

EFFICIENT WORD SEGMENTATION DRIVEN BY UNCONSTRAINED HANDWRITTEN PHRASE RECOGNITION

Jaehwa Park, Venu Govindaraju, and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition
Department of Computer Science and Engineering
SUNY, Buffalo
Amherst, NY 14260

ABSTRACT

An efficient system which finds the best match between an input image and a lexicon is presented. To capture writing style of spacing between words and characters prime stroke analysis based on statistical methods is introduced. A method for estimating bound on number of characters without actual recognition is also presented. For system efficiency, before actual recognition, classified groups of word segments and eligible subset of lexicons are generated as hypotheses. The hypotheses are verified and ordered by a lexicon driven word recognizer. We have tested our approach in the street name recognition/interpretation for US mail stream. Experimental results are encouraging.

1. INTRODUCTION

The practical environment of unconstrained phrase recognition is different from that of word recognition - one or more words exist within the input image and the number of words and word break points are generally unknown. If the phrase recognition problem is extended to allow partial matching of lexicon to image, the task becomes challenging as contextual information is seldom available.

The general approaches to solve the problem follow these steps; *word segmentation*, *word recognition* and *postprocessing*. Word segmentation precedent to recognition is primarily related to geometric relationship between components. But separating an unconstrained handwritten phrase into words is not so simple, since handwritten phrase lacks the uniform spacing which is normally found in machine-printed phrases. To overcome this problem, several distance measuring methods [1] [2] have been proposed. These methods are simple, efficient and accurate in estimating word gaps but they have been weak in adapting to global spacing style and often cause over or under word segmentations.

One successful approach to phrase recognition has been introduced in [3], which combines word segmentation using recognition. The word break point detector is implemented by a neural network. An exhaustive combination of word segments is submitted to a word recognizer with given a lexicon. Statistical characteristics of character segmentation are used to limit number of combinations. Missing word break points and efficiency were cited as cause for problems.

In this paper, we describe an efficient method of phrase recognition extending ideas in [3]. Rule based stroke segmentation and a statistical prime stroke analysis are implemented to detect the dominant stroke frequency characterized by the global spacing style of a writer. This analysis is also applied to word gap classification which provides numeric confidences and estimation of number of characters. For efficiency reasons, the word gap confidence and estimated number of characters of each word segment are used as basic parameters to generate a subset of possible lexicons as lexeme hypotheses and these hypotheses are tested and ordered by a general lexicon driven word recognizer.

2. WORD SEGMENTATION

The primary goal of word segmentation is to detect possible word gaps in order to drive a word recognizer with proper subset of lexicons.

Even though in handwriting text various inter-character and inter-word spacing styles exist from different writing style and flourishes, the assumption that space between two adjacent words is larger than that of characters is acceptable in general. The underlying problem of unconstrained word segmentation is not just the different ways of measuring distance between components but the fact of accurate inter-character gap is hard to detect.

A writer's spacing style is generally ruled by the strokes or character bodies which are located within the primary horizontal writing band (imaginary) established throughout lines of handwritten text. The others, such as descenders, ascenders and decorative strokes are not ruled by the inter-character spacing style either.

If we define a prime stroke as a stroke which has effective vertical components laid on the primary writing band (exclude flourishing and ligature), we can generalize spacing style in terms of the prime strokes. Dominant inter-character spacing style can be judged in terms of prime spatial periods. Further, if the dominant prime stroke frequency is known, word break points can be detected based on the spatial period of stroke frequency assuming that word spacing is larger than character spacing.

Based on the above assumptions, an analytic handwriting spacing model can be considered as one where

- A character consists of one or more pieces of strokes called *primitives*.

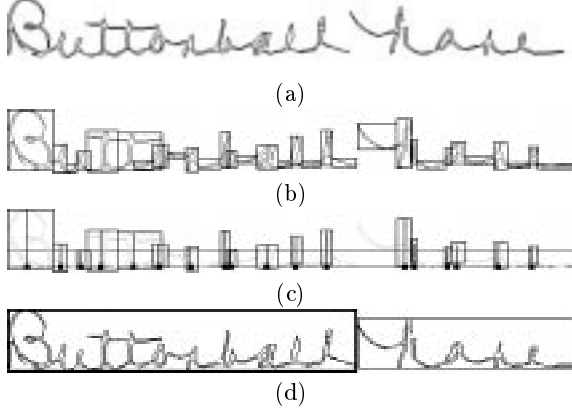


Figure 1: Word segmentation, (a) slant and skew normalized contours, (b) segmented primitives, (c) stroke interval estimation, (d) word segments

- The primitives are classified into two categories, *prime*: a basic stroke of characters and *auxiliary*: such as ligature, hat or punctuation mark.
- A character has at least one prime primitive.
- The distance between two adjacent prime primitives of a character is equal or less than the dominant prime primitive frequency.
- Word gaps are larger than the spatial period of prime frequency.

As an image representation, the conventional contour method is used. However for efficient resource usage, we use piecewise linearized contour model instead of a pixel based contour

2.1. Stroke Segmentation

Connected components are split into more detailed symbolic primitives. Stroke segmentation points are determined using features of convexity, concavity and ligature by tracing contours. Considering geometric relations among convexities, concavities, and ligatures, primary segmentation points are determined by a heuristic rule based algorithm. To extract more accurate geometric relations among strokes, the primary segmentation points are expanded along ligature strokes. These expansions can separate flourish and ligature strokes from basic primitives.

2.2. Stroke Primitive Classification

After stroke segmentation, each split stroke primitive is classified into *Prime*, *Hat*, *Tail* and *Ligature* in order to separate auxiliary strokes from prime strokes.

Two reference lines, *centerline* and *baseline* are used to establish the primary writing band. The merging ratio of bounding box to primary writing band is used as the reference to classify *Hats* and *Tails*. Vertical component run ratio, ligature stroke ratio and number of convexities and concavities throughout the trace of primitive contours are used as features to classify ligature primitives.

2.3. Word Break Detection

Assume M primitives are generated in a given phrase image and N primitives are classified as prime type. Let i be index of prime primitives in $0 \sim N - 2$ and (x_i, y_i) be the center point of mass of i^{th} prime primitive.

The interval of the i^{th} prime primitive d_i is $d_i = x_{i+1} - x_i$. Let the average and variance of intervals be \bar{d} and v^2 respectively. \bar{d} is the average of mixture of inter-stroke, inter-character and inter-word intervals. \bar{d} can not be directly taken as the period of prime frequency. Following the assumption of our model that inter-word gap is much larger than the inter-character gap and that the number of inter-word gaps is significantly smaller than inter-character and inter-stroke gaps, the prime interval samples can be refined by using \bar{d} and v .

Let $T(d)$ be a filtering function as

$$T(d) = \begin{cases} d & \bar{d} - \alpha v \leq d \leq \bar{d} + \alpha v \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

The prime period p is defined by the average value of filtered prime intervals as

$$p = \frac{1}{K} \sum_{i=0}^{N-2} T(d_i) \quad (2)$$

where K is the number of filtered prime intervals. Since handwriting does not guarantee uniform spacing among strokes or characters, error rate e is defined by variance of filtered prime intervals to support the confidence of the extracted prime period.

$$e^2 = \frac{1}{K} \sum_{i=0}^{N-2} (T(d_i) - p)^2 \quad (3)$$

Finally, a metric of word gap confidence c_i in i^{th} prime primitive is defined in terms of the *Maharanobis Distance* as

$$c_i = \begin{cases} 0 & d_i \leq p \\ \frac{d_i - p}{e} & \text{otherwise.} \end{cases} \quad (4)$$

Possible word break points are detected by simple thresholding after computation of all gap confidences.

3. PHRASE RECOGNITION

3.1. Prescreening

The number of characters in an image segment is estimated by counting the number of times the distance between prime primitives is larger than the spatial period of estimated prime frequency. We skip the prime primitives whose distance to the previous neighbor is equal or less than prime period from left to right. But the accuracy of the estimated number strongly depends on the writing style because of flourishes and non-uniform character size or spacing.

Error range is defined to capture the estimation accuracy. Let m be a total number of prime primitives and n

be an estimated number of character. The error range r is defined as

$$r = \frac{ne}{p} \quad (5)$$

The boundary of estimated number of characters can be

$$n_{max} = \begin{cases} n + \beta r & \text{if } n + \beta r < m \\ m & \text{otherwise.} \end{cases} \quad (6)$$

$$n_{min} = \begin{cases} n - \gamma r & \text{if } n - \gamma r > 1 \\ 1 & \text{otherwise.} \end{cases} \quad (7)$$

where β and γ are control parameters.

Using the estimated bound of number of characters and number of word segments, a subset of lexicons can be selected by the following conditions

1. the number of lexemes in a lexicon is smaller than the number of word segments.
2. total number of characters in a lexicon is equal or less than the sum of n_{max} .
3. total number of characters in a lexicon is equal or greater than the sum of n_{min} .
4. if the number of words is the same as the number of image segments, the number characters of each word is within the bound of n_{max} and n_{min} in the corresponding image segment.

3.2. Hypothesis Generation

If a candidate of word gap is significantly large than others, then we can be certain that the gap is a valid word break. As another restriction, if anchors are selected among candidate of word breaks, the matching complexity can be lower.

For the purpose of additional restriction to generate hypotheses, word break point candidates are categorized into two groups; *hard word point* and *soft word point*. The gap confidence is used for the classification. The combination of image segments crossing hard word points is not allowed as hypotheses. In Figure 2, thick lined box describes a valid word block and the thin lined box shows soft word block.

To explain the generation of hypotheses, let us assume that the inputs are the images shown in Figure 2 and raw lexicons (Table 2) in the application of street name recognition. To overcome missing or incorrect prefix and suffix, each lexicon entry is expanded in all conceivable ways. Also, we assume all additional information can be attached at the end of the street name.

Table 1 shows the boundary of estimated number of characters. All possible combination of image segments and lexeme are generated within the restriction of class of word break and boundary of estimated number of characters.

In the example, the lexeme "Steamboat" can not match with the image segment 0, because the number of characters is out of bound. Also since a valid word break point is in between the image segment 0 and 1, both segments can not be grouped to match "Steamboat". So the lexicons which containing "Steamboat" at the first position are rejected.

After generating lexeme hypotheses, phrase (lexicon) hypotheses are generated as a sequence of lexeme hypotheses index. A phrase can have multiple hypotheses sharing common lexeme hypotheses.



Figure 2: word segments

image segment index	0	1	2	3	4
no. of primitives m	4	12	5	6	3
esti. no. of char. n	3	9	5	5	2
error bound r	1	3	2	2	1
n_{max}	4	12	5	6	3
n_{min}	1	2	1	1	1

Table 1: Estimated Number of Characters

3.3. Hypothesis Verification

Because the hypotheses consist of word segments and a subset of possible lexemes, a lexicon driven word recognizer is preferred for hypotheses verification rather than character based word recognizer.

Dynamic programming based matching through character segments and character array of lexicons is well known technique to find the best match between a word image and lexicon [4].

Since the average of an individual character's matching score is used in the word recognition confidence value, mismatches between characters and character segments are compensated by other good matches. Therefore, a longer string has a better chance of matching.

Finally the best match between the entire group of image segments has an advantage over individual lexeme based matching in hypothesis verification. Each possible group of image segments is submitted to the word recognizer with eligible subsets of phrase strings. The word recognition scores are retained as the hypotheses confidence.

4. EXPERIMENT

The phrase recognition method described is applied to street name recognition. The street name images are collected from live mail pieces and the raw lexicons are obtained from postal databases. Since the testing set comes from live mail pieces, the size of lexicon of each images is not fixed - it ranges from 1 to 100. Among our test sets 9 % of images have additional unwanted segments such as apartment numbers (Table 3) in the input image.

Table 4 shows the word segmentation performance on the test set. Our word segmentation algorithm misses actual word segmentation point in about 2% of all images, maintaining a perfect word segmentation in 48% and over segmentation points in 31% of images.

The recognition performance of phrase mode recogni-

Pre	Name	Suffix
	LAKE	AVE
	STEAMBOAT	DR
W	PUTNAM	AVE

Table 2: Raw lexicon

Total Images	1342
Images with additional unwanted segments	116 (8.64%)
Raw Lexicons (average)	22.0
Expanded Lexicons (average)	77.2

Table 3: Experiment Parameters

Number of Images	1342	100%
Perfect Segmentation	643	47.9%
1 Over Segmentation	418	31.1%
2 Over Segmentation	181	13.5%
3 or More segmentation	70	5.2%
Fail(under-segmentation)	30	2.2%

Table 4: Segmentation Performance

tion using our segmentation method is compared to that of word mode recognition. The word mode recognition is a simplistic approach where a street name image is treated as a single word by ignoring word spacing.

Table 5 shows generated hypotheses matching complexity before hypotheses are submitted to the word recognizer. By the effect of prescreening and clever generation of hypotheses, the phrase method is superior.

Figure 3 shows correct ratio and error rate changes by adjusting threshold of confidence respectively in both word mode and phrase mode recognition. Phrase mode recognition achieves higher correct ratio maintaining lower or same error rate. Figure 4 shows the reject-error curves.

5. CONCLUSION

In this paper, we describe an efficient method of phrase recognition in unconstrained environment - lexicons can provide only a subset of the true string and a phrase image can have additional stings associated with the core string. Prime stroke analysis is performed to get word break points and estimate the number of characters in a given phrase image. Over-segmentation followed by dynamic matching between one or more group of image segments and single lexeme is taken as a phrase recognition strategy. Experiments show enhanced recognition performance of our method. However, because of poor accuracy in estimating the number of characters in image segments, relatively large number of hypotheses are generated. More accurate estimation method and hypotheses generation method based on features and choosing optimal parameters which maximize system performance remains as future work.

Accepted Lexicons	62.0
Lexeme Hypotheses	59.2
Phrase Hypotheses	98.5

Table 5: Generated Hypotheses - All numbers are average values

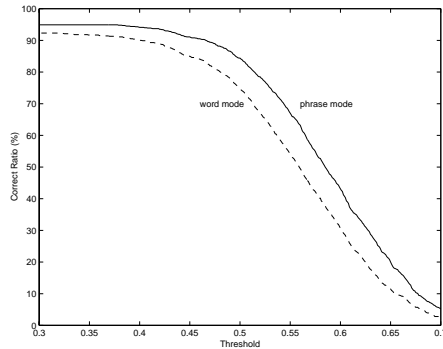


Figure 3: Correct Ratio Comparison

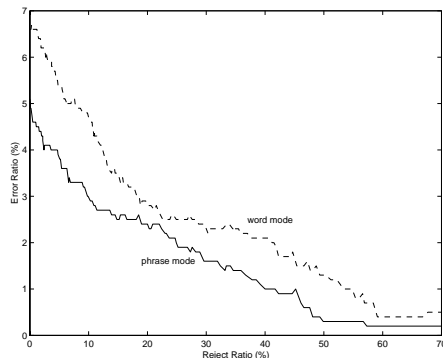


Figure 4: Reject to Error Characteristic Comparison

6. REFERENCES

- [1] G. Seni and E. Cohen. External word segmentation of off-line handwritten lines. *Pattern Recognition*, 27(1):44–52, 1994.
- [2] U. Mahadevan and R. C. Nagabushnam. Gap matrices for word separation in handwritten lines. *Int. Conf. on Document Analysis and Recognition, Montreal, Canada*, pages 124–127, 1995.
- [3] G. Kim and V. Govindaraju. Handwritten phrase recognition as applied to street name images. *Pattern Recognition*, 31(1):41–51, 1998.
- [4] G. Kim and V. Govindaraju. A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 19(4):366–379, 1997.