

한국어 多音節 單語의 초성, 중성, 종성단위의 音節間 조건부 확률

(Conditional Probability of a 'Choseong', a 'Jungseong', and a 'Jongseong'
Between Syllables in Multi-Syllable Korean Words)

李 在 弘,* 李 在 鶴**

(Jae Hong Lee and Jae Hak Lee)

要 約

한국어 단어는 음절이 모여서 이루어지고 단어를 이루는 음절은 그 발생의 확률적 성질에 따라 확률변수로 간주된다. 음절은 확률변수로 간주되는 초성, 중성, 종성의 세 단위로 이루어진다. 한국어 단어 안에서의 음절들의 발생의 상관관계는 음절간 조건부 확률로 표시된다. 그런데 음절간 조건부 확률의 표본공간이 매우 커서 나열할 수 없어서 대신에 초성, 중성, 종성 단위의 음절간 조건부 확률을 고려한다. 한국어 다음절 단어에 있어서 다음절 단어의 발생 빈도수와 종류에 따른 음절의 길이 분포를 계산한다. 다음절 단어의 발생 빈도수를 고려한 음절별 누적 빈도수 통계로부터 음절의 초성, 중성, 종성 단위의 단위별 발생확률과 단위간 조건부 확률을 계산한다. 다음절 단어의 길이와 단어내에서의 음절의 위치에 따른 초성, 중성, 종성 단위의 발생확률과 조건부 확률을 계산한다. 그리고 한국어 다음절 단어내에서 인접 음절사이의 초성-초성, 중성-중성, 종성-종성, 종성-초성 단위의 음절간 조건부 확률을 계산한다.

Abstract

A Korean word is composed of syllables. A Korean syllable is regarded as a random variable according to its probabilistic property in occurrence. A Korean syllable is divided into 'choseong,' 'jungseong,' and 'jongseong' which are regarded as random variables. We can consider the conditional probability of syllable as an index which represents the occurrence correlation between syllables in Korean words. Since the number of syllables is enormous, we use the conditional probability of a 'choseong,' a 'jungseong,' and a 'jongseong' between syllables as an index which represents the occurrence correlation between syllables in Korean words. The length distribution of Korean words is computed according to frequency and to kind. From the cumulative frequency of a Korean syllable computed from multi-syllable Korean words, all probabilities and conditional probabilities are computed for the three random variables. The conditional probabilities of 'choseong' - 'choseong,' 'jungseong' - 'jungseong,' 'jongseong' - 'jongseong,' 'jongseong' - 'choseong' between adjacent syllables in multi-syllable Korean words are computed.

*正會員, **學生會員 서울大學校 電子工學科

(Dept. of Elec. Eng., Seoul Nat'l Univ.)

接受日字: 1991年 7月 3日

(※ 본 논문은 峨山社會福祉事業財團의 1989年度 研究費 支援에 의하여 研究되었음.)

I. 서 론

언어는 인간의 의사소통의 가장 중요한 수단으로서 언어가 문자화 된 것이 글이다. 글은 기호들이 연속적으로 나열되어 형성된 기호열로 볼 수 있다.

개의 단어에 대한 빈도순서와 잣기로부터 다음절 단어에 있어서 빈도수와 종류에 따른 음절 길이의 분포와 음절의 발생확률을 계산하고, 전체 음절의 초성, 중성, 중성 단위간의 조건부 확률을 계산하고 다음절 단어의 길이와 단어내에서 음절의 위치에 따른 초성, 중성, 중성 단위의 발생확률과 조건부 확률을 계산하며 다음절 단어의 초성, 중성, 중성 단위의 음절간 조건부 확률을 계산한다.

1. 한국어 다음절 단어의 음절의 길이 분포와 음절별 발생확률

한국어 단어의 빈도수와 종류에 따른 음절의 길이 분포를 표1,2에 보인다. 음절의 길이를 나타내는 확률변수를 L이라 하자. 여기서 L=10은 길이가 10인 단어와 그 이상인 단어를 나타낸다. 한국어 단어는 1음절 단어의 빈도가 가장 높고 2음절 단어의 종류가 가장 많음을 표1과 2에서 볼 수 있다. 1음절 단어의 발생빈도가 가장 높은 것은 조사가 있기 때문이다. 한국어 단어는 약 40.3%가 2음절 단어이고 5음절 이상의 단어는 약 0.5%만이 발생됨을 표2에서 볼 수 있다.

표 1. 한국어 다음절 단어의 빈도수에 따라 음절의 길이 분포

Table 1. Word length distribution of multi-syllable korean words according to frequency.

L	1	2	3	4	5	6	7	8	9	10
p(l)	.421480	.413039	.099342	.061079	.004240	.000730	.000071	.000011	.000004	.000003

표 2. 한국어 다음절 단어의 종류에 따른 음절의 길이 분포도

Table 2. Word length distribution of multi-syllable korean words according to kind.

L	1	2	3	4	5	6	7	8	9	10
p(l)	.034310	.402982	.285375	.218236	.043458	.013802	.001230	.000392	.000143	.000071

확률변수 X, Y, Z의 표본공간의 크기는 각각 $|A_x|=19$, $|A_y|=21$, $|A_z|=28$ 이므로 가능한 초성, 중성, 중성의 결합가지수는 $19 \times 21 \times 28 = 11,172$ 개이다. 그런데 "우리말 말수 사용의 잣기 조사"에서 조사된 단어에는 모두 1,535개의 음절만이 사용되었다. 이는 모든 가능한 초성, 중성, 중성의 결합의 약 13.7%만이 발생됨을 뜻한다. 한국어 단어에 있어서 발생하는 음절의 종류는 모든 가능한 초성, 중성, 중성의 결

합에 비하여 매우 적음을 알 수 있다. 발생하는 음절을 발생확률이 높은 순서로 나열했을 때 가장 많은 발생하는 음절10개와 그것의 발생확률을 표 3에 보인다. 위의 자료에서 동사와 형용사는 기본형으로 수록되어 있으므로 음절 '다'와 '하'가 다른 음절에 비하여 발생빈도가 높다. 실제로 발생하는 1,535개의 음절중 가장 많이 발생하는 10개의 음절이 약 36.0%를 차지한다.

표 3. 음절의 발생확률
Table 3. Syllable probability.

S	다	이	히	에	가	의	음	는	리	그
p(s)	.152224	.045519	.044541	.023764	.018903	.018675	.018591	.013200	.012473	.012175

2. 다음절 단어의 길이와 단어내에서의 음절의 위치에 따른 초성, 중성, 중성 단위의 발생확률과 조건부 확률

전체 단어에서 발생 빈도를 고려하여 음절별 누적 빈도수를 구하였다. 이로부터 음절의 초성, 중성, 중성 단위의 발생확률과 초성, 중성, 중성 단위의 조건부 확률을 계산하여 각각 표4와 표5에 보인다. 초성에서는 'ㅇ'의 발생빈도가 가장 높고 그 다음으로 'ㄷ', 'ㄱ'의 순서로 발생빈도가 높음을 표 4(a)에서 볼 수 있고 중성에서는 'ㅏ'의 발생빈도가 가장 높고 그 다음으로 'ㅣ', 'ㅡ'의 순서로 발생빈도가 높음을 표 4(b)에서 볼 수 있고 중성에서는 '공백소'의 발생빈도가 가장 높고 그 다음으로 'ㄴ', 'ㄱ'의 순서로 발생빈도가 높음을 표 4(c)에서 볼 수 있다. 음절에서 초성, 중성, 중성 단위간의 조건부 확률에는 중성이 Y일 때 초성이 X일 조건부 확률 $p(x|y)$, 중성이 Z일 때 초성이 X일 조건부 확률 $p(x|z)$, 초성이 X일 때 중성이 Y일 조건부 확률 $p(y|x)$, 중성이 Z일 때 중성이 Y일 조건부 확률 $p(y|z)$, 초성이 X일 때 중성이 Z일 조건부 확률 $p(z|x)$, 중성이 Y일 때 중성이 Z일 조건부 확률 $p(z|y)$ 등이 있다. 이 여섯 종류의 조건부 확률을 모두 보이려면 그 양이 방대하므로 중성이 Y일 때 초성이 X일 조건부 확률 $p(x|y)$ 만을 표 5에 보인다.

단어의 길이에 따른 음절의 초성, 중성, 중성 단위의 발생확률과 초성, 중성, 중성 단위의 조건부 확률을 계산하여 각각 표6과 표7에 보인다. 단어의 길이가 5 이상인 단어는 0.05% 뿐이므로 단어의 길이가 4 이하인 경우만 고려한다. 여섯 종류의 조건부 확률을 모두 보이려면 그 양이 방대하므로 중성이 Y

표 4. 초성, 중성, 중성의 발생확률
Table 4. Probabilities of a choseong, a jungseong, and a jongsung.

Table with 2 rows and 21 columns. Row 1: X (초성) with Korean characters. Row 2: p(x) with numerical probabilities.

(a) 초성의 발생확률 P(x)

Table with 2 rows and 21 columns. Row 1: Y (중성) with Korean characters. Row 2: p(y) with numerical probabilities.

(b) 중성의 발생확률 P(y)

Table with 2 rows and 21 columns. Row 1: Z (중성단위) with Korean characters. Row 2: p(z) with numerical probabilities.

(c) 중성의 발생확률 P(z)

표 5. 중성이 Y일 때 초성이 X일 조건부 확률 p(x|y)
Table 5. Conditional probability p(x|y).

Large table with 21 rows (Y) and 21 columns (X). Each cell contains a conditional probability value.

표 6. 단어의 길이에 따른 초성, 중성, 중성의 발생확률
Table 6. Probabilities of a choseong, a jungseong, and a jongsung as a function of word length.

Table with 2 rows and 21 columns. Row 1: X (초성) with Korean characters. Row 2: p(x|i) for i=1, 2, 3, 4.

(a) 단어의 길이에 따른 초성의 발생확률 p(x|i)

Table with 2 rows and 21 columns. Row 1: Y (중성) with Korean characters. Row 2: p(y|i) for i=1, 2, 3, 4.

(b) 단어의 길이에 따른 중성의 발생확률 p(y|i)

Table with 2 rows and 21 columns. Row 1: Z (중성단위) with Korean characters. Row 2: p(z|i) for i=1, 2, 3, 4.

(c) 단어의 길이에 따른 중성의 발생확률 p(z|i)

일 때 초성이 X일 조건부 확률 $p(x|y)$ 만을 표7에 보인다.

단어내에서의 음절의 위치에 따른 음절의 초성, 중성, 종성 단위의 발생확률과 초성, 중성, 종성 단위의 조건부 확률을 계산하여 각각 표8과 표9에 보인다. 각 단어에서 첫번째 음절과 두번째 음절 그리고 끝 음절에 대해서 계산하였다. 여섯 종류의 조건부 확률을 모두 보려면 그 양이 방대하므로 중성이 Y일 때 초성이 X일 조건부 확률 $p(x|y)$ 만을 표9에 보인다.

3. 초성, 중성, 종성 단위의 음절간 조건부 확률
음절간의 상호 발생의 상관관계를 표시하는 척도로 조건부 확률 $p(s_n|s_p)$ 이 사용된다. 조건부 확률 $p(s_n|s_p)$ 는 앞 음절이 S_p 일때 그 다음 음절이 S_n 일 확률이다. 그런데 위의 자료에서 실제로 발생하는 음절수는 1,535이므로 앞 음절 S_p 와 그 다음 음절 S_n 간의 조건부 확률 $p(s_n|s_p)$ 을 구하면 $1,535 \times 1,535$ 개의 확률이 구해지므로 표시가 힘들어 진다. 따라서 음절간 상호 발생의 상관관계를 표시하는 척도로 음절간 초성, 중성, 종성 단위의 조건부 확률을 고려한다.

표 7. 음절의 길이에 따른 조건부 확률 $p(x|y)$

Table 7. Conditional probability $p(x|y)$ as a function of word length.

Table with 19 columns (X, Y, and 18 consonants) and 19 rows (Y, and 18 consonants). It contains conditional probability values for word length L=1.

(a) 음절의 길이 L = 1

Table with 19 columns (X, Y, and 18 consonants) and 19 rows (Y, and 18 consonants). It contains conditional probability values for word length L=2.

(b) 음절의 길이 L = 2

표 7. 음절의 길이에 따른 조건부 확률 p(x|y)

Table 7. Conditional probability p(x|y) as a function of world length.

X \ Y	ㅏ	ㅑ	ㅓ	ㅕ	ㅗ	ㅛ	ㅜ	ㅠ	ㅡ	ㅣ	ㅚ	ㅜ	ㅟ	ㅝ	ㅞ	ㅟ	ㅠ	ㅡ	ㅢ	ㅣ	ㅤ	
ㅏ	.059497	.004945	.064719	.514158	.012978	.017202	.051917	.025684	.003821	.029313	.003031	.044923	.038601	.000775	.008252	.000153	.002985	.004419	.112626			
ㅑ	.033354	.023793	.139145	.272051	.009926	.047185	.035726	.051491	.012334	.183557	.001168	.035252	.031931	.010656	.020728	.001679	.014269	.004635	.071124			
ㅓ	.003181	-	.006893	-	-	.093849	-	.001060	.000795	.004242	-	.732768	.000265	.000265	-	-	-	-	.156681			
ㅕ	.072607	.013666	.015888	.033459	.038054	.122038	.049785	.084085	.004208	.144843	.000488	.147502	.213983	.002440	.027939	.003854	.007877	.006076	.011209			
ㅗ	.023976	.017982	.040000	.190092	.014067	.066300	.021896	.019450	.003180	.266055	.005627	.023853	.195719	.005505	.056147	-	.002936	.016758	.030459			
ㅛ	.132827	.002249	.021945	-	-	.152097	.103222	.132340	.006565	.000851	-	.319088	.000061	.000061	.004134	.000304	-	.057082	.067173			
ㅜ	.686309	-	-	-	.043420	-	-	-	-	-	.227293	-	-	-	-	-	-	.029685	.013292			
ㅠ	.124787	.017170	.053251	.108458	.009525	.087924	.155353	.074059	.003383	.090955	.001604	.132473	.039034	.006101	.016153	.011382	.042886	.008018	.021081			
ㅡ	.445741	-	-	-	-	-	.002139	-	-	-	.000389	.094710	.022170	-	.000972	-	-	-	.433878			
ㅣ	.400000	.024138	-	.055172	-	-	-	-	-	.327586	.003448	.048276	-	-	-	.062069	.024138	-	.055172			
ㅚ	.105704	.033838	.004189	.105382	-	.077022	.009668	-	-	.077345	.000967	.157912	.011602	-	.046729	-	.012246	-	.357396			
ㅜ	.313543	-	.001131	-	-	.109132	.010178	-	.001131	.000565	-	.374329	.002262	-	-	-	-	.170766	.016964			
ㅟ	.104675	.026107	.034779	.037917	.004101	.030074	.173853	.160721	.019026	.096943	.000627	.137056	.070298	.004146	.046208	.000314	.008605	.036483	.008067			
ㅝ	.165628	.001099	-	.001099	-	-	-	-	-	-	.827043	-	-	-	-	-	-	-	.005130			
ㅞ	.038462	.480769	-	.057692	-	-	-	-	-	-	.057692	-	-	.019231	.057692	.269231	-	-	.019231			
ㅟ	.106586	.010912	.001364	.074045	.011691	-	-	-	-	.023772	-	.623149	.004287	-	.099182	.008769	.002338	-	.033905			
ㅠ	.075644	-	.034871	.001073	-	.175429	-	-	-	.003219	.000536	.667918	.002146	-	-	-	-	-	.015021	.024142		
ㅡ	.247250	.022159	.036509	.155188	.009463	.318262	.000078	.001266	.016746	.036159	.012988	.061647	.015636	.000058	.004323	.002122	.016960	.010164	.033024			
ㅢ	-	-	.001684	.000561	.046042	-	-	-	-	-	.029759	.807412	-	-	-	-	-	.003369	-	.111173		
ㅣ	.112582	.010145	.085628	.006479	.000422	.172302	.025029	.038823	.000932	.074876	.005275	.163300	.153289	.004774	.050840	.012721	.000923	.011657	.070023			

(c) 음절의 길이 L = 3

X \ Y	ㅏ	ㅑ	ㅓ	ㅕ	ㅗ	ㅛ	ㅜ	ㅠ	ㅡ	ㅣ	ㅚ	ㅜ	ㅟ	ㅝ	ㅞ	ㅟ	ㅠ	ㅡ	ㅢ	ㅣ	ㅤ
ㅏ	.042119	.002099	.023868	.455599	.002942	.015619	.013769	.019461	.001266	.024768	.001270	.041308	.018906	.000615	.009099	.000968	.007786	.003272	.315266		
ㅑ	.064887	.028066	.124247	.142999	.001118	.031233	.039242	.058367	.008196	.253089	.000683	.017013	.051785	.001366	.016020	.001676	.022353	.013971	.123686		
ㅓ	.026134	-	.028107	-	-	.168639	-	-	.000493	.016272	-	.671105	.000986	-	-	-	-	-	.088264		
ㅕ	.046156	.006458	.014949	.011283	.084506	.226452	.011017	.036127	.003571	.098771	.000893	.265732	.142705	.002982	.022546	.004692	.002621	.002203	.016335		
ㅗ	.025790	.008344	.052592	.059418	.036410	.087231	.020986	.002028	.026549	.203793	-	.076359	.257901	.001770	.052592	.006068	.017699	.001264	.045006		
ㅛ	.201881	.000812	.011704	-	-	.182126	.163182	.067248	.003518	.000609	-	.193086	.013531	-	.026182	.003315	-	.075096	.057709		
ㅜ	.653763	-	-	-	-	.089606	-	-	-	-	.158423	-	-	-	-	-	-	.093190	.005018		
ㅠ	.105903	.012853	.029515	.142639	.007188	.084603	.039320	.122466	.003909	.086226	.004538	.143136	.103584	.005333	.021664	.001226	.033192	.020803	.031900		
ㅡ	.264102	.000236	-	-	-	-	-	-	.000944	-	-	.150578	.016993	.000472	.001180	-	-	-	.565494		
ㅣ	.098039	-	-	.068627	.002451	-	-	-	-	.183824	-	.007353	-	-	-	.639706	-	-	-		
ㅚ	.036688	.002291	.002291	.806739	-	.004852	-	-	-	.015768	.000135	.016712	.004717	.002156	.011995	.000135	.020216	-	.075337		
ㅜ	.092551	-	-	-	.015185	.015185	-	-	.005435	.000480	-	.727781	.000160	.000160	.000160	-	-	-	.134910	.007992	
ㅟ	.092000	.007280	.010882	.034022	.010325	.021468	.132298	.199079	.008617	.108416	.002043	.093782	.148083	.002451	.080486	.000223	.009025	.017308	.022211		
ㅝ	.090226	-	-	-	-	-	-	-	-	-	.899033	-	-	-	.003222	-	-	-	.007519		
ㅞ	.035714	-	-	-	-	-	-	-	-	-	.821429	-	-	.035714	-	-	-	-	.107143		
ㅟ	.157243	.003770	.021002	.050619	.047388	-	-	-	.002693	-	.369952	.018848	.000539	.181475	.018848	.001077	-	-	.126548		
ㅠ	.057779	-	.000696	-	-	.151758	.001392	.001044	-	.001740	.003481	.754264	-	-	-	-	-	-	.027845		
ㅡ	.361691	.041141	.027291	.186817	.035934	.090278	.023446	.004322	.002348	.059312	.029469	.046143	.016402	.001940	.008167	.005921	.030762	.006704	.021914		
ㅢ	-	-	.003739	-	.010684	-	-	-	-	-	.002671	.854701	-	-	-	-	-	-	.128205		
ㅣ	.069319	.003448	.081541	.004643	.000205	.136402	.036000	.035419	.001178	.107640	.004967	.245152	.168852	.013860	.042623	.013024	.000717	.024512	.010498		

(d) 음절의 길이 L = 4

앞음절 초성이 X_p일 때 바로 다음 음절의 초성이 X_n일 조건부 확률 p(x_n|x_p)를 표 10(a)에 보인다. 조건부 확률 p(ㄷ|ㄱ)과 p(ㅅ|ㄱ)는 각각 초성의 발생 확률 p(ㄷ)와 p(ㅅ)보다 크고 조건부 확률 p(ㄴ|ㄱ)와 p(ㅇ|ㄱ)는 각각 초성의 발생 확률 p(ㄴ)와 p(ㅇ)보다 작다. 조건부 확률 p(ㄷ|ㄷ)와 p(ㄷ|ㄷ)는 각각 초성의 발생 확률 p(ㄷ)과 p(ㄷ)보다 크고 조건부 확률 p(ㄱ|ㄷ)와 p(ㄴ|ㄷ)는 각각 초성의 발생 확률 p(ㄱ)와 p(ㄴ)보다 작다. 조건부 확률 p(ㄷ|ㄷ)와 p(ㅎ|ㄷ)는 각각 초성의 발생 확률 p(ㄷ)와 p(ㅎ)보다 크고 조건

부 확률 p(ㄱ|ㄷ)와 p(ㄷ|ㄷ)는 각각 초성의 발생 확률 p(ㄱ)와 p(ㄷ)보다 작다. 조건부 확률 p(ㄷ|ㅅ)는 초성의 발생 확률 p(ㄷ)보다 크고 조건부 확률 p(ㅇ|ㅅ)는 초성의 발생 확률 p(ㅇ)보다 작다. 조건부 확률 p(ㄷ|ㅎ)는 초성의 발생 확률 p(ㄷ)보다 매우 크다.

앞음절 중성이 Y_p일 때 바로 다음 음절의 중성이 Y_n일 조건부 확률 p(y_n|y_p)를 표 10(b)에 보인다. 조건부 확률 p(ㅏ|ㅏ)는 중성의 발생 확률 p(ㅏ)보다 크고 조건부 확률 p(ㅑ|ㅏ)는 중성의 발생 확률 p(ㅑ)보다 작다. 조건부 확률 p(ㅑ|ㅑ)는 중성의 발생 확률 p(ㅑ)

표 9. 단어내에서의 음절의 위치에 따른 조건부 확률 p(x|y)
Table 9. Conditional probability p(x|y) as a function of syllable's location in a word.

Table with 21 columns (ㄱ to ㅎ) and 21 rows (ㄱ to ㅎ) showing conditional probabilities for syllable transitions. The table is a lower triangular matrix where the diagonal elements are 1.0.

(b) 두번째 음절

Table with 21 columns (ㄱ to ㅎ) and 21 rows (ㄱ to ㅎ) showing conditional probabilities for the second syllable transitions. The table is a lower triangular matrix.

(c) 마지막 음절

표 10. 음절간 초성, 중성, 종성 단위의 조건부 확률
Table 10. Conditional probabilities for a chosong, a jungseong, and a jongseong between syllables.

Table with 21 columns (ㄱ to ㅎ) and 21 rows (ㄱ to ㅎ) showing conditional probabilities between syllables. The table is a lower triangular matrix.

(a) 조건부 확률 p(xn|xp)

표 10. 음절간 조성, 중성, 중성 단위의 조건부 확률

Table 10. Conditional probabilities for a chosong, a jungseong, and a jongseong between syllables.

Table with 19 columns and 19 rows of conditional probabilities for syllable transitions. Values range from 0.00000 to 0.14115.

(b) 조건부 확률 p(y_n | y_p)

Table with 26 columns and 26 rows of conditional probabilities for syllable transitions. Values range from 0.00000 to 0.04119.

(c) 조건부 확률 p(z_n | z_p)

Table with 19 columns and 19 rows of conditional probabilities for syllable transitions. Values range from 0.00000 to 0.06560.

(d) 조건부 확률 p(x_n | z_p)

다 크다. 조건부 확률 $p(L|L)$ 는 중성의 발생확률 $p(L)$ 보다 크고 조건부 확률 $p(\text{공백소}|L)$ 는 중성의 발생확률 $p(\text{공백소})$ 보다 작다. 조건부 확률 $p(\text{공백소}|\text{공백소})$ 는 중성의 발생확률 $p(\text{공백소})$ 보다 크고 조건부 확률 $p(L|\text{공백소})$ 와 $p(\text{리}|\text{공백소})$ 는 각각 중성의 발생확률 $p(L)$ 와 $p(\text{리})$ 보다 작다. 조건부 확률 $p(\text{공백소}|리)$ 는 중성의 발생확률 $p(\text{공백소})$ 보다 크고 조건부 확률 $p(\text{리}|리)$ 는 중성의 발생확률 $p(\text{리})$ 보다 작다.

앞음절 중성이 Z_p 일 때 바로 다음 음절의 초성이 X_n 일 조건부 확률 $p(x_n|z_p)$ 를 표 10(d)에 보인다. 조건부 확률 $p(\text{리}|\text{공백소})$ 와 $p(\text{리}|\text{공백소})$ 각각 초성의 발생확률 $p(\text{리})$ 와 $p(\text{리})$ 보다 크고 조건부 확률 $p(\text{리}|\text{공백소})$ 와 $p(\text{리}|\text{공백소})$ 는 각각 초성의 발생확률 $p(\text{리})$ 와 $p(\text{리})$ 보다 작다. 조건부 확률 $p(\text{리}|리)$ 와 $p(\text{리}|리)$ 는 각각 초성의 발생확률 $p(\text{리})$ 와 $p(\text{리})$ 보다 작다. 조건부 확률 $p(\text{리}|리)$ 는 0.97보다 크다.

IV. 결 론

한국어 다음절 단어에 있어서 빈도수와 종류에 따른 음절의 길이 분포를 계산하였다. 한국어 단어는 빈도수를 고려했을 때 1음절 단어의 발생확률이 가장 높고 단어의 종류로는 2음절 단어가 가장 많다. 전체 음절의 초성, 중성, 중성 단위의 조건부 확률을 계산하였다. 다음절 단어의 길이와 단어내에서의 음절의 위치에 따른 초성, 중성, 중성 단위의 발생확률과 조건부 확률을 계산하였다. 한국어 단어내에서 인접음절의 초성-초성, 중성-중성, 중성-중성, 중성-중성 사이의 음절간 조건부 확률을 계산하였다. 이 결과로부터 한국어 단어에는 음절간의 상관관계가 있음을 알 수 있다.

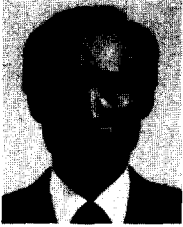
1956년 문교부에서 조사한 “우리말 말수 사용의 잦기 조사”를 한국어 단어의 빈도분포로 사용하였는데 한국어 단어의 발생빈도가 30년동안 근본적으로

는 변하지 않았다는 점을 고려하면 이 논문에서 계산된 한국어 다음절 단어의 초성, 중성, 중성 단위의 음절간 조건부 확률은 한국어 음성인식 및 합성, 자연언어처리, 암호법, 언어학, 음성학, 한글부호 표준화 연구등에 이용될 것으로 기대된다.

參 考 文 獻

- [1] C. Shannon, "Prediction and entropy of printed English," *Bell Syst., Tech. J.*, vol. 30, pp. 50-64, Jan. 1951.
- [2] 이재홍, 오상현, "한글 음절의 초성, 중성, 중성 단위의 발생확률, 엔트로피, 평균상호정보량" 전자공학회논문지, 제26권, 제 9 호, pp. 1-9, 1989년 9월.
- [3] 이주근, 최홍문, "한국어 음절의 entropy에 관한 연구" 전자공학회지, 제11권, 제3호, pp. 15-21, 1974년 6월.
- [4] 안수길, 안지환, "공백소를 포함한 한글 자소 발생 확률과 엔트로피" 전자공학회지, 제17권, 제 2호, pp. 23-28, 1980년 4월.
- [5] 남궁건, 한글 낱말의 발생빈도 분포의 엔트로피에 관한 연구. 석사학위논문, 서울대학교, 1979.
- [6] 과학기술처, 컴퓨터 관련 표준화 규격. 과학기술처, 1982.
- [7] 한국표준연구소, 한글 정보처리 표준화연구, 과학기술처, 1986.
- [8] 한국표준연구소, 한자부호 표준시안 작성을 위한 연구. 과학기술처, 1987.
- [9] 문교부, 우리말 말수 사용의 잦기 조사. 문교부 1956.
- [10] 정우상, "국민학교 교과서 어휘 연구" 국어연구소 연구보고서 제1집, pp. 651-872, 1987.
- [11] R. Manprino, "Printed Portuguese entropy-statistical calculation," *IEEE Trans. Inform, Theory*, vol. IT-16, p.122, Jan. 1970.

著者紹介



李在弘(正會員)

1953年 12月 7日生. 1976年 2月 서울대학교 전자공학과 졸업. 공학사학위 취득. 1978年 2月 서울대학교 대학원 전자공학과 졸업, 공학석사학위 취득. 1986年 8月 미국 미시간 대학교 전기공학과 및 컴퓨터 과학과 졸업. 공학박사학위 취득. 1978年~1981年 해군사관학교 전자공학과 교관, 전임 강사. 1987年 2月~현재 서울대학교 전자공학과 조교수, 부교수. 1991年 1月~현재 AT & T 벨 연구소 방문교수. 주관심분야는 채널코우딩, 확산대역 시스템, 디지털 통신이론, 위성통신, 이동통신, 음성합성 등임.



李在鶴(學生會員)

1967年 5月 1日生. 1989年 2月 서울대학교 전자공학과 졸업. 공학사학위 취득. 1991年 2月 서울대학교 대학원 전자공학과 졸업. 공학석사학위 취득. 1991年 3月~현재 서울대학교 대학원 전자공학과 대학원 박사과정 재학중. 주관심분야는 디지털 통신이론, 채널코우딩, 이동통신, 위성통신 등임.