

특징점을 이용한 필기체 한글인식에 관한 연구

권중강 · 이정천 · 김성태 · 김재균

한국과학기술원 · 전기및 전자공학과

Received on October 24, 1986, in final form on November 26, 1986

Accepted for publication on December 10, 1986

Journal of the Korean Institute of Electrical Engineers, Vol. 34, No. 11, pp. 863-866, 1986

Keywords: Feature point, Handwritten Korean character recognition

요약

In this paper, we propose a new method for recognizing handwritten Korean characters. The method is based on the feature points of the characters. The feature points are defined as the points which are located at the corners of the characters. The feature points are used to identify the characters. The method is simple and efficient. It can be used for the recognition of handwritten Korean characters.

1. 서론

최근 막대한 양의 문자정보를 컴퓨터에 자동 입력시키기 위해서는 효과적인 문자인식 과정이 필요하다. 이의 일환으로 국내에서 한글인식에 관한 연구가 진행되고있다. 그러나 대부분이 인쇄체 중심으로 연구되어 왔으며(1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100), 필기체의 경우는 최근에 들어 연구되고있다. 필기체는 영문과와 달라서 직선성이 적어, 따라서 효과적인 한글인식을 위하여 한글을 선 혹은 선의 선이 만나서 생기는 꼭지점들로 표시할 필요가 있다. 이와 같은 특징을 쉽게 추출하기 위한 일반적인 방법으로는 문자 영상의 이진화, 이진화 후 임계값을 가진 선의 선화 과정이 필요하다. 그러나 선화 과정을 받은 선의 선소요하며 경우에 따라서는 선화 생길 수 있는 불필요한 꼭지점들을 제거하는 과정이 필요하다. 이 필기체 인식 부분자료의 변형이 심하며, 날소가 많은 경우 분리하는 과정이 필요하다.

본 논문에서는 선화 과정을 거치지 않고 문자영상의 경계를 이용하여 선으로 부호화하여 경계의 점인 부분을 추출하고 이를 효과적으로 표현하는 방법인 굵은 글자의 분리방법 그리고 이를 이용한 필기체 한글인식 알고리즘을 제안한다. 본 알고리즘의 전체 흐름도는 그림 1과 같다.

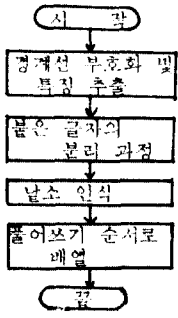


그림 1 한글 인식을 위한 전체적인 흐름도

2. 한글글자의 특징점 추출과 표현

문자영상의 경계선이 보이는 부분의 특징점이 되고, 특징점의 선은 경계선의 굵은 선의 방향에 따라 다음과 같이 4가지로 정의하며 그림 1을 보라.

- 1. 내부 경계에 점선부분의 특징점
- 2. 내부 경계에 굵은 선의 끝부분에 점인 부분의 특징점
- 3. 외부 경계에 굵은 선의 끝부분에 점인 부분의 특징점
- 4. 외부 경계 선 꼭지점을 가리는 꼭지점

한글의 인지는 내부 경계의 내부 경계선의 굵은 선으로 인식될 수 있고, 외부 경계의 내부 경계의 꼭지점들의 인식으로 보일 수 있다. 꼭지점들은 그 꼭지점들의 위치로 구분된다.

$$C = \langle E, H \rangle$$

$$E = \{E_i, i=1, \dots, N_E\}$$

$$H = \{H_j, j=1, \dots, N_H\}$$

$$E_i = [EV_{ik}, k=1, \dots, N_V(E_i)]$$

$$H_j = [HV_{jm}, m=1, \dots, N_V(H_j)]$$

$$EV_{ik} = \langle \text{type}, \text{position} \rangle$$

$$HV_{jm} = \langle \text{type}, \text{position} \rangle$$

한글의 글자 내부 경계의 굵은 선
내부 경계의 굵은 선 외부 경계의 굵은 선
외부 경계의 굵은 선

EV_{ik} : i번째 내부 경계의 k번째 꼭지점
 HV_{jm} : j번째 외부 경계의 m번째 꼭지점

예를 들면 '삼'의 굵은 선의 경우, 내부 경계의 내부 경계의 내부 경계로 인식된 선의 끝부분의 내부 경계의 인식은 그림 2와 같다.

3. 낱소인식 및 분리 방법

한글은 필기자의 습관에 따라 글자의 변형이 심하며, 낱소의 상태에 따라 크게 세가지로 분류할 수 있다.

- 1) 글자의 낱소들이 모두 떨어져 있는 경우
- 2) 글자의 낱소가 붙어있는 경우
- 3) 글자의 낱소의 획이 모두 떨어져 있는 경우

특히 2)와 3)의 경우는 분리하는 과정과 연결하는 과정이 필요하므로 인식하는데 어려움이 있다. 본 논문에서 제안하는 방법은 1), 2)의 경우에만 적용된다.

3.1 낱소인식

위에서 분류한 3가지 중 1)의 경우는 꼭지점의 개수, 꼭지점의 type, type 포함 여부 등을 저장한 look-up table(LUT)를 통과하면 쉽게 인식할 수 있다. 본 논문에서는 LUT에 포함여부와 꼭지점 type을 이용하여 LUT와 같이 5개의 bin으로 분류하고, 입력 문자 영상이 위의 5개의 bin에 속하지 않으면 붙은 경우로 간주하고 인식을 한다. 각 bin의 자소들은 꼭지점의 상대적 위치와 방향을 이용하여 인식할 수 있다.

각 bin의 자모	N_b	N_0	N_1	N_2	N_3
L0 ---	0	0	0	2	0
L1 ㄱ, ㅋ, ㆁ, ㄷ, ㅌ, ㄴ	0	2	0	0	1
L2 ㅈ, ㅊ, ㅊ, ㅊ, ㅊ	0	4	0	4	2
L3 ㅊ, ㅊ	1	4	0	4	0
L4 ㅊ	0	4	0	4	0
L5 ㅊ, ㅊ, ㅊ	1	x	>0	x	x
L6 ㅈ, ㅊ, ㅊ, ㅊ, ㅊ	0	>0	>0	2	0
L7 ㅈ, ㅊ, ㅊ, ㅊ, ㅊ	0	>0	>0	>2	>0

Table 1. 한글 기본자모

N_b : bin의 수

N_0, N_1, N_2, N_3 : 각각 type이 '0', '1', '2', '3'인 꼭지점의 수

x: don't care

3.2 분리과정

한글의 구조상 한글자모가 붙는 경우는 다음 5가지 경우로 분류한다.

- 1) 꼭지점과 꼭지점이 아닌 부분이 붙는 경우
- 2) 꼭지점과 꼭지점이 붙으면서 일부부분 변형된 형태로 남아 있고 나머지는 사라지는 경우
- 3) 꼭지점과 꼭지점이 붙으면서 변형된 형태로 남아 있는 경우
- 4) 꼭지점과 꼭지점이 아닌 부분이 붙으면서 LUT에 생기는 경우
- 5) 꼭지점과 꼭지점이 붙으면서 꼭지점이 모두 사라지는 경우

위의 경우들의 예는 그림 4에 나타나 있다. 그리고 본 논문에서는 1)~5)까지만 고려한다.

낱소 인식과정에서 1)~5)의 아니면 다음과 같이 4가지 경우를 차례로 조사하여 붙은 낱소를 분리한다.

1)의 경우: 이와 같은 경우는 꼭지점과 꼭지점이 아닌 부분이 붙어서 생성된 꼭지점의 개수는 항상 '0'이 되어야 한다. 즉, 이 경우인 모든 꼭지점을 찾아서 이들중 가장 인접한 두 꼭지점 사이의 우선순위를 값을 식으로 부터 계산하여 고 값이 아닌 가장 큰 값을 갖는 꼭지점 쌍을 분리 하면된다. 그리고 이 경우에 따라서 LUT와 LUT의 중간 경계점에 있는 인접한 꼭지점을 삼입한다. 여기서 LUT는 인접한 두 꼭지점중 먼저 추출된 꼭지점이고 LUT는 나중에 추출된 꼭지점이다.

2)의 경우: 이 경우와 같은 경우는 LUT에 속하는 경우

3)의 경우: 이 경우와 같은 경우는 LUT에 속하는 경우

4)의 경우: 이 경우와 같은 경우는 LUT에 속하는 경우

5)의 경우: 이 경우와 같은 경우는 LUT에 속하는 경우

3.3 풀어쓰기 순서로 배열

LUT를 사용하지 않았기 때문에 자모의 순서가 임의로 배열될 수 있다. 이를 풀어쓰기 순서로 배열하기 위하여 낱소들의 중심좌표 좌하의 상관관계를 이용한다. 낱소가 좌우이면 중심좌표의 변형은 수평 관계, 수직관계, 그리고 대각선 관계로 나눌 수 있다. 먼저 낱소들을 좌표의 상관 순서로 배열한 뒤 수평관계일 때는 좌표가 수직관계일 때는 좌표가 같은 낱소가 먼저오도록 배열한다. 대각선 관계일 때는 좌표의 값이 연속된 두 개의 낱소 중 앞의 낱소의 좌표가 뒤의 낱소 좌표보다 크면 낱소의 순서를 바꾼다.

좌/우	1	2	3	4
상/하	1	2	3	4
대/소	1	2	3	4
대/소	1	2	3	4

이후에 이 대각선 관계의 낱소들의 좌표순서

마지막으로 'ㅇ'과 'ㅇ'을 'ㅎ'으로 'ㅈ'과 'ㅈ'을 'ㅈ'으로 알아보게 하는 등의 한글글자 단위로 조합하는 과정을 거친다. 예를들면 '홍'이 'ㅇ'과 'ㅇ'과 'ㅇ'에서 'ㅎ'과 'ㅇ'과 'ㅇ'으로 알아보는 과정을 거친다.

5. 실험 및 결과

제안한 방법의 유용성을 조사하기 위하여 15명의 사람들에게 100자씩 글을 적게하여 실험하였다. computer는 IBM 3600을 사용하였고 10 characters로 20 groups을 작성하였다. 시간은 글자의 크기 즉 character code의 길이에 비례한다. 오인식의 대부분은 한글글자의 변형으로 문장의 내용을 보고 추정할수 있는 경우와 필기습관에 의해 꼭지점이 추출되지 않는 경우 그리고 'ㅁ'과 'ㅇ'의 구별이 힘든 경우였다.

6. 결론

본 논문에서는 세션화 과정을 거치지 않고, 또 computer를 사용하지 않는 방법으로 시간을 줄이고 개인의 습관에 의한 변형된 글자에 대해서도 높은 인식률을 보이는 방법을 제안하였다. 보다 인식률을 높이기 위하여 꼭지점을 정확히 추출하는 방법에 대한 연구가 필요하다 그리고 한글인식의 실용화를 위해서는 글자의 획이 떨어져 있는 경우와 문맥에 의한 글자의 추정이 가능하도록 하는 연구가 있어야 하겠다.

6. 참고문헌

- [1] 김태곤, "Synthetic 법에 의한 한글의 패턴인식에 관한 연구", 전자공학회지 vol. 14, NO. 2, pp. 10-16, 1977.12.
- [2] 최병욱, "한글인식에 있어서의 자소추출", 전자공학회지 vol. 10, NO. 2, pp. 30-43, 1973.12.
- [3] 이주근, 남궁재찬, 김영진, "한글 인식을 위한 Character분리와 인식에 관한 연구", 전자공학회지 vol. 10, NO. 3, pp. 11-18, 1973.6.
- [4] 박종욱, 이주근, "Machine Recognition에 의한 필리핀의 한글 인식", 전자공학회지 vol. 20, NO. 2, pp. 1-9, 1980.9.
- [5] 이정현, "Character Coding Code를 이용한 새로운 Character Coding 알고리즘에 관한 연구", 한국과학기술논문집 제14권 1980.

그림 1. 추출된 꼭지점들

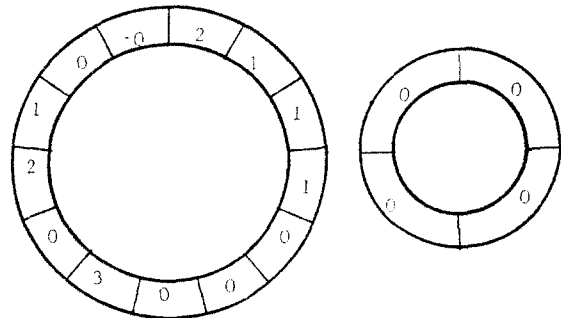


그림 2. 'ㅁ'자의 외부경계와 내부경계의 원순열

직갈갈며밥

- 1) 2) 3) 4) 5)

그림 3. 한글이 없는 경우의 예

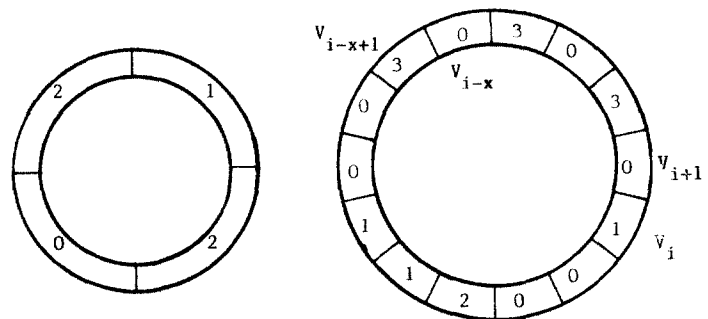
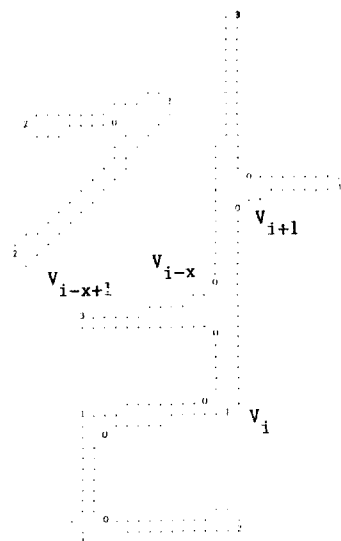


그림 4. 3)의 경우를 설명하기 위한 예