

# 우편번호 체계에서 사용중인 한글의 빈도수 조사

김 민기, 권 영빈  
중앙대학교 컴퓨터공학과

## A Frequency Measure of Hangeul in Korean Zip Code

Min-Ki Kim and Young-Bin Kwon  
Dept. of Computer Science & Engineering, Chung-Ang University

### 요 약

계약이 없이 자유롭게 쓴 오프라인 필기체 한글을 인식하는 문제는 응용분야에 따른 도메인의 정보를 이용함으로써 보다 쉽게 접근할 수 있다. 본 연구는 오프라인 필기체 한글 인식을 위한 한 도메인으로 우편봉투를 대상으로 하였을 때, 우편번호가 할당된 지명과 건물명을 대상으로 글자의 종류와 빈도수를 통계 분석하였다. 분석 결과 가능한 한글 조합 11,172자중 403자만이 쓰이고 있음을 알았다. 이러한 정보는 자소 분할이 어려운 오프라인 필기체 한글 인식에 있어, 문자 단위 정합을 사용했을 때 인식속도 및 인식률 향상에 기여 할 것으로 생각된다.

### I. 서 론

사회가 고도화 되어감에 따라 정보량이 폭발적으로 증가하게 되었고 이에 따른 정보의 수집 및 효율적인 관리를 위하여 자료를 자동으로 입력,처리,관리 및 출력하는 전산화가 필요하게 되었다. 그러나 컴퓨터를 비롯한 사무자동화 기기들의 발전에도 불구하고 자료 입력은 대부분 키보드를 통한 수작업으로 이루어지고 있는 실정이다.

문서의 자동 입력을 위해서는 문서 구조 분석과 문자 인식이 선행되어야 한다 [13,14,15,16]. 국내의 문자인식, 특히 한글인식에 관한 연구를 살펴보면 문서의 자동 입력에 필요한 오프라인 인쇄체 인식에 대한 연구[4,5,6]와 최근 대두되고 있는 펜 컴퓨터에 활용되는 온라인 필기체 한글인식에 대한 연구[1,2]가 활발히 진행되어 상품화가

가능하거나 이미 상품화한 수준에 와있다. 그러나 오프라인 필기체 한글인식은 아직도 연구가 미흡한 상태이다.

오프라인 필기체 한글인식은 한글의 문자수가 방대하고 문자간의 유사성이 심할 뿐 아니라 각 개인간의 필기 형태에 따른 변형때문에 상당한 어려움을 내포하고 있다. 따라서 오프라인 필기체 한글인식을 위해서는 다소 필기 형태에 제약을 가하거나 응용될 분야의 도메인 정보를 활용하는 방법[3,11,12]이 있다. 본 연구에서는 오프라인 필기체 한글인식의 응용중의 한 방법으로 우편봉투에 나타난 주소 인식을 생각하였다. 그리고, 우편번호 체계를 분석하여 나타나는 지명과 건물명을 구성하는 글자의 종류와 빈도수를 먼저 조사하게 되었다.

이 연구는 편지봉투의 인식에 따른 자동 분류를 위한 기초 연구의 일환으로 수행되었다. 우편번호 체계에 존재하는 글자의 종류 및 특성을 분석하여 인식기의 구성에 도움을 줄 수 있도록 하였다.

## II. 우리나라 우편번호 체계

우편물의 구분작업을 능률적으로 처리하기 위하여 세계 각국은 우편번호제를 채택하여 행선지를 코드화하고 구분작업의 기계화를 추진하고 있다. 더욱이 우편봉투내의 주소를 자동으로 인식할 수 있는 기기의 개발에도 많은 심혈을 기울이고 있어, 실용화된 제품도 많이 개발되어 판매되고 있다[9,10].

우편번호제는 1962년 서독에서 제도화된 것이 처음이며 우리나라에서는 70년 7월 1일을 기해 실시 되었다[7]. 처음 제정된 우편번호는 60년대말의 전국 주요 철도역을 기점으로 한 교통망에 기준을 두어 우편물을 배달하는 우체국에 3 또는 5자리 숫자를 부여한 것이었다. 그러나 70년 이후 경제의 고도성장에 따라 정보의 교환량이 급속도로 증가하여 이의 처리에 대응할 수 있는 우편번호 제도의 체계화가 요구되었다. 배달 우체국 단위로 부여된 그동안의 우편번호는 모든 주소를 나타낼 수 없기 때문에 우편물의 신속한 대량처리라는 측면에서 비효율적인 면이 많았다. 이러한 불편을 없애고 전산화에 적절히 대처할 수 있는 새로운 우편번호 체계가 제정되어 88년 2월 1일부터 시행되고 있다[7].

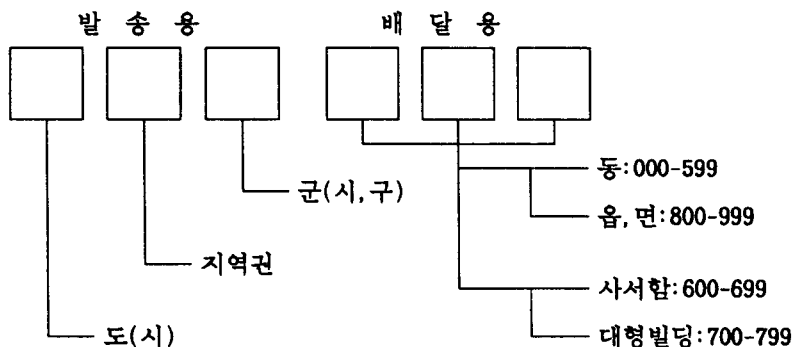


그림 1. 우편번호의 구성

우리나라의 우편번호는 발송용과 배달용으로 구분되어 도(시), 지역권, 군(시·구)의 번호가 3자리로 나열되어 있고 배달용으로 동·읍(면)을 3자리로 구분하여 전체가 6자리 숫자로 구성되어 있다(그림 1). 이 밖에 1일 평균 1,000통(면단위의 1일 평균 배달물량에 해당) 이상의 우편물이 배달되는 다량배달처(대형 빌딩)에는 고유번호가 부여되어 있다[8].

### III. 우리나라 주소의 구성

우리나라 주소는 행정구역의 구조에 따라 계층적으로 이루어져 있다(그림 2). 최상위 계층은 도(특별시·직할시)이고 다음 계층은 군(시·구)이다. 구,시 아래에는 동이, 군 아래에는 읍(면)이 존재하나 성남시 분당구, 안양시 만안구 등 규모가 커진 시는 그 아래 구를 두고 있다. 다음으로 동·읍(면),리(동)으로 이루어져 있다.

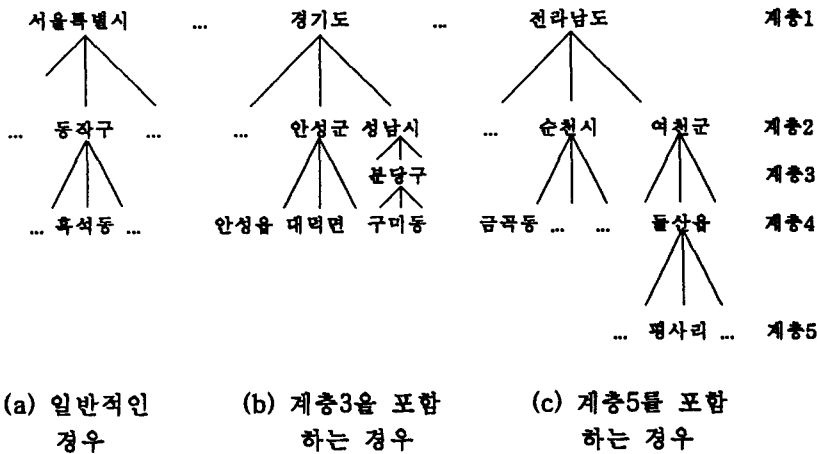


그림 2. 계층적 주소 체계

우편물 자동 분류를 목적으로 주소를 분류하면 크게 들로 나눌수 있다. 주소 전반부는 계층1과 계층2로 이루어지며 성남시, 안양시 등 7개(92년 10월 1일 현재)의 규모가 큰 시는 계층3을 포함하게 된다(그림 2-(b)). 주소 후반부는 계층4로 이루어져 있으나 배달상의 편의를 위해 별도의 우편번호가 할당되는 경우는 계층5를 포함하게 된다(그림 2-(c)). 주소 전반부는 우편번호 상위 3자리 발송용에 대응되고 주소 후반부는 하위 3자리 배달용에 대응된다.

#### IV. 글자의 종류와 빈도수 조사

주소에 나타나는 글자의 종류와 빈도수를 조사하기 위해서 먼저 체신부에서 발행한 우편번호부[7]를 기초로 우편번호 사전을 구성하였다. 사전은 주소 체계에 따라 계층적으로 설계하였고 사서함에 관한 우편번호는 별도로 구성하였다.

주소 전반부는 계층1,2,3을 하나로 했으며, 주소 후반부는 우편번호 하위 3자리가 할당되는 동·읍(면), 사서함, 대형 빌딩을 기준으로 하고 우편 배달의 편의성에 따라 낙도, 오지와 같은 지리적 여건으로 별도의 우편번호가 할당되는 예외적인 경우는 리(동)을 포함시켰다.

오프라인 필기체 한글 주소 인식에 활용할 목적으로, 주소 전반부와 후반부에 나타나는 글자를 대상으로 종류와 빈도수를 조사하였다. 주소 전반부에 나타나는 글자를 대상으로 조사하여 <표 1>의 결과를 얻었다. 주소 후반부는 지명, 건물명, 사서함으로 나누어 분석하였다<표 2 : 사용되는 글자의 총괄표는 부록1. 참조>. <표 1>과 <표 2>에는 내용별로 나누었을 때 중복하여 나타나는 것을 모두 표현하였다. 물론 <표 1>의 각 항에 나타나는 글자들은 <표 2> '전체'란의 각 항에 나타나는 글자들에 모두 포함된다. 그 결과 가능한 한글 조합 11,172자 중에서 우편번호가 할당된 지명과 건물명에 사용되는 글자는 총 403자에 불과하다는 것을 발견하였으며 음절별로 나타날 수 있는 글자의 수효도 각기 다른 것을 알게 되었다. 또한 우리나라 우편번호 체계에서 나타난 403자를 한글의 6 타입으로 분류 하여 빈도수를 표시하면 <표 3>과 같이 나타나고 있다.

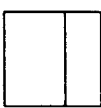
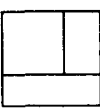

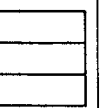
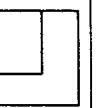
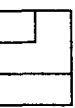
표 1. 주소 전반부에 나타나는 글자

첫 음절에 나타나는 글자	106자
둘째 음절에 나타나는 글자	79자
나타나는 모든 글자의 수효(중복제외)	143자

표 2. 주소 후반부에 나타나는 글자

	지명	건물명	사서함	전체
첫 음절에 나타나는 글자	290자	117자	55자	304자
둘째 음절에 나타나는 글자	276자	119자	47자	294자
나타나는 모든 글자의 수효 (중복제외)	355자	246자	106자	403자

표 3. 403자에 대한 Type별 분류

한글 6 Type						
빈도수	58개	169개	42개	102개	19개	13개

## V. 결 론

주소에 나타나는 글자의 종류와 빈도수를 조사해본 결과 우편물을 발송지별로 분류하기 위해서는 주소의 전반부만을 인식하면 되기 때문에 143자의 한글을 인식하면 되고, 배달지별로 분류하기 위해서는 403자만을 인식하면 됨을 알았다. 신도시나 새로운 건물등이 생겨난다 하더라도 예외적인 경우를 제외하고는 기존의 지명, 건물명에 나타나는 글자들을 사용할 것으로 예측되기 때문이다.

자소 분할이 어려운 오프라인 필기체 한글 인식에 있어 문자 단위 정합은 자소 분할이 필요없고 필기 형태에 대한 제약을 가하지 않기 때문에 가장 일반적인 방법으로 볼 수 있다. 그러나, 정합해야 할 부류의 수가 많아 자소 단위 정합이나 영역 단위 정합에 비해 상대적으로 인식 속도가 느리고 높은 인식률을 기대할 수 없었다.

오프라인 필기체 한글 인식의 한 도메인으로 우편봉투를 대상으로 할때 상기에 조사된 결과를 활용하여 문자 단위 정합을 수행한다면 인식속도는 물론 인식률을 향상시킬수 있을 것으로 생각된다. 또한 계층적으로 설계된 사전 정보를 활용하고 우편번호를 이용하여 후처리를 수행할 경우 오인식 및 미인식을 크게 줄일 수 있을 것이다.

## VI. 참고문헌

- [ 1 ] 권 오성, 권 영빈, “동적인 선분생성을 이용한 온라인 한글 필기 인식”, 한국정보과학회 봄 학술발표 논문집, 제 20권 1호, pp. 151-154, 1993년 4월.
- [ 2 ] 신 봉기, 김 진형, “통계적 방법에 의한 온라인 한글 필기 인식”, 제4회 한글 및 한국어 정보처리 학술발표 논문집, pp. 533-542, 1992년 10월.
- [ 3 ] 이 성환, 박 정선, “통계적 특징 추출 방법을 이용한 샘플체 필기 한글의 오프라인 인식”, 제4회 한글 및 한국어 정보처리 학술대회 논문집, pp. 237-248, 1992년 10월.
- [ 4 ] 최 동혁, 류 성원, 강 현철, 박 규태, “계층 구조 신경망을 이용한 한글 인식”, 대한전자공학회 논문지, 제 28권 B편 제 11호, pp. 1-7, 1991년 11월.
- [ 5 ] 김 화룡, 방 승양, “네오코그니트론 방식의 신경 회로망을 이용한 한글 문자 인식”, 한국정보과학회 논문지, 제 18권 제 3호, pp. 280-289, 1991년 5월.
- [ 6 ] 정 주성, 김 춘석, 박 충규, 황 회용, “윤곽선 추적에 의한 고딕체 한글의 신속 인식에 관한 연구”, 전기학회 논문지, 제 37권 제 8호, pp. 579-587, 1988년 8월.
- [ 7 ] 학원세계대백과사전(전32권), 제 22권, pp. 437-441, 학원출판사, 1993년 1월.
- [ 8 ] 체신부, “우편번호부(1992년 10월 1일 현재)”, pp. 2-63, 1992.
- [ 9 ] Mauricette Feuillas, Serge Hugot, Emmanuel Miette, “Architecture of a Multigrey Level Videocoding System”, JET POSTE 93, pp. 495-502, June 1993.
- [ 10 ] T. Ono, T. Shima, M. Yoshioka, “Automated Material Handling Systems in ‘Osaka Letter Post Office’ and ‘Osaka Parcel Post Office’”, JET POSTE 93, pp. 465-472, June 1993.
- [ 11 ] Michel Gilloux, Jean-Michel Bertille, Manuel Leroux, “Recognition of Handwritten Words in a Limited Dynamic Vocabulary”, JET POSTE 93, pp. 148-155, June

1993.

- [ 12 ] Venu Govindaraju, Ajay Shekhawat, Sargur N. Srihari, "Interpretation of Handwritten Addresses in US Mail Stream", IWFHR III, pp. 197-206, May 1993.
- [ 13 ] Kyoong Ha Lee, Kie-Bum Eom, R. L. Kashyap, "Character Recognition Based on Attribute-Dependent Programmed Grammar", IEEE Trans. on PAMI, Vol 14, No. 11, pp. 1122-1128, Nov. 1992.
- [ 14 ] Sargur N. Srihari, "High-Performance Reading Machines", proceedings of the IEEE, Vol. 80, No. 7, pp. 1120-1132, July 1992.
- [ 15 ] Shuichi Tsujimoto, Haruo Asada, "Major Components of a Complete Text Reading System", proceedings of the IEEE, Vol. 80, No 7, pp. 1133-1149, July 1992.
- [ 16 ] Joachim Kreich, Achim Luhn, Gerd Maderlechner, "An Experimental Environment for Model Based Document Analysis", ICDAR 91, pp. 50-58, 1991.

## 부록 1. 우편번호 체계에서 사용중인 한글의 빈도수

동성	4596	면학	1389	대천	520	산도	497	가곡	449	리정	394	서읍	357	신원	355	사화	285	남문	281
국우	281	교안	281	전천	280	도수	237	곡평	229	정양	229	읍장	209	원부	205	화송	202	문암	201
구함	200	안덕	199	상월	193	용주	186	북현	181	중포	179	지하	178	금청	169	내삼	169	암봉	163
빌의	163	체안	157	건월	156	방주	147	촌현	146	계포	145	림하	143	이청	141	호삼	140	로봉	136
소두	134	물덕	134	광월	134	명주	134	창현	131	고포	130	영하	125	석청	121	연호	121	일당	118
아함	117	진안	116	홍월	108	인주	94	공현	94	한포	92	오하	90	당석	89	연기	87	일운	87
비소	86	선안	84	자월	84	미주	82	보현	81	계포	81	경하	80	관당	77	기강	73	운마	73
등두	70	유안	69	목월	69	양주	66	옥현	66	무포	66	울하	64	태관	63	풍강	62	죽마	58
아함	57	외안	56	회월	52	여주	52	초현	51	백포	50	충하	50	량관	50	해강	50	노마	49
황비	48	시안	48	항월	48	개주	47	업현	45	은포	44	조하	44	룡관	42	만강	42	본마	41
등소	41	군안	40	내월	40	행주	39	효현	38	반포	38	복하	38	종관	37	모강	37	단마	37
길등	37	감안	36	좌월	36	법주	36	야현	36	홍포	35	과하	35	세관	33	직강	32	왕마	31
간길	31	접안	30	거월	30	농주	29	역현	29	임포	29	망하	28	세관	28	직강	28	왕마	26
합은	25	병안	25	탄월	24	갈주	24	래현	24	생포	24	실하	23	상관	23	음강	23	치마	23
은토	22	다안	22	달월	22	라주	21	예현	21	적포	21	통하	21	상관	21	어강	21	파마	21
공공	21	옥안	20	둔월	20	락주	20	류현	20	향포	20	심하	19	낙관	18	통강	18	민마	18
담공	18	후안	18	괴월	17	불주	17	센현	17	스포	17	위하	16	울관	16	철강	16	철마	16
담공	16	귀안	16	능월	15	배주	15	순현	15	십포	15	재하	15	잠관	15	타강	15	검마	14
담공	14	독안	14	륜월	14	응주	14	팔현	14	십포	14	혜하	14	저관	14	캠강	14	퍼마	14
추담	13	춘안	13	륜월	13	응주	13	팔현	13	십포	13	혜하	13	저관	14	캠강	14	퍼마	14
완승	12	나안	12	령월	11	척주	11	옥현	11	회포	11	각하	10	담관	10	승강	10	와마	10
승려	10	악안	9	막월	9	변주	9	속현	9	응포	9	입하	9	작관	9	축강	9	탐마	9
려목	8	손안	8	악월	8	익주	8	총현	8	판포	8	격하	7	김관	7	녹강	7	랑마	7
목별	7	번안	7	술월	7	압주	7	염현	7	중포	7	차하	7	환관	7	논강	6	련마	6
별린	6	빈안	6	침월	6	찰주	6	채현	6	대포	5	록하	5	료관	5	묘강	5	박마	5
린패	5	빈안	5	색월	5	험주	5	훈현	5	갑포	4	권하	4	누관	4	돈강	4	롯데	4
패녕	4	필안	4	색월	4	협주	4	요현	4	잔포	4	카하	4	택관	4	돌강	4	퇴마	4
녕섬	4	돌안	4	할월	4	협주	4	횡현	4	곤포	3	곳하	3	극관	3	남강	3	빙마	3
섬휘	3	연안	3	드월	3	두주	3	릭현	3	릭포	3	립하	3	밀관	3	분강	3	빙마	3
휘냉	3	연안	3	엄월	3	절주	3	출현	3	코포	3	키하	3	더관	3	형강	3	활마	3
냉소	3	유안	3	견월	2	결주	2	골현	2	곳포	2	규하	2	균관	2	그강	2	난마	2
소걸	2	늘안	2	니월	2	력주	2	렴현	2	론포	2	통하	2	팔관	2	삼강	2	새마	2
걸느	2	숙안	2	애월	2	죽주	2	출현	2	표포	2	피하	2	핑관	2	현강	2	혜마	2
느몽	1	겨안	1	검월	1	괘주	1	골현	1	필포	1	날하	1	닐관	1	넝강	1	뉴마	1
몽열	1	늑안	1	독월	1	동주	1	란현	1	탤포	1	릴하	1	래관	1	특강	1	맹마	1
열운	1	밤안	1	봉월	1	삭주	1	살현	1	채포	1	최하	1	알관	1	영강	1	알마	1
운탁	1	엽안	1	울월	1	음주	1	웅현	1	왜포	1	육하	1	육관	1	위강	1	윈마	1
탁힐	1	접안	1	존월	1	집주	1	징현	1	책포	1	취하	1	취관	1	추강	1	컴마	1
힐	1	탕안	1	탈월	1	트주	1	류현	1	평포	1	품하	1	프관	1	피강	1	혜마	1