content for object category recognition. The idea is to represent objects as histograms of feature counts.[13][14] This representation quantizes the continuous high-dimensional space of image features to a manageable vocabulary of “visual words”. This is achieved, for instance, by grouping the low-level features collected from an image corpus into a specified number of clusters using an unsupervised algorithm such as *K-Means* (for other methods of generating the vocabulary see (Jurie and Triggs, 2005)). One can then map each feature extracted from an image onto its closest visual word and represent the image by a histogram over the vocabulary of visual words. We learn a set of visual words per class and aggregate them across classes to form the vocabulary. In our experiments, we learned 5 visual words per class for English leading to a vocabulary of size 310. For Kannada, we learn 3 words per class, resulting in a vocabulary of 1971 words.

4.1 Features

We evaluated six different types of local features. Not only did we try out shape and edge based features, such as Shape Context, Geometric Blur and SIFT, but also features used for representing texture, such as filter responses, patches and Spin Images, since these were found to work well in (Weinman and Learned Miller, 2006). We explored the most commonly used parameters and feature detection methods employed for each descriptor, with a little tuning, as described below. Shape Contexts (SC) (Belongie et al., 2002) is a descriptor for point sets and binary images. We sample points using the Sobel edge detector. The descriptor is a log-polar histogram, which gives a $\theta \times n$ vector, where θ is the angular resolution and n is the radial resolution. We used $\theta = 15$ and $r = 4$. Geometric Blur (GB) (Berg et al., 2005) is a feature extractor with a sampling method similar to that of SC, but instead of histogramming points, the region around an interest point is blurred according to the distance from this point. For each region, the edge orientations are counted with a different blur factor. This soothes the problem of hard quantization and allows its application to gray scale images. Scale Invariant Feature Transform (SIFT) (Lowe, 1999) are extracted on points located by the Harris Hessian-Laplace detector, which gives affine transform parameters. [15] The feature descriptor is computed as a set of orientation histograms on (4×4) pixel neighborhoods. The orientation histograms are relative to the key-point orientation. The histograms contain 8 bins each, and each descriptor contains a 4×4 array of 16 histograms around the key-point. This leads to feature vector with 128 elements. Spin image (Lazebnik et al., 2005), (Johnson and Herbert, 1999) is a two-dimensional histogram encoding the distribution of image brightness values in the neighborhood of a particular reference point. [16] The two dimensions of the histogram are d , distance from the center point, and i , the intensity value. We used 11 bins for distance and 5 for intensity value, resulting in 55-dimensional descriptors. The same interest point locations used for SIFT were used for spin images. Maximum Response of filters (MR8) (Varma and Zisserman, 2002) is a texture descriptor based on a set of 38 filters but only 8 responses. This filter is extracted densely, giving a large set of 8D vectors. Patch descriptor (PCH) (Varma and Zisserman, 2003) is the simplest dense feature extraction method. For each position, the raw $n \times n$ pixel values are vectorized, generating an n^2 descriptor. We used 5×5 patches.

CONCLUSIONS

In this paper, we tackled the problem of recognizing characters in images of natural scenes. We introduced a database of images of street scenes taken in Bangalore, India and showed that even commercial OCR systems are not well suited for reading text in such images. Working in an object categorization framework, we were able to improve character recognition accuracy by 25% over an OCR based system. The best result on the *English Img* database was 55.26% and was obtained by the multiple kernel learning (MKL) method of (Varma and Ray, 2007) when trained using 15 *Img* samples per class. This could be improved further if we were not to be case sensitive. Nevertheless, significant improvements need to be made before an acceptable performance level can be reached. Obtaining and annotating natural images for training purposes can be expensive and time consuming. We therefore explored the possibility of training on hand-printed and synthetically generated font data. The results obtained by training on hand-printed characters were not encouraging. This could be due to the limited variability amongst the writing styles that we were able to capture as well as the relatively small size of the training set. On the other hand, using synthetically generated fonts, the performance of nearest neighbor classification based on Geometric Blur features was extremely good. For equivalent size training sets, training on fonts using a NN classifier could actually be better than training on the natural images themselves. The performance obtained when training on all the font data was nearly as good as that obtained using MKL when trained on 15 natural image samples per class. This opens up the possibility of harvesting synthetically generated data and using it for training. As regards features, the shape based features, Geometric Blur and Shape Context, consistently outperformed SIFT as well as the appearance based features. This is not surprising since the appearance of a character in natural images can vary a lot but the shape remains somewhat consistent.

We also presented preliminary results on recognizing Kannada characters but the problem appears to be extremely challenging and could perhaps benefit from a compositional or hierarchical approach given the large number of visually distinct classes.

REFERENCES

- [1]. Belongie, S., Malik, J., and Puzicha, J. (2002). Shape matching and object recognition using shape contexts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [2]. Berg, A. C., Berg, T. L., and Malik, J. (2005). Shape matching and object recognition using low distortion correspondence. In *Proc IEEE Conf on Computer Vision and Pattern Recognition*, San Diego CA, June 20-25.

- [3]. Brown, M. S., Sun, M., Yang, R., Yun, L., and Seales, W. B. (2007). Restoring 2d content from distorted documents. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. Clark, P. and Mirmehdi, M. (2002). Recognising text in real scenes. *International Journal on Document Analysis and Recognition*, 4:243–257. Jin, Y. and Geman, S. (2006). Context and hierarchy in a probabilistic image model. In *Proc IEEE Conf on Computer Vision and Pattern Recognition*, New York NY, June 17-22.
- [4]. Johnson, A. E. and Herbert, M. (1999). Using spin images for efficient object recognition in cluttered 3d scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(5):433–449.
- [5]. Jurie, F. and Triggs, B. (2005). Creating efficient codebooks for visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*.
- [6]. Kise, K. and Doermann, D. S., editors (2007). *Proceedings of the Second International Workshop on Camera-based Document Analysis and Recognition CBDAR*, Curitiba, Brazil. <http://www.imlab.jp/cbdar2007/>.
- [7]. Krempp, A., Geman, D., and Amit, Y. (2002). Sequential learning of reusable parts for object detection. Technical report, Computer Science Department, Johns Hopkins University.
- [8]. Kumar, S., Gupta, R., Khanna, N., Chaudhury, S., and Joshi, S. (2007). Text extraction and document image segmentation using matched wavelets and mrf model. *IEEE Transactions on Image Processing*, 16(8):2117–2128.
- [9]. Lazebnik, S., Schmid, C., and Ponce, J. (2005). A sparse texture representation using local affine regions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(8):1265–1278.
- [10]. Le Cun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- [11]. Lowe, D. G. (1999). Object recognition from local scale-invariant features. In *Proc 7th Int Conf on Computer Vision*, Corfu, Greece.
- [12]. Pal, U., Sharma, N., Wakabayashi, T., and Kimura, F. (2007). Off-line handwritten character recognition of devnagari script. In *International Conference on Document Analysis and Recognition (ICDAR)*, pages 496–500, Curitiba, PR, Brazil. IEEE.
- [13]. Plamondon, R. and Srihari, S. N. (2000). On-line and offline handwriting recognition: A comprehensive survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(1):63–84.
- [14]. Tu, Z., Chen, X., Yuille, A. L., and Zhu, S. C. (2005). Image parsing: Unifying segmentation, detection, and recognition. *International Journal of Computer Vision*,
- [15]. Marr Prize Issue. Varma, M. and Ray, D. (2007). Learning the discriminative power-invariance trade-off. In *Proceedings of the IEEE International Conference on Computer Vision*, Rio de Janeiro, [16] Brazil. Varma, M. and Zisserman, A. (2002). Classifying images of materials: Achieving viewpoint and illumination independence. In *Proceedings of the 7th European Conference on Computer Vision*, Copenhagen, Denmark, volume 3, pages 255–271. Springer-Verlag. Varma, M. and Zisserman, A. (2003). Texture classification: Are filter banks necessary? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.