

한글 음절의 초성, 중성, 종성 단위의 발생 확률, 엔트로피 및 평균상호정보량

(Entropy and Average Mutual Information for a 'Choseong', a 'Jungseong', and a 'Jongseong' of a Korean Syllable)

李在弘*, 吳相縣**

(Jae Hong Lee and Sang Hyun Oh)

要約

한글의 음절은 그 발생의 확률적 성질에 따라 확률변수로 간주된다. 한글의 음절은 확률변수들로 간주되는 초성, 중성, 종성으로 나누어진다. 이 세 확률변수에 대하여 모든 가능한 결합 확률과 조건부 확률이 한글의 음절별 누적 빈도수 통계로부터 계산된다. 세 확률변수에 대하여 모든 가능한 결합 엔트로피와 조건부 엔트로피가 결합확률과 조건부 확률로부터 계산된다. 또한 세 확률변수에 대하여 모든 가능한 평균상호정보량이 계산된다. 이들 엔트로피와 평균상호정보량으로부터 한글 음절의 성질을 정보이론적으로 분석한다.

Abstract

A Korean syllable is regarded as a random variable according to its probabilistic property in occurrence. A Korean syllable is divided into a 'choseong,' a 'jungseong,' and a 'jongseong' which are regarded as random variables. From the cumulative frequency of a Korean syllable all possible joint probabilities and conditional probabilities are computed for the three random variables. From the joint probabilities and the conditional probabilities all possible joint entropies and conditional entropies are computed for the three random variables. Also all possible average mutual informations are calculated for the three random variables. Average mutual information between two random variables has its biggest value between choseong and jungseong. Average mutual information between a random variable and other two random variables has its biggest value between jungseong and choseong-jongseong.

*正會員, **準會員, 서울大學校 電子工學科 (Dept. of Elec. Eng., Seoul Nat'l Univ.)

**準會員, (株) 金星社. (GoldStar Co., Ltd.)

接受日字: 1989年 5月 29日

(※ 이 논문은 1988년도 문교부 지원 한국 학술진흥재단의 자유공모과제 학술연구조성비에 의하여 연구되었음.)

I. 서론

언어는 인간의 의사전달의 가장 중요한 수단이다. 언어가 문자화 된 것이 글이다. 글은 기호들이 연속적으로 나열된 기호열로 볼 수 있다. 기호열은 정보를 내포하는데 기호열이 내포하는 정보량을 계량하려는 시도가 있어왔다. 즉, Shannon이 정보이론에 기초하여¹⁾, 영어 기호열의 정보량 척도로서 엔트로피(entropy)를 제안한 이래 몇몇 언어에 대하여 그 엔트로피

가 계산되었다.^[2,3] 한글에 있어서는 Shannon의 계산방법을 이용한 자소단위의 엔트로피에 관한 연구가 몇 차례 있었다.^[4-6] 이들 연구에서는 한글의 24자모를 영어의 26자의 알파벳에 대응되는 것으로 보고 엔트로피를 구하였다.

한글은 표음문자로서 모아쓰기를 한다. 모아쓰기를 함으로써 초성, 중성, 종성이 모여서 한 개의 도형적 습字 즉, ‘음절’(以下 ‘음절’로 표기)을 구성하는데 이러한 도형상 특성은 다른 언어에는 찾아볼 수 없는 것이다. 이러 특성을 고려할때, 초성, 중성 또는 종성을 단위로 한 발생의 불확실성(uncertainty) 즉, 엔트로피와 초성이 중성에 대하여 제공하는 정보량, 초성과 중성이 중성에 대하여 제공하는 정보량등의 계산분석 결과는 한국어 음성인식 및 합성, 자연언어처리, cryptography, 언어학, 음성학 등에 유용할 것이다. 그러나 이러한 점을 고려한 논문은 몇 년전까지는 전무하였으며 최근 부분적인 연구 결과가 발표된 바 있다.^[7]

한글의 엔트로피 계산에는 한글의 낱말 또는 음절 단위의 빈도 분포가 필요하다. 1945년 이래 한글의 낱말 또는 음절 단위의 빈도 분포의 체계적 조사는 빈약하였다.^[8,9] 조사의 내용에 있어서도 다른 언어에 대한 보다 체계적인 분포 조사와 통계적 분석에 비하여 빈약하였다.^[10] 한편, 한글사전의 표제어에 나타난 음소의 빈도분포를 구한 연구가 있었으나,^[11] 사전에 나타난 각 표제어의 사용빈도가 서로 다르기 때문에 사전 표제어에 나타난 음소의 빈도분포와 실제 사용되는 음소의 빈도분포는 차이가 있을 것이다.

이 논문에서는 먼저 한글의 음절별 발생빈도수 분포로부터 구한 음절별 발생 확률로부터 초성, 중성, 중성 단위의 조건부 발생확률(conditional probability)을 구하고 그로부터 초성, 중성, 중성 단위의 조건부 엔트로피를 구한다. 또한 초성, 중성, 중성 단위간의 평균상호정보량(average mutual information)을 구하고 비교한다. 이 결과로부터 초성, 중성 또는 중성이 무슨 자소인가를 아는 것이 같은 음절 내의 다른 단위(초성, 중성, 종성)가 무슨 자소인가에 대하여 얼마만한 정보를 제공하는지를 분석한다.

II. 한글의 엔트로피와 정보량

한글의 도형적 습字(이하 ‘음절’)를 구성하는 초성, 중성, 종성은 각각 가능한 자소집합으로부터 어떤 확률로 발생하는 확률변수(random variable)로 간주할 수 있다. 초성, 중성, 종성을 표시하는 확률변수를 각각 X, Y, Z라 정의하고 그 가능한 자소의 집합 즉, 표본 공간(sample space)을 각각 A_x, A_y, A_z

초성 x의 발생확률, 중성 y의 발생확률, 종성 z의 발생확률을 각각 $p(x), p(y), p(z)$ 라 하자. 그러면 A_x, A_y, A_z 는 다음과 같이 주어진다.

$$A_x = \{ㄱ, ㅋ, ㆁ, ㄷ, ㅌ, ㄹ, ㅁ, ㅂ, ㅅ, ㅆ, ㅈ, ㅊ, ㅌ, ㅍ, ㅎ\}$$

$$A_y = \{나, ㄹ, ㅈ, ㅊ, ㄷ, ㅌ, ㅋ, ㆁ, ㄱ, ㆁ, ㄷ, ㅌ, ㄹ, ㅁ, ㅂ, ㅅ, ㅆ, ㅈ, ㅊ, ㅌ, ㅍ, ㅎ\}$$

$$A_z = \{공백소, ㄱ, ㅋ, ㆁ, ㄷ, ㅌ, ㄹ, ㅁ, ㅂ, ㅅ, ㅆ, ㅈ, ㅊ, ㅌ, ㅍ, ㅎ\}$$

공백소(blank)는 자소가 결여된 상태 즉, 확률변수의 outcome이 null인 상태를 표시하며 중성에만 발생할 수 있고 초성, 중성에는 발생하지 않는다. 초성과 중성에 공통으로 사용되는 자소의 집합은 다음과 같이 주어진다.

$$A_x \cap A_z = \{ㄱ, ㅋ, ㆁ, ㄷ, ㅌ, ㄹ, ㅁ, ㅂ, ㅅ, ㅆ, ㅈ, ㅊ, ㅌ, ㅍ, ㅎ\}$$

한 음절 내의 초성, 중성, 중성간의 발생의 상관관계를 분석하기 위하여 먼저 모든 가능한 엔트로피, 결합 엔트로피, 조건부 엔트로피를 구한다. 엔트로피의 유형에는 초성의 엔트로피 $H(X)$, 초성, 중성 (X, Y) 의 결합 엔트로피 $H(X, Y)$, 초성, 중성, 종성 (X, Y, Z) 즉, 음절의 엔트로피 $H(X, Y, Z)$, 중성 Y가 주어질 때 초성의 조건부 엔트로피 $H(X|Y)$, 중성 Z가 주어질 때 초성, 중성, (X, Y) 의 조건부 엔트로피 $H(X, Y|Z)$, 중성, 종성 (Y, Z) 가 주어질 때 초성의 조건부 엔트로피 $H(X|Y, Z)$ 가 있으며 각각 다음의 식으로 구해진다.^[11,12]

$$H(X) = -\sum_x p(x) \log p(x) \tag{1}$$

여기서 $p(x)$ 는 초성 x의 발생확률이다.

$$H(X, Y) = -\sum_{x,y} p(x, y) \log p(x, y) \tag{2}$$

여기서 $p(x, y)$ 는 초성, 중성 (x, y) 의 결합확률이다.

$$H(X, Y, Z) = -\sum_{x,y,z} p(x, y, z) \log p(x, y, z) \tag{3}$$

여기서 $p(x, y, z)$ 는 초성, 중성, 종성 즉, 음절 (x, y, z) 의 결합확률이다.

$$H(X|Y) = -\sum_{x,y} p(x, y) \log p(x|y) \tag{4}$$

여기서 $p(x|y)$ 는 중성이 y일 때 초성이 x일 조건부 확률이다.

$$H(X, Y|Z) = -\sum_{x,y,z} p(x, y, z) \log p(x, y|z) \tag{5}$$

여기서 $p(x, y|z)$ 는 중성이 z일 때 초성이 x이고 중성이 y일 조건부 확률이다.

H(X|Y, Z) = -sum_{x,y,z} p(x,y,z) log p(x|y,z) (6)

여기서 p(x|y, z)는 중성이 y이고 중성이 z일때 초성이 x일 조건부 확률이다.

이러한 유형의 엔트로피들 간에 다음의 관계가 성립한다.^{12,13)}

H(X, Y) = H(X) + H(Y|X) (7)

H(X, Y, Z) = H(X, Y) + H(Z|X, Y) (8)

H(X, Y, Z) = H(X) + H(Y, Z|X) (9)

이 엔트로피로부터 평균상호정보량이 다음과 같이 구해진다.

I(X; Y) = H(X) - H(X|Y) (10)

I(X, Y; Z) = H(X, Y) - H(X, Y|Z) (11)

I(X; Y, Z) = H(X) - H(X|Y, Z) (12)

평균상호정보량간에 다음의 관계가 성립한다.^{12, 13)}

I(X; Y) = I(Y; X) (13)

I(X; Y, Z) = I(Y, Z; X) (14)

III. 계산 및 결과

초성, 중성, 종성의 발생 확률, 결합 확률 및 조건부 확률을 1985년 유 재원이 조사한 1,547개 한글 음절의 빈도수¹⁹⁾로부터 계산하고, 결합 확률과 조건부 확률로부터 초성, 중성, 종성의 결합 엔트로피와 조건부 엔트로피 및 평균 상호 정보량을 계산한다.

1. 초성, 중성, 종성의 발생 확률, 결합 확률과 조건부 확률

음절(x, y, z)의 발생 확률 p(x, y, z)는 '전체 음절별 누적 빈도수'로부터 계산된다. p(x, y, z)의 계산값은 음절(x, y, z)에 따른 데이터 양이 방대하여 여기에 신지 않는다.

음절별 발생빈도수 분포 통계로부터 음절(x, y, z)의 발생 확률 p(x, y, z)가 계산되고, p(x, y, z)로부터 초성의 발생 확률 p(x)가 다음 식으로 주어진다.

p(x) = sum_{y,z} p(x, y, z) (15)

중성의 발생 확률 p(y)와 종성의 발생 확률 p(z)도 같은 방법으로 구해진다. p(x), p(y)와 p(z)의 계산값을 표 1에 보인다.

표 1. 초성, 중성, 종성의 발생 확률 Table 1. Probabilities of a choseong, a jungseong, and a jongseong.

Table with 2 rows and 19 columns. Row 1: x, ㄱ, ㅋ, ㆁ, ㄷ, ㄸ, ㄴ, ㄹ, ㄲ, ㄳ, ㅃ, ㅅ, ㅆ, ㅈ, ㅉ, ㅊ, ㅊ, ㅌ, ㅍ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ, ㅡ, ㅣ, ㅞ. Row 2: p(x), .12519, .03043, .04511, .07900, .02138, .09746, .06384, .07112, .00986, .07769, .01289, .11429, .09593, .01419, .03107, .00980, .01779, .00986, .06365

(a) 초성의 발생 확률 p(x)

Table with 2 rows and 19 columns. Row 1: y, ㅏ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ, ㅡ, ㅣ, ㅞ, ㅟ, ㅠ, ㅡ, ㅢ, ㅣ, ㅤ, ㅥ, ㅦ, ㅧ, ㅨ, ㅩ, ㅪ, ㅫ, ㅬ, ㅭ, ㅮ, ㅯ, ㅰ, ㅱ, ㅲ, ㅳ, ㅴ, ㅵ, ㅶ, ㅷ, ㅸ, ㅹ, ㅺ, ㅻ, ㅼ, ㅽ, ㅾ, ㅿ, ㅿ, ㅿ, ㅣ. Row 2: p(y), .24227, .05231, .00910, .00009, .13293, .01872, .02091, .00156, .09486, .00609, .00170, .00762, .00266, .12027, .00209, .00063, .01110, .00131, .08739, .00192, .18468

(b) 중성의 발생 확률 p(y)

Table with 2 rows and 19 columns. Row 1: z, ㅏ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ, ㅡ, ㅣ, ㅞ, ㅟ, ㅠ, ㅡ, ㅢ, ㅣ, ㅤ, ㅥ, ㅦ, ㅧ, ㅨ, ㅩ, ㅪ, ㅫ, ㅬ, ㅭ, ㅮ, ㅯ, ㅰ, ㅱ, ㅲ, ㅳ, ㅴ, ㅵ, ㅶ, ㅷ, ㅸ, ㅹ, ㅺ, ㅻ, ㅼ, ㅽ, ㅾ, ㅿ, ㅿ, ㅿ, ㅣ. Row 2: p(z), .50180, .08148, .00125, .00027, .06531, .00051, .00056, .00400, .12396, .00228, .00032, .00080, .00005, .00012, .00002, .00044, .05238, .02106, .00192, .03148, .00035, .08990, .00473, .00346, .00019, .00707, .00363, .00221

(c) 종성의 발생 확률 p(z)

음절(x, y, z)의 발생확률 p(x, y, z)로부터 초성 x와 중성 y의 결합확률 p(x, y)가 다음 식으로 주어진다.

$$p(x, y) = \sum_z p(x, y, z) \tag{16}$$

초성 x와 중성 z의 결합확률 p(x, z)와 중성 y와 중성 z의 결합확률 p(y, z)도 같은 방법으로 구해진다. p(x, y), p(x, z)와 p(y, z)의 계산값은 여기에 실지 않는다.

음절(x, y, z)의 발생확률 p(x, y, z)로부터 중성이 y일 때 초성이 x일 조건부 확률 p(x|y)가 다음 식으로 주어진다.

$$p(x|y) = \frac{p(x, y)}{p(y)} = \frac{\sum_z p(x, y, z)}{\sum_x \sum_z p(x, y, z)} \tag{17}$$

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)와 중성이 y일 때 중성이 z일 조건부 확률 p(z|y)도 같은 방법으로 구해진다. p(x|y), p(x|z), p(y|x), p(y|z), p(z|x)와 p(z|y)의 계산값을 각각 표 2에 보인다.

음절(x, y, z)의 발생확률 p(x, y, z)로부터 중성이 z일 때 초성이 x이고 중성이 y일 조건부 확률 p(x, y|z)가 다음 식으로 주어진다.

$$p(x, y|z) = \frac{p(x, y, z)}{p(z)} = \frac{p(x, y, z)}{\sum_{x, y} p(x, y, z)} \tag{18}$$

중성이 y일 때 초성이 x이고 중성이 z일 조건부 확률 p(x, z|y)와 초성이 x일 때 중성이 y이고 중성이 z일 조건부 확률 p(y, z|x)도 같은 방법으로 구해진다. p(x, y|z), p(x, z|y)와 p(y, z|x)의 계산값은 초성 x, 중성 y, 중성 z에 따른 데이터 양이 방대하여 여기에 실지 않는다.

음절(x, y, z)의 발생확률 p(x, y, z)로부터 중성이 y이고 중성이 z일 때 초성이 x일 조건부 확률 p(x|y, z)가 다음식으로 주어진다.

$$p(x|y, z) = \frac{p(x, y, z)}{p(y, z)} = \frac{p(x, y, z)}{\sum_x p(x, y, z)} \tag{19}$$

초성이 x이고 중성이 z일 때 중성이 y일 조건부 확률 p(y|x, z)와 초성이 x이고 중성이 y일 때 중성이 z일 조건부 확률 p(z|x, y)도 같은 방법으로 구해진다. p(x|y, z), p(y|x, z)와 p(z|x, y)의 계산값은 초성 x,

표 2. 초성, 중성, 중성의 조건부 확률 p(x|y)型

Table 2. Conditional probabilities for a choseong, a jungseong, and a jongseong p(x|y) type.

x \ y	ㄱ	ㄴ	ㄷ	ㄹ	ㅁ	ㅂ	ㅅ	ㅇ	ㅈ	ㅊ	ㅋ	ㆁ	ㆁ	ㆁ	ㆁ	ㆁ	ㆁ	ㆁ	ㆁ	ㆁ
ㄱ	.09321	.02278	.06781	.06854	.01959	.05080	.06998	.09230	.00945	.06881	.01683	.06249	.09941	.01700	.02505	.00714	.01707	.00945	.16977	
ㄴ	.10760	.03030	.08054	.15006	.03354	.06450	.08070	.09399	.02252	.13725	.00940	.02090	.03581	.01215	.03565	.00535	.01847	.02252	.04132	
ㄷ	.04007	.02516	.08947	-	-	.04680	.23299	.02237	.01491	-	-	.48183	-	-	-	.01678	-	.01491	.02703	
ㄹ	.36364	-	-	-	-	-	-	-	-	-	-	.63636	-	-	-	-	-	-	-	
ㅁ	.19622	.02583	.02615	.06964	.02857	.09062	.05408	.06919	.01282	.07276	.01333	.11504	.11090	.01588	.03539	.00829	.02315	.01282	.02066	
ㅂ	.08741	.02491	.03623	.09194	.04710	.21286	.07201	.03578	-	.12726	.00091	.07246	.10145	.00543	.02627	.00408	.01721	-	.03668	
ㅅ	.18119	.00567	.04824	-	-	.11958	.08837	.14552	.04621	.00122	-	.24970	.00203	-	.01216	.01621	-	.04621	.02513	
ㅇ	.41848	-	.00543	-	-	.20109	-	-	-	-	-	.32609	-	.00543	-	.01087	-	-	.01087	
ㅈ	.12331	.06268	.05409	.13674	.02087	.04639	.06000	.07871	.00546	.10674	.00734	.08185	.06197	.01854	.02669	.02266	.03949	.00546	.03143	
ㅊ	.25905	.04596	-	.00418	.00557	.00279	-	.00975	-	.00696	.00418	.28412	.02646	.01114	-	.01671	-	-	.32312	
ㅋ	.22500	.12000	-	.07000	.01000	-	-	-	-	.02500	.04500	.32500	.02500	-	-	.02500	.01000	-	.12000	
ㆁ	.05784	.04894	.01335	.20912	.00779	.01446	.00556	.01335	-	.23804	.00556	.21691	.04116	.00556	.00334	.00222	.02892	-	.08788	
ㆁ	.03185	-	.00955	-	-	.04777	.00637	-	.07006	.00955	-	.65605	-	-	-	-	-	-	.07006	
ㆁ	.09240	.03996	.03538	.07612	.02347	.04257	.15104	.12363	.01318	.07464	.00959	.08549	.10058	.01734	.03073	.00430	.01818	.01318	.01219	
ㆁ	.05285	.07317	.00407	.00813	-	-	.01626	.00407	-	.00407	-	.74797	-	-	-	.04472	-	-	.04472	
ㆁ	.05405	.29730	.01351	.06757	-	-	-	-	-	.05405	.37838	.05405	-	-	-	.01351	.02703	-	.04054	
ㆁ	.23511	.00992	.02386	.24656	.04885	.00076	.00076	-	.03511	-	.18626	.08779	-	.01069	.02824	.01832	-	-	.06794	
ㆁ	.01290	-	.05806	-	-	.07097	-	-	.07742	-	.67097	-	-	-	-	-	-	-	.07742	
ㆁ	.20652	.05898	.05199	.17664	.04792	.15191	.00204	.00058	.00146	.08798	.01630	.10069	.00621	.00126	.00407	.02406	.02716	.00146	.03250	
ㆁ	-	-	.12335	-	.04846	-	-	-	-	.05727	.29075	-	-	-	-	-	-	.00881	.47137	
ㆁ	.09924	.01019	.02084	.00812	.00252	.21316	.03323	.04434	.00578	.06885	.01611	.18094	.18012	.01708	.05935	.00413	.00280	.00578	.02217	

(a) 조건부 확률 p(x|y)

중성 y, 종성 z에 따른 데이터 양이 방대하여 여기에 실지 않는다.

2. 초성, 중성, 종성의 엔트로피와 결합 엔트로피 초성 X의 엔트로피 $H(X)$ 가 식(1)과 식(15)으로 부터 구해지며 중성 Y의 엔트로피 $H(Y)$ 와 종성 Z의 엔트로피 $H(Z)$ 도 같은 방법으로 구해진다. $H(X), H(Y)$ 와 $H(Z)$ 의 계산값을 표 3에 보인다. 표 3에서 초성의 엔트로피가 가장 크고, 중성의 엔트로피가 가장 작음을 알 수 있다.

표 3. 초성, 중성, 종성의 엔트로피
Table 3. Entropy for a choseong, a jungseong, and a jongseong.

엔 트 로 피	bits
H (X)	3.897
H (Y)	3.118
H (Z)	2.518

초성-중성 (X, Y)의 결합 엔트로피 $H(X, Y)$ 가 식(2)와 식(16)로부터 구해지며 초성-중성(X, Z)의 결합 엔트로피 $H(X, Z)$ 와 중성, 종성 (Y, Z)의 결합 엔트로피 $H(Y, Z)$ 도 같은 방법으로 구해진다. $H(X, Y), H(X, Z)$ 와 $H(Y, Z)$ 의 계산값을 표 4에 보인다. 표 4에서 초성-중성의 엔트로피가 가장 크고, 중성-종성의 엔트로피가 가장 작음을 알 수 있다.

표 4. 초성, 중성, 종성의 결합 엔트로피 $H(X, Y)$ 형
Table 4. Joint entropy for a choseong, a jungseong, and a jongseong. $H(X, Y)$ type.

엔 트 로 피	bits
H (X, Y)	6.715
H (X, Z)	6.229
H (Y, Z)	5.459

음절(X, Y, Z)의 엔트로피 $H(X, Y, Z)$ 가 식(3)에 의해 구해지며 그 값은 8.679이다.

3. 초성, 중성, 종성의 조건부 엔트로피
중성 Y가 주어질 때 초성 X의 조건부 엔트로피 $H(X|Y)$ 가 식(4)와 (17)로부터 구해진다. 종성 Z가 주어질 때 초성 X의 조건부 엔트로피 $H(X|Z)$, 초성 X

가 주어질 때 중성 Y의 조건부 엔트로피 $H(Y|X)$, 종성 Z가 주어질 때 중성 Y의 조건부 엔트로피 $H(Y|Z)$, 초성 X가 주어질 때 종성 Z의 조건부 엔트로피 $H(Z|X)$ 와 중성 Y가 주어질 때 종성 Z의 조건부 엔트로피 $H(Z|Y)$ 도 같은 방법으로 구해진다. $H(X|Y), H(X|Z), H(Y|X), H(Y|Z), H(Z|X)$ 와 $H(Z|Y)$ 의 계산값을 표 5에 보인다. 표 5에서 중성이 주어질 때 초성의 조건부 엔트로피가 가장 크고 중성이 주어질 때 종성의 조건부 엔트로피가 가장 작음을 알 수 있다.

표 5. 초성, 중성, 종성의 조건부 엔트로피 $H(X|Y)$ 형
Table 5. Conditional entropy for a choseong, a jungseong, and a jongseong $H(X|Y)$ type.

엔 트 로 피	bits
H (X Y)	3.598
H (X Z)	3.712
H (Y X)	2.836
H (Y Z)	2.947
H (Z X)	2.350
H (Z Y)	2.341

중성 Z가 주어질 때 초성-중성(X, Y)의 조건부 엔트로피 $H(X, Y|Z)$ 가 식(5)와 (18)로부터 구해진다. 중성 Y가 주어질 때 초성-중성 (X, Z)의 조건부 엔트로피 $H(X, Z|Y)$ 와 초성 X가 주어질 때 중성-종성 (Y, Z)의 조건부 엔트로피 $H(Y, Z|X)$ 도 같은 방법으로 구해진다. $H(X, Y|Z), H(X, Z|Y)$ 와 $H(Y, Z|X)$ 의 계산값을 표 6에 보인다. 표 6에서 중성이 주어질 때 초성-중성의 조건부 엔트로피가 가장 크고 초성이 주어질 때 중성-종성의 조건부 엔트로피가 가장 작음을 알 수 있다.

표 6. 초성, 중성, 종성의 조건부 엔트로피 $H(X, Y|Z)$ 형
Table 6. Conditional entropy for a choseong, a jungseong, and a jongseong $H(X, Y|Z)$ type.

엔 트 로 피	bits
H (X, Y Z)	6.161
H (X, Z Y)	5.561
H (Y, Z X)	4.799

중성-중성 (Y, Z)가 주어질 때 초성 X의 조건부 엔트로피 $H(X|Y, Z)$ 가 식(6)과 (19)로부터 구해지며 초성-중성(X, Z)가 주어질 때 중성 Y의 조건부 엔트로피 $H(Y|X, Z)$ 와 초성-중성 (X, Y)가 주어질 때 중성 Z의 조건부 엔트로피 $H(Z|X, Y)$ 도 같은 방법으로 구해진다. $H(X|Y, Z)$, $H(Y|X, Z)$ 와 $H(Z|X, Y)$ 의 계산값을 표 7 에 보인다. 표 7 에서 중성과 중성이 주어질 때 초성의 엔트로피가 가장 크고, 초성과 중성이 주어질 때 중성의 조건부 엔트로피가 가장 작음을 알 수 있다.

표 7. 초성, 중성, 중성의 조건부 엔트로피 $H(X|Y, Z)$ 型.

Table 7. Conditional entropy for a choseong, a jungseong, and a jongseong $H(X|Y, Z)$ type.

엔 트 로 피	bits
$H(X Y, Z)$	3.220
$H(Y X, Z)$	2.449
$H(Z X, Y)$	1.967

4. 초성, 중성, 중성간의 평균상호정보량

초성 X와 중성 Y간의 평균상호정보량 $I(X;Y)$ 이 식(10)으로부터 구해지며 초성 X와 중성 Z간의 평균상호정보량 $I(X;Z)$ 와 중성 Y와 중성 Z간의 평균상호정보량 $I(Y;Z)$ 도 같은 방법으로 구해진다. $I(X;Y)$, $I(X;Z)$ 와 $I(Y;Z)$ 의 계산값을 표 8 에 보인다. 표 8 에서 초성과 중성간의 평균상호정보량이 가장 크고, 초성과 중성 간의 평균상호정보량이 가장 작음을 알 수 있다.

이것은 한 음절 내의 초성이 무슨 자소인가를 아는 것이 같은 음절 내의 중성이 무슨 자소인가에 대해

표 8. 초성, 중성, 중성간의 평균상호정보량 $I(X;Y)$ 型.

Table 8. Average mutual information for a choseong, a jungseong, and a jongseong $I(X;Y)$ type.

평균상호정보량	bits
$I(X;Y)$	0.282
$I(X;Z)$	0.168
$I(Y;Z)$	0.177

여 비교적 큰 정보를 제공하는 반면, 초성이 무슨 자소인가를 아는 것은 같은 음절 내의 중성이 무슨 자소인가에 대하여 비교적 작은 정보를 제공함을 뜻한다.

초성 X와 중성-중성 (Y, Z)간의 평균상호정보량 $I(X;Y, Z)$ 가 식(11)로부터 구해지며 중성 Y와 초성-중성 (X, Z) 간의 상호평균정보량 $I(Y;X, Z)$ 와 중성 Z와 초성-중성 (X, Y)간의 평균상호정보량 $I(Z;X, Y)$ 도 같은 방법으로 구해진다. $I(X;Y, Z)$, $I(Y;X, Z)$ 와 $I(Z;X, Y)$ 의 계산값을 표 9 에 보인다. 표 9 에서 중성과 초성-중성 간의 평균상호정보량이 가장 크고, 중성과 초성-중성간의 평균상호정보량이 가장 작음을 알 수 있다.

이것은 한 음절 내의 중성이 무슨 자소인가를 아는 것이 같은 음절 내의 초성-중성이 무슨 자소들인가에 대하여 비교적 큰 정보를 제공하는 반면, 중성이 무슨 자소인가를 아는 것은 같은 음절 내의 초성-중성이 무슨 자소들인가에 대하여 비교적 작은 정보를 제공함을 뜻한다.

표 9. 초성, 중성, 중성간의 평균상호정보량 $I(X;Y, Z)$ 型.

Table 9. Average mutual information for a choseong, a jungseong, and a jongseong $I(X;Y, Z)$ type.

평균상호정보량	bits
$I(X;Y, Z)$	0.659
$I(Y;X, Z)$	0.668
$I(Z;X, Y)$	0.554

IV. 결 론

한글의 음절을 초성, 중성, 중성 단위로 나누고 각각 확률변수로 간주하고 음절의 발생확률을 음절별 발생 빈도수 분포로부터 계산하였다. 음절의 발생확률로부터 초성, 중성, 중성의 발생확률, 결합확률과 조건부 확률을 계산하였다.

초성, 중성, 중성의 발생확률, 결합확률과 조건부 확률로부터 초성, 중성, 중성의 엔트로피를 계산하고 초성, 중성, 중성의 결합에 대한 결합 엔트로피를 계산하고 초성, 중성, 중성 간의 조건부 엔트로피를 계산하였다. 각 유형의 조건부 엔트로피에 있어서 중성이 주어질 때 초성의 조건부 엔트로피, 중성이 주어질 때 초성-중성의 조건부 엔트로피, 중성

-중성이 주어질 때 초성의 엔트로피가 각각 가장 크다.

또한 초성, 중성, 종성간의 평균상호정보량을 계산하였다. 두 확률변수 간의 평균상호정보량은 초성과 중성 간에 가장 크고 초성과 종성 간에 가장 작으며, 한 확률변수와 다른 두 확률변수간의 평균상호정보량은 중성과 초성-중성 간에 가장 크고 중성과 초성-중성 간에 가장 작다.

이 논문에 실린 초성, 중성, 종성의 발생확률, 조건부 확률, 엔트로피, 평균상호정보량의 계산값이 한국어 음성인식 및 합성, 자연언어처리, cryptography, 언어학, 음성학 등에 이용될 것을 기대한다.

參 考 文 獻

[1] C. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, pp. 379-413, 623-656, July and Oct. 1948.

[2] C. Shannon, "Prediction and entropy of printed English," *Bell Syst. Tech. J.*, vol. 30, pp. 50-64, Jan. 1951.

[3] R. Manfrino, "Printed Portuguese entropy - statistical calculation," *IEEE Trans. Inform. Theory*, vol. IT-16, p. 122, Jan. 1970.

[4] 이주근, 최홍문, "한국어 음절의 entropy에 관

한 연구," *전자공학회지*, 제 11 권, 제 3 호, pp. 15-21, 1974년 6 월.

[5] 안수길, 안지환, "공백소를 포함한 한글 자소 발생 확률과 엔트로피," *전자공학회지*, 제 17 권, 제 2 호, pp. 23-28, 1980년 4 월.

[6] 남궁건, 한글 낱말의 발생빈도 분포의 엔트로피에 관한 연구. 석사학위 논문, 서울대학교, 1979.

[7] 이재홍, 오상현, "한글의 초성, 중성, 종성 단위의 조건적 발생확률과 엔트로피," *한국통신학회 추계학술발표회 논문집*, pp. 53-56, 1987년 11월.

[8] 문교부, 우리말 말수 사용의 찾기 조사. 문교부, 1956.

[9] 유재원, 우리말 역순사전. 정음사, 1985.

[10] J.B. Carroll, P. Davies, and B. Richman, *Word Frequency Book*, Houghton Mifflin, 1971.

[11] 이용주, 김경태, 조철우, 이성구, "대용량 발음사전 표제어에 나타난 음소의 통계적 성질," *대한전자공학회 교환 및 통신연구회 합동학술발표회 논문집*, pp. 117-121, 1987년 11월.

[12] R. McEliece, *The Theory of Information and Coding*, Addison-Wesley, 1977.

[13] R. Gallager, *Information Theory and Reliable Communication*, Wiley, 1968.

著 者 紹 介

李 在 弘 (正會員) 第25卷 第3號 參照.
현재 서울대학교 전자공학과
조교수.

吳 相 縣 (準會員) 第25卷 第3號 參照.
서울대학교 전자공학과
석사학위 취득.
현재 (주)금성사 근무