

발음사전 표제어중의 음소의 통계적 성질
- 음성 DB용 단어선정을 위하여 -

○ 이 용주, 김경태, 조철우, 이태원
한국전자통신연구소 *고려대학교

On the statistics of Korean Phonetic Dictionary
- Basic Survey to make corpus of Korean Speech DB -

Y. J. LEE, K. T. KIM, C. W. JO, T. W. RHEE
ETRI * Korea Univ.

Abstract

Statistical information about spoken Korean was obtained. The data are the results of analyzing the Korean phonetic dictionary. This is one of the basic survey to make phoneme balanced corpus of Korean Speech Data Base(KSDB).

의 우리말 발음사전의 표제어에 대한 음소 또는 음소열 관련의 통계적인 조사결과에 대하여 논하였다.

1. 서론

2. 분석자료 및 방법

불특정 발성자의 대어휘. 연속어인식이나 규칙합성을 위해서는 각종환경하에서의 음운변동을 파악할 필요가 있으므로 다수 발성자의 다양한 음운환경을 가진 음성데이터를 음소와 같은 미세단위로 labelling한 음성데이터베이스(DB)가 필요하다. 따라서, 각국에서는 본격적인 음성연구의 기반으로 이 음성 DB의 구축을 조직적으로 시도하고 있다. (1)(2)(3)

- 대상자료: 한국어 발음사전(9) (8575 단어)
- 입력방법: 발음사전의 표제어와 함께 한글로 발음표기된 단어를 2byte code로 컴퓨터에 입력.

필자들도 한국어 음성 DB의 필요성을 절감하고 이에 대한 기초적인 검토를 수행중에 있다. (1)(4)(5)

- 음소 및 음소열의 선정:
음소는 일반적인 음소분류(10)를 따랐으나, 중모음은 반모음과 단모음으로 분리하였고 음절의 어두자음과 어말자음은 음소분류에 상관없이 서로 구분하여 다루었다. (표1참조)
2음소열은 위와같이 분류된 음소단위의 쌍으로 구성하였고
3음소열은 단어의 첫머리와 끝어오는 경우도 같이 고려하였다.

일반적으로 다양한 음운환경을 가진 음성 DB용 대상 단어를 선정하기 위해서는 사전의 표제어, 중요어, 고빈도 단어 등을 대상으로 2음소열 또는 3음소열의 종류수, Entropy 등을 매개체로한 여러가지 통계적인 선정방법을 사용하고 있다. (6)(7)

(예: bCV는 단어첫머리에 오는 CV음절로서 공백을 음소처럼 취급하여 3음소열로 설정하였다)

특히 사전의 표제어를 대상으로 할 경우에는 먼저 이 표제어를 발음표기하여 컴퓨터에 입력한후 음소 또는 음소열의 발생 빈도등 각종 통계적인 자료를 도출해야만 한다. (13)

- 제한사항:
본 조사에서는 장모음은 고려되지 않았고 어두의 /의/가 /으이/로 발음표기 되어있어서 이를 그대로 이용하였으므로 별도의 음소로 다룰수 없었음을 먼저 밝혀둔다.

표1. 음소 구분

일반적인 분류	본 조사에서의 분류
모 음 18	모 음 9
중모음 12	반모음 2
반모음 2	어두자음 19
자 음 19	어말자음 7
계 43 종류	계 37 종류

이러한 통계적 자료는 음성 DB 용 대상단어 선정을 위해서 만이 아니고 자연언어처리, 음성인식 및 합성에 기초자료로서도 의미가 있으며 사전 데이터베이스(8) 자체를 우리말의 대량 데이터로 다루므로서 언어학적, 음성학적으로 우리말이 가진 성질에 관해서 유익한 자료를 얻을 수도 있다. 본고에서는 음성 DB용 단어 선정을 위한 기초조사로서 우선 8500단어 규모

앞으로 장모음등 음운환경의 보다 자세한 검토 및 대상어휘의 확대로 DB용 대상 단어 선정만이 아니라 minimal pair 나 음운결합에 따른 각종 신뢰성 있는 통계자료를 얻기 위해서 10만 단어 이상이 수록된 발음사전을 컴퓨터에 입력할 계획이다.

감사의 글

본 연구는 ETRI의 기초연구의 일환으로 이루어진 것이다. 사전입력을 위해 수고해준 이 성구군과 조 화신양에게 감사의 뜻을 표한다.

참고 문헌

1. 이용주, 김경태, "연구 및 평가용 음성데이터베이스의 개발 동향과 제안" 전자통신 8권 3호 (86. 10)
2. 이용주, 김경태, "음성이해 연구의 동향" 전자통신 9권 1호 (87. 3)
3. 일본전자공업진흥협회, "일본어 정보처리 표준화에 관한 조사연구보고서 (85.)"
4. 강철희 외, "ETRI의 음성정보처리분야 연구 활동 소개" 음성통신 및 처리기술 Workshop논문집 (86. 5. 31)
5. 이용주 외, "단어음성데이터의 수집 및 DB 구성 시스템" 전자공학회 추계학술대회 논문집 Vol. 9 No. 2 (86. 12)
6. 鹿野, "일본 음향학의 강연 논문집 3-3-10 (84. 3)
7. 速水 외, "일본 음향학의 강연논문집 2-4-7 (85. 10)
8. 木廣山 외, "일본 전자통신학회 논문집 Vol J68-D No. 12 (85. 12)
9. 전영우, 표준한국어 발음사전, 한국방송사업단 (84. 2)
10. 허용, 국어음운론, 정음사 (85. 2)
11. 문헌10 P. 213 의 참고문헌을 계 인용
정철; "국어 음소 배열의 연구"
12. 남광우외, 한국어 표준 발음사전 한국 정신문화연구원 (84. 11)
13. Denes, "On the statistics of spoken English" JASA 35-6 (63. 6)

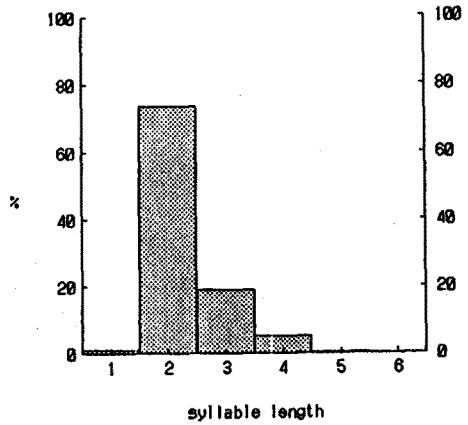


그림1. 음절길이별 단어분포

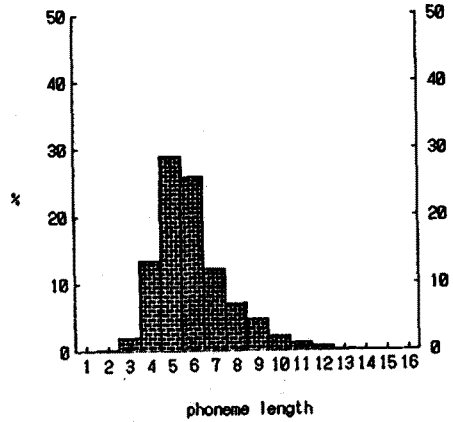


그림2. 음소길이별 단어분포

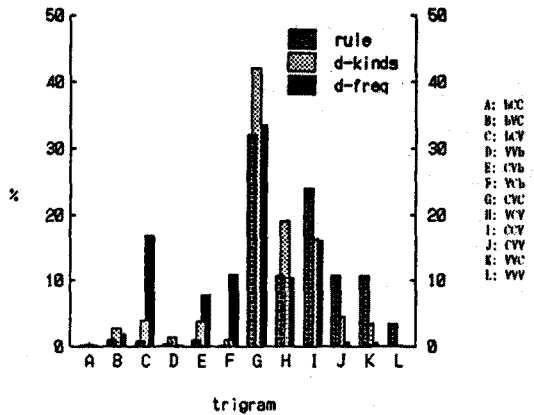


그림3. 3음소열, 종류수 및 발생빈도

표4. 2음소열의 빈도분포

DIG

이진	이진																										Total
	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	
1/1	447	141	235	20	393	16	285	208	238	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/2	302	145	22	4	259	1	47	20	64	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/3	702	220	46	12	271	15	79	84	296	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/4	125	228	16	5	53	10	38	21	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/5	88	19	39	3	61	4	21	53	177	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/6	309	102	39	3	146	2	334	11	193	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/7	373	112	83	0	171	1	288	3	162	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/8	585	84	369	0	104	0	19	3	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/9	585	84	369	0	104	0	19	3	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/10	66	13	62	9	21	6	24	1	68	407	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/11	316	78	459	76	187	7	205	24	370	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/12	155	35	157	23	57	4	22	42	58	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/13	167	41	1	1	2	22	2	14	14	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/14	127	41	1	1	2	102	26	14	14	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/15	440	201	60	2	200	89	133	45	54	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/20	15	3	15	1	5	0	12	4	86	204	79	216	79	27	0	0	0	0	0	0	0	0	0	0	0	0	0
1/21	289	3	1315	243	300	0	211	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/22	802	26	220	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/23	5388	1369	3461	5	18	2737	198	2150	747	2640	1243	421	1148	1813	193	609	461	555	200	1223	283	520	512	567	180	0	0
1/24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/29	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/31	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/32	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/33	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/34	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/35	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/36	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/38	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/39	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/41	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/42	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/43	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/44	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/45	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/46	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/47	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/48	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/49	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/51	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/52	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/53	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/54	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/55	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/56	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/57	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/58	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/59	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/60	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/61	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/62	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/63	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/64	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/65	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1/66	0	0	0	0	0																						