

A Study on Grapheme and Grapheme Recognition Using Connected Components Grapheme for Machine-Printed Korean Character Recognition

Kyong-Ho Lee *

Abstract

Recognition of grapheme is a very important process in the recognition within 'Hangul(Korean written language)' letters using phoneme recognition. It is because the success or failure in the recognition of phoneme greatly affects the recognition of letters. For this reason, it is reported that separation of phonemes is the biggest difficulty in the phoneme recognition study.

The current study separates and suggests the new phonemes that used the connective elements that are helpful for dividing phonemes, recommends the features for recognition of such suggested phonemes, databases this, and carried out a set of experiments of recognizing phonemes using the suggested features.

The current study used 350 letters in the experiment of phoneme separation and recognition. In this particular kind of letters, there were 1,125 phonemes suggested. In the phoneme separation experiment, the phonemes were divided in the rate of 100%, and the phoneme recognition experiment showed the recognition rate of 98% in recognizing only 14 phonemes into different ones.

▶ Keyword : character segmentation, grapheme extraction, grapheme recognition

1. Introduction

문자 인식은 지난 수십 년 동안 많이 연구되어 온 분야이나, 인식을 향상 위해 지금도 여전히 많은 연구들이 나오고 있다. 오현아 등은 문자 인식 향상을 위해 대상 문자열에 회전 알고리즘을 적용하여 인식을 향상시키고 있으며[1], 진문용 등은 인식을 향상 위해 휴리스틱 분할 알고리즘을 적용하여 향상을 꾀하고 있다[2]. 김강산 등은 구글의 Open API를 이용하여 한글 문자 인식 향상을 위한 연구를 하였으며[3], 정규수는 한글 문자 주변에 심벌이 위치한 상황에서 인식을 향상 위해 한글 문자 템플릿에 의존하지 않는 점진적 좌측 방향으로의 블러 투사를 제시하였다[4]. 박경화 등은 한글 문자의 다양성 확대를 위해 다중 컬럼 딥 인공 신경망을 이용하여 한글 필기체 인식을 시도하는[5] 등 아직도 많은 연구들이 지속되고 있다. 이렇게 연구가 지속되는 이유는 한글 문자의 특수성 때문으로 한글은 영문자 알파벳과 달리 19개의 초성과 21개의 중성, 27

개의 중성들이 2차원적 조합으로 되어 11,172개의 문자가 구성되기 때문에 인식해야할 대상이 많으며, 자소 간 구분이 작은 획 하나에 의해 구분될 정도로 유사성이 높아 인식에 대한 처리 과정이 알파벳에 비해 복잡하기 때문이다[6,7]. 또한 아직도 영문자 알파벳 등 보다 인식이 낮기 때문에 많은 연구가 지속되고 있다.

그동안 연구된 한글 인식 방법들은 한글을 작성하는 방법에 따라 '인쇄체 인식'과 '필기체 인식'으로 분류할 수 있고 또 인식하는 단위에 따라 '문자 단위 인식'과 '자소 단위 인식'으로 분류할 수 있다. 문자 단위 인식은 문자 단위를 통으로 인식 대상으로 보는 것이며, 자소 단위 인식은 문자를 구성하고 있는 자소를 분리해 자소 단위로 인식을 하는 것으로 초성, 중성, 종성을 이루고 있는 소수의 자소 단위를 인식한 후 인식한 자소 정보와 문자를 구성하는 형식 정보를 이용하여 글자를 인식하는 방법이다. 두 가지 다 전처리 작업으로 입력된 영상에서 먼저 기울어짐 보정을 하고, 텍스트 영역 추출과 수평 분할, 수직

*First Author: Kyong-Ho Lee, Corresponding Author: Kyong-Ho Lee

*Kyong-Ho Lee(khlee@halla.ac.kr), School of Information & Communication, Broadcasting Engineering, Halla University

*Received: 2016. 08. 09, Revised: 2016. 08. 30, Accepted: 2016. 09. 20.

분할 등의 반복을 통해 문자 단위로 분할 한 후, 문자 단위에서는 각각의 방법으로 인식을 한다. 또한 대부분의 경우 추가되는 전 처리 작업으로 추출해 낸 문자의 크기를 자신들이 정한 규격으로 조절한 후 인식을 한다. 문자 단위 인식의 경우 정해진 규격 크기를 웨이블릿 변환을 하여 인식하거나, 신경망을 통한 인식에서는 문자를 구성하는 화소를 입력으로 하여 인식을 한다[8-11]. 그런데 현재 한글에서 구성할 수 있는 문자의 종류는 11,172자로 문자 단위 인식에서는 인식 대상 문자 수에 따라 인식률과 인식 속도, 메모리 요구량이 크게 달라지는 것으로 알려져 있다. 관련 논문들을 보면 실험에 수만 자를 인식한 경우 우라도 인식 대상에 중복되는 문자들 때문에 수만 자가 되는 것이고, 문자 별로 분류해 보면 그 가짓수는 조합형 11,172자에 터무니없이 적으며, 사용 빈도가 높은 완성형 2,350자에도 턱없이 부족한 수이다.

자소 단위 인식은 문자에서 자소를 분리하여 분리한 일정한 수의 자소만을 인식하므로 인식해야할 가지 수가 적어 성능의 변화가 상대적으로 적다고 알려져 있다. 물론 자소를 인식한 후에는 문자 구성 형식 정보, 오토마타, 신경망 등과 같은 적당한 방법으로 문자를 인식한다. 자소 단위 인식에서 자소의 분할 연산은 매우 중요한 과정이다. 그림 1은 현재 한글에서 사용하는 자소들을 사용 위치에 따라 분류한 것이다. 자소 단위 인식에서 자소 분리를 하다가 자소의 일부분이 소실될 수도 있고, 자소 구성 영상 이외의 잡영이 포함 될 수도 있다. 자소의 범위를 작게 잡으면 자소의 일부가 소실되고 소실을 막기 위하여 자소의 영역을 크게 잡으면 잡영이 들어오는데 소실과 잡영을 처리하는데 애를 먹기도 한다[7].

구분	내 용	갯 수
초성	ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅇ ㅈ ㅊ ㅋ ㅌ ㅍ ㅎ ㄱ ㄷ ㅂ ㅅ ㅈ	19
중성	ㅏ ㅑ ㅓ ㅕ ㅗ ㅛ ㅜ ㅠ ㅡ ㅣ ㅖ ㅗ ㅛ ㅜ ㅠ ㅖ ㅗ ㅛ ㅜ ㅠ ㅡ ㅣ	21
종성	ㄱ ㄴ ㄷ ㄹ ㅁ ㅂ ㅅ ㅇ ㅈ ㅊ ㅋ ㅌ ㅍ ㅎ ㄱ ㅅ ㄴ ㅈ ㄹ ㅂ ㅅ ㄴ ㄷ ㄹ ㅂ ㅅ ㄴ ㄷ ㄹ ㅂ ㅅ	27

Fig. 1. Korean consonants and vowels according to the location

자소 분할에 애를 먹는 가장 큰 이유로는 한국어 문법에서 분류하고 있는 초성, 중성, 종성의 개념에서 벗어나지 못하고, 그 개념을 가지고 구성 형식 정보로 추정된 크기를 자소 추출 원도우로 적용하여 인식하려는 방법이 자소 영역을 잘못 잡아 자소 정보 소실 또는 잡영 추가 문제를 일으킨다. 자소 분할 인식에서는 대부분의 경우 한글을 아래와 같이 6개의 패턴으로 이루어 졌다고 보고 분할을 시도하고 있다[7]. 그러나 수많은 글자들은 인쇄하였을 때 자소가 서로 붙어 인쇄되고 있어 위 패턴에 맞추어 기계적으로 분류하기가 매우 어렵다.

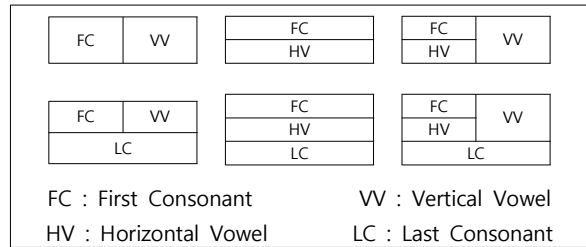


Fig. 2. Korean character syntactic form

맑은 고딕으로 많이 쓰이는 문자로 판명된 완성형 2350자를 인쇄하여 분석해 보면 649자(약 28%)의 문자가 자소가 서로 붙어 있어 자소별 기계적 분할이 어려울 것으로 사료되며, 앞에서 지적한 바와 같이 자소를 분류하기 위한 범위 설정 크기가 알맞지 않으면 인식해야할 자소 일부의 손실과 원치 않는 잡영의 유입 문제가 발생되어 자소 인식률에 많은 영향을 준다.

컴퓨터와 같은 기계를 이용한 문자 인식 자동화는 사람이 발화하는 음운을 기본으로 하는 초성, 중성, 종성 틀이 아닌 기계를 중심의 새로운 자소 인식 방법이 제공되어야 한다. 본 연구에서는 인쇄체 문자에서 자동 자소 단위 한글 인식에 쉬운 자소를 제안하고, 자동 자소 인식 방법을 제안하였다. 본 연구의 가치는 한국어 연구에서 음운론을 기본으로 하는 초성, 중성, 종성의 관점에서의 기존의 연구 방법과 다르게 자동 자소 인식을 유리하게 하기 위한 노력으로 한국어 음운학에서의 기본 자소 뿐 아니라 이들이 결합된 다양한 복합 자소도 새로운 자소로 제안하였고, 이런 자소들을 자동으로 인식하는 방법을 제안하였다. 본 연구에서 사용한 인쇄체는 맑은 고딕으로 하였으며, 본 연구에서 새로운 자소 추출과 새로운 자소에 인식을 위해 분석한 글자는 조합형 11,172자 모든 글자를 대상으로 하지 않고 많이 쓰이는 빈도로 구성된 완성형 2,350자를 인쇄 출력하고 이 인쇄물로 연구를 수행하였다.

II. 본 론

2.1 연구개요

본 연구는 차량 운행 시 안전과 원활한 소통을 위하여 차도 표면에 글자를 쓰는 작업인 도로 노면 문자 도색[13]을 위한 자동화 장비를 만드는 연구 과정에서 도로노면 문자도색을 위한 정자체의 인쇄 한글 인식의 필요로 본 연구를 추가 수행하였다. 정부 관련 기관의 관련 자료를 분석하여 규정상 문제없는 문자체로 인식 대상 문자체를 맑은고딕으로 선택하고 이로 연구를 수행하였다. 최종 목표는 자소 단위 인식 방법을 통한 문자 인식이다. 자소 인식 방법을 선택한 이유는 통 문자 인식 방법은 인식 대상 수가 많아지면 인식률은 떨어지고, 인식 속도는 느려지며, 메모리 요구량 매우 커지기 때문이다.

자소 단위 문자 인식에서 자소의 분할 연산은 매우 중요한 과정이다. 기존 연구에서는 사람이 발화하는 음운을 기본으로

의 종류, 글자 최인접 사각형 크기 정보와 이 안에서의 3x3 위치 정보가 저장되었다.

0 0 0 0 1 1 0 0 0	0 0 1 0 1 0 0 0 0	0 1 0 0 1 0 0 0 0	1 0 0 0 1 0 0 0 0
0 0 0 1 1 0 0 0 0	0 0 0 0 1 0 1 0 0	0 0 0 0 1 0 0 1 0	0 0 0 0 1 0 0 0 1

Fig. 11. Mask for end point

2선 모인 점은 위와 같은 패턴을 5x5 행렬로 구성하되 중심 점은 무관조건으로 하고, 2패턴의 or 결합 형태로 2선 모임 점을 추출하되 패턴들을 동쪽 방향에서 부터 반시계 원형으로 나열하였을 경우 1 또는 2칸 거리의 패턴들의 결합만 허용하였다. 또한 끝점과 같은 형태의 정보들이 추출되어 저장되었다.

3선 이상 모임 점은 아래 패턴을 번호순으로 원형 나열하고, 16패턴 중 각각이 양 방향으로 거리가 3이상 차이 나는 3개 또는 4개의 패턴을 or 연산으로 결합하여 3선 모임 점과 4선 모임 점을 찾았고 앞에서와 같이 관련 정보를 저장하였다.

(1)	(2)	(3)	(4)
(5)	(6)	(7)	(8)
(9)	(10)	(11)	(12)
(13)	(14)	(15)	(16)

Fig. 12. Mask for three lines or more meeting point

요소	구역별 특징 점 분포표			
	끝점	2선 모임	3선 이상	고립점
o				1
□		1 1		1
┌	1		1	
L	1			

ㅂ	1 1							
ㅅ		1 1						
ㅆ	1 1				1 1			
ㅈ			1 1					
ㅊ				1 1				
ㅋ	1 1				1 1			
ㆁ							1	
ㆂ	1 1				1 1			
ㆃ								1
ㆄ								
ㆅ								
ㆆ								
ㆇ								
ㆈ								
ㆉ								
ㆊ								
ㆋ								
ㆌ								
ㆍ								
ㆎ								
㆏								
㆐								
㆑								
㆒								
㆓								
㆔								
㆕								
㆖								
㆗								
㆘								
㆙								
㆚								
㆛								
㆜								
㆝								
㆞								
㆟								
ㆠ								
ㆡ								
ㆢ								
ㆣ								
ㆤ								
ㆥ								
ㆦ								
ㆧ								
ㆨ								
ㆩ								
ㆪ								
ㆫ								
ㆬ								
ㆭ								
ㆮ								
ㆯ								
ㆰ								
ㆱ								
ㆲ								
ㆳ								
ㆴ								
ㆵ								
ㆶ								
ㆷ								
ㆸ								
ㆹ								
ㆺ								
ㆻ								
ㆼ								
ㆽ								
ㆾ								
ㆿ								

토			2			1	1													
			1							1										
	1		1								1									
토			2	1																1
			1	1																
	1		1																	
관			2																	
	1		1																	
	1	2		1																
			1	1																
하			1																	
	1		1																	
	1																			
		1	1	1	1	1	1													
			1																	
	2		1																	
			1	2																
대			1																	
	1	1	1																	
	1		1																	
	1																			
사																				
	1		1																	
	1	1	1																	
쌍 (0)																				
	1		1																	
	1	1	1																	
국																				
	2																			
	2		1																	
			1																	
규																				
	2																			
	1		1																	
			1	1																
규																				
	1	1																		
	1		1																	
	1		1																	
분																				
	1	2	1																	
	1		1																	
부																				
	1	1																		
	1		1																	
	1	1	1																	
분																				
	1	2	1																	
	1																			
방																				
	1																			
	2		1																	
			1	1	2	1														
배																				
	2																			
	2		1																	
			1	2																
방																				
	1	1																		
	2		1																	
			1	2																
표 (L)																				
	1		1																	
	1		1																	
	1		1																	
쌍																				
	1		1																	
	1	2	1																	
표 (L)																				
	1		1																	
	1		1																	
	1		1																	
꽃																				
	1		1																	
	1		1																	
	1		1																	
꽃																				
	1		1																	
	1		2																	
			1	1																
			1	1																

L:하단부분 I:고립점있음.

Fig. 13. Per grapheme features distribution

저장된 자소별 추출 특징 점은 앞의 설명과 같이 저장하되 끝 점, 2선 모임 점, 3선 이상 모임 점과 고립점은 자소의 가로와 세로를 3대 3으로 구역으로 나누었을 때 각 추출된 위치에 따라 그림 13과 같이 저장하여 인식에 이용하였다. 이렇게 저장하였을 때 저장 모습이 자소의 모습을 유추할 수 있었다.

2.5 자소 인식 실험

본 연구에서 천연색 영상을 회색 영상으로 바꾸는 작업은 [13]을 이용하였고, 회색 영상을 2진화 영상으로 바꿀 때 경계값은 127로 하였고, 기울어짐 보정과 텍스트 영역 추출, 문자 분리는 [11]을 참고하였다. 글자의 크기는 20포인트 글자 크기로 가정하여 글자의 높이가 100픽셀이 되도록 영상 크기를 조정하였다. 분리된 문자에서 자소 추출은 자소 범위 내에서 연결 성분을 이용한 추출로 수행하였고, 앞의 그림 6에서 보인 바와 같이 글자에서 자소는 완벽히 분리 추출됨을 확인할 수 있었다.



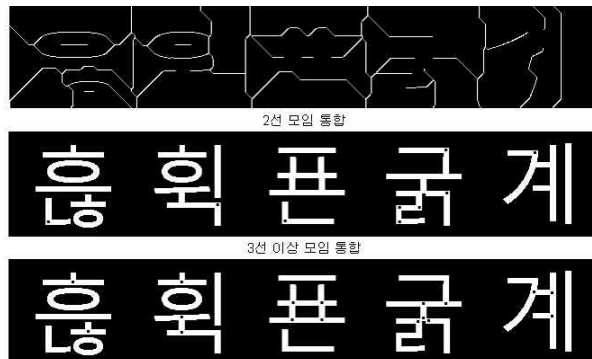


Fig. 14. Grapheme recognition test image

그림 14는 인식 실험을 위한 변환 과정과 특징점 추출과정이다. 자소 인식 실험에는 조합형 11,172개 중 난수 발생을 통해 무작위로 문자를 추출하도록 하여 350개 글자로 실험을 수행하였고 350개 글자에는 중 자소 1,125개가 있었으며, 실험에서는 자소 1125개 중 8개의 특징 점을 찾지 못하였고, 3개의 잘못된 특징 점을 추출하였다. 잘못된 특징 점을 찾는 경우는 오른쪽 빠침 부분이 직선으로 세선화 되지 못하고 한 번 절곡되어 세선화 될 경우에 나타났으며, 찾지 못한 특징 점은 2선 모임점이 부드럽게 세선화 될 경우 발생하였다.

자소 인식은 그림 13에서 보인 특징 점 정보를 데이터베이스화하여 끝 점의 수, 2선 모임 점의 수, 3선 이상 모임 점의 수, 고립점의 개 수 정보로 1차 분류한 후 각 특징 점의 분포를 이용하여 2차 분류하였다. 데이터베이스의 저장 값에 일치하지 않을 경우 일치하지 않은 특징 점과 데이터베이스의 특징 점의 상대 거리가 작은 쪽의 자소로 판단되도록 하였고, 실험에서는 1,125개의 자소 중 14개를 다른 자소로 판단하였다.

앞의 상황을 참고하건데 세선화 과정 후 글자 구조를 감안한 적절한 세선화 글자의 스무딩 과정을 고안하거나 또 다른 연구에서와 같이 자소 분리를 한다면 세선 후 자소 분리 과정을 연구해 보면 어떤 방법이 도출 될 것으로 사료된다.

III. Conclusions

본 연구에서는 문자 인식 방법 중에서 자소 인식을 이용한 문자 인식 방법에서의 자소 인식의 어려움을 해소하기 위한 새로운 자소 인식에 관한 연구를 수행하였다.

기존 자소 인식을 이용한 문자 인식 연구에서는 자소 인식에 한국어 음운학에서 분리한 초성, 중성, 종성의 자소들을 인식 단위로 하여 인식을 하려고 시도한 것과 다르게 본 연구에서는 인쇄체 문자의 연결 성분을 그대로 인정한 새로운 자소를 구성 하였고 제안하여 자소 분리에 매우 유리함을 보였다. 또한 이 자소들을 인식하기 위한 특징 점들을 한글의 구조를 감안하여 끝 점, 2선 모임 점, 3선 이상 모임 점, 고립점으로 제안하고 각

특징 점들을 추출하는 방법을 제시하였다. 또한 특징 점의 추출 위치도 함께 저장하고 본 연구에서 제안한 자소들을 인식하는 실험을 수행하였다. 본 연구에서는 비록 높은 인식률을 얻기는 하였으나 제안한 방법의 알고리즘 수행시간을 줄이려는 노력과 오인식과 미인식 특징 점들을 위한 개선 노력이 필요하다.

또한 본 연구의 후속 연구로는 본 연구에서의 인식한 자소들 오 문자를 인식하기 위한 연구가 필요하다.

REFERENCES

- [1] Hyuna Oh, EuGene Rhee, "Enhancement of Car Licence Plate and Security with Rotation Algorithm", Journal of Security Engineering, Vol. 13, No. 2, pp. 83~90, Apr. 2016
- [2] Moon Yong Jin, Jong Bin Park, Dong Suk Lee, "Real-Time Vehicle Licence Plate Recognition System Using Adaptive Heuristic Segmentation Algorithm", KIPS Tr. Software and Data Eng., Vol. 3, No. 9, pp. 361-368, Mar. 2014
- [3] Kang-San Kim, Seok-Cheon Park, Seok-Ho Oh, "Suggestion of Enhanced Korean Character Recognition Technique Using Google Tesseract Open API", Proceeding of Korean Society for Internet Information, Vol. 16, No. 1, Spring. 2015.
- [4] Kyusoo Chung, "Text Area Detection of Road Sign Images based on IRBP Method", Journal of Intelligent Transportation System, Vol. 13, No. 6, Dec. 2014.
- [5] Kyung-Wha Park, Byoung-Hee Kim, Dong-Sig Han, Seong-Ho Son, Woo-Yung Kang, Byoung-Tak Zhang, "Handwritten Hangeul Recognition using Multi-column Deep Neural Networks", Proceeding of KIPS Spring Conference, Vol. 26, No. 1, 2016.
- [6] Min-Soo Kim, Eun-Young Kang, Woo-Sung Kim, Sun-Hwa Han, Jin-Hyung Kim, "A Study on Implementation of Printed Character Recognition System And Performance Evaluation", Korea Information Processing Society, Vol. 7, No. 11, pp. 3584-3591, Nov. 2000.
- [7] Kil Taek Lim, Ho Yon Kim, "A Study on Machine Printed Character Recognition Based on Character Type Classification", The Institute of Electronics and Information Engineers, Vol. 40, No. 5, pp. 26-39, Sep. 2003.
- [8] Kil-Taek Lim, Gi-Seok Kim, "Reestimation of

- Recognition Result of MLP Classifier for Machine Printed Hangul - Feasibility Study on Softmax Method", Journal of Information & Electronic Technology, Vol. 6, pp. 93-105, 2007.
- [9] Kil Taek Lim, Ho Yon Kim, "A Study on Machine Printed Character Recognition Based on Character Type Classification", Journal of IEIE, Vol. 40, No. 5, pp. 266-279, Sep. 2003.
- [10] Kil Taek Lim, Seon Hwa Jeong, Seung Ick Jang, Ho Yon Kim, "An Implementation Method of the Character Recognizer for the Sorting Rate Improvement of an Automatic Postal Envelope Sorting Machine", Journal of Korea Society of Industrial Information System, Vol. 12, No. 4, Dec. 2007.
- [11] Duk-Ryong Lee, Woo-Youn Kim, Il-Seok Oh, "A Hangul Document Image Retrieval System Using Rank-based Recognition", Journal of The Korea Contents Association, Vol. 5, 2005.
- [12] Kyong-Ho Lee, Jae-Joon Seong, "Study on Automation about Painting the Letters to Road Surface", Proceeding of Korea Society of Computer Information (Summer), Vol. 24, No. 1, pp. 113-116, Summer 2016.
- [13] A. A. Gooch, S. C. Olsen, J. Tumblin, and B. Gooch, "Color2gray: Saliency-preserving color removal", ACM Transactions on Graphics, vol. 24, no. 3, pp. 634 - 639, 2005.

Authors



Kyong-Ho Lee received the B.S. degree in Computer Science from Korea National Open University and the M.S. degree in Information and Communication Engineering from Korea Advanced Institute of Science and Technology and the Ph.D. degrees in Electronic Engineering from Dankook University, Korea, in 1991, 1994 and 2008, respectively. Dr. Lee is currently a Professor in the School of Information & Communication, Broadcasting Engineering, Halla University. He is interested in HCI.