

한글의 초성, 중성, 종성 단위의 조건적 발생확률과 엔트로피

이 재홍 오 상현[○]
 서울대학교 전자공학과

On the Conditional Probabilities and the Entropies of 'Choseong',
 'Jungseong', and 'Jongseong' of a Korean Syllable

Jae Hong Lee Sang Hyun Oh

Department of Electronics Engineering
 Seoul National University

Abstract

We regard 'choseong', 'jungseong', and 'jongseong' of a Korean syllable as random variables. We calculate the conditional probabilities and the entropies of three units. We obtain mutual information between units.

1. 서론

언어는 인간의 의사소통의 가장 중요한 수단으로서 언어가 문자화된 것이 글이다. 글은 기호들이 연속적으로 나열되어 형성된 기호열로 볼 수 있다. 이렇게 형성된 기호열은 정보를 내포하며 그 척도로서 엔트로피가 사용된다 [1], [2].

Shannon의 정보이론에 근거한 언어의 엔트로피가 몇몇 언어에 대하여 연구되었다 [3]. 한글에 있어서는 음절 단위와 자소 단위의 엔트로피에 관한 연구가 있었다 [4], [5], [6]. 또 엔트로피 연구의 자료가 되는 한글의 음절 및 낱말 단위의 빈도 분포가 조사된 바 있다 [7], [8].

한글은 다른 언어와는 달리 초성, 중성, 종성의 각 자소가 모여서 한 음절을 구성하고 이 음절이 모여서 낱말을 구성하는 특성을 가지고 있다. 이러한 특성에 비추어 볼 때 초성, 중성, 종성을 단위로 하여 엔트로피를 계산하는 것은 의미있는 일이다.

이 논문에서는 한글의 음절별 발생빈도 통계로부터 초성, 중성, 종성 단위의 조건적

발생확률 (conditional probability)을 구하고 그것을 사용하여 초성, 중성, 종성 단위의 조건적 엔트로피 (conditional entropy)를 구한다. 또한 초성, 중성, 종성 단위 간의 평균 상호정보량 (average mutual information)을 계산하고 비교한다..

2. 한글의 엔트로피와 정보량

한 음절을 이루는 초성, 중성, 종성은 각각 그 자소들의 집합에서 어떤 확률을 가지고 선택되는 불규칙 변수(random variable)로 볼 수 있다. 초성, 중성, 종성을 표시하는 불규칙 변수를 각각 X, Y, Z라 하고 그 표본공간(sample space) 즉 자소들의 집합을 각각 Ax, Ay, Az라 하자. 그러면 Ax, Ay, Az는 다음과 같이 주어진다.

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

$$A_y = \{ ㅏ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ, ㅡ, ㅚ, ㅜ, ㅞ, ㅟ, ㅛ, ㅝ, ㅟ, ㅜ, ㅞ, ㅟ, ㅛ, ㅝ, ㅟ, ㅜ, ㅞ, ㅟ, ㅛ, ㅝ, ㅟ \}$$

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

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

중성에는 공백소가 발생할 수 있는데, 공백소는 중성에 자음이 없는 경우를 표시한다. 그러나 초성에는 공백소가 발생하지 않는다.

한 음절 내의 단위 간의 발생의 상관관계를

알아보기 위하여 먼저 조건적 발생확률과 조건적 엔트로피를 구한다. 초성 X가 주어질 때 중성 Y의 조건적 엔트로피 $H(Y|X)$ 는 다음 식으로 주어진다.

$$H(Y|X) = -\sum_x p(x) H(Y|X=x) \\ = -\sum_x p(x) \sum_y p(y|x) \log p(y|x) \quad (1)$$

초성의 엔트로피 $H(X)$ 는 다음 식으로 주어진다.

$$H(X) = -\sum_x p(x) \log p(x) \quad (2)$$

식 (1)과 (2)로부터 초성과 중성 간의 평균 상호정보량 $I(X;Y)$ 는 다음 식으로 주어진다.

$$I(X;Y) = H(X) - H(X|Y) \quad (3)$$

같은 방법으로 중성과 중성, 초성과 중성간의 평균 상호정보량이 구해진다.

3. 계산 및 결과

초성, 중성, 중성 간의 조건적 발생확률을 계산하는데 사용된 자료는 '우리말 역순 사전'에 조사된 전체 음절별 누적 빈도수를 자료로 사용하였다[8]. 음절별 빈도수로부터 $p(x,y,z)$ 가 계산되고, $p(x,y,z)$ 로부터 $p(y|x)$ 가 다음 식에 의하여 구해진다 [9],[10].

$$p(y|x) = \frac{p(x,y)}{p(x)} = \frac{\sum_z p(x,y,z)}{\sum_{y,z} p(x,y,z)} \quad (4)$$

계산된 $p(y|x)$ 은 표 1과 같다. 같은 방법으로 계산된 $p(x|y)$, $p(z|y)$, $p(y|z)$, $p(z|x)$, $p(x|z)$ 는 각각 표 2에서 6까지와 같다.

초성이 특정한 값 x 로 주어졌을 때의 중성의 엔트로피 $H(Y|X=x)$ 와 이것을 x 에 대하여 평균한 조건적 엔트로피 $H(Y|X)$ 는 식 (1)로부터 구해지고 계산 결과는 표 7과 같다. 표 7에서 중성의 엔트로피는 초성이 'ㅇ'일 때 가장 크고 초성이 'ㅎ'일 때 가장 작음을 볼 수 있다. 같은 방법으로 계산된 $H(X|Y)$, $H(Z|Y)$, $H(Y|Z)$, $H(Z|X)$, $H(X|Z)$ 는 각각 표 8에서 12까지와 같다.

초성의 엔트로피 $H(X)$ 는 식 (2)로부터 구해지고 $H(Y)$, $H(Z)$ 도 같은 방법으로 구해지며 계산 결과는 표 13과 같다. 표 13에서 초성이 그 발생의

불확실성이 가장 크고 중성의 불확실성이 가장 작음을 볼 수 있다.

표 13. 각 단위의 엔트로피
table 13. entropy of the units

엔트로피	단위
$H(X)$	3.38180
$H(Y)$	3.11659
$H(Z)$	2.54731

표 14. 평균 상호정보량
table 14. mutual information

상호 정보량	단위
$I(X;Y)$	0.22504
$I(Y;Z)$	0.03654
$I(X;Z)$	0.11704

초성과 중성 간의 평균 상호정보량 $I(X;Y)$ 는 식 (3)으로부터 구해지고 $I(Y;Z)$, $I(X;Z)$ 도 같은 방법으로 구해지며 계산 결과는 표 14와 같다. 표 14에서 평균 상호정보량은 초성과 중성 간에 가장 크고 초성과 중성 간에 가장 작음을 볼 수 있다. 즉 초성이 중성에 대하여 공급하는 정보량이 가장 크고 초성이 중성에 대하여 공급하는 정보량이 가장 작음을 볼 수 있다.

4. 결론

한글의 음절을 초성, 중성, 중성의 단위로 나누고, 이들을 각각 불규칙 변수로 간주하였다. 각 단위의 조건적 발생확률을 계산하고, 이것을 사용하여 조건적 엔트로피를 계산하였다. 또 각 단위의 엔트로피를 계산하여, 초성의 엔트로피가 가장 크고 중성의 엔트로피가 가장 작음을 보였다. 같으로 각 단위들 간의 평균 상호정보량을 비교하여 초성과 중성의 상호 정보량이 가장 큼을 보였다.

인용 문헌

- [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] 이 주근, 최 흥문, "한국어 음절의 entropy에 관한 연구," 전자공학회지, 제11권, 제3호, pp. 15-21, 1974년 6월

- [5] 안수길, 안지환, "공백소를 포함한 한글 자소발생 확률과 엔트로피," 전자공학회지, 제17권, 제2호, pp. 23-28, 1980년 4월
- [6] 남궁건, 한글 낱말의 발생빈도 분포의 엔트로피에 관한 연구, 석사 학위 논문, 서울대학교, 1979.

- [7] 문교부, 우리말 말수 사용의 잦기 조사, 문교부 1956.
- [8] 유재현, 우리말 역순사전, 정음사 pp. 858-877, 1985
- [9] R. McEliece, The Theory of Information and Coding. Addison-Wesley, 1977.
- [10] R. Gallager, Information Theory and Reliable Communication. Wiley, 1968.

표

표 2. 조건적 발생확률 $p(x|y)$
table 2. conditional probability $p(x|y)$

X \ Y	ㅏ	ㅑ	ㅓ	ㅕ	ㅗ	ㅛ	ㅜ	ㅠ	ㅡ	ㅣ
ㅏ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅑ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅓ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅕ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅗ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅛ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅜ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅠ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅡ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅣ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
P(y)	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

표 1. 조건적 발생확률 $p(y|x)$
table 1. conditional probability $p(y|x)$

X \ Y	ㅏ	ㅑ	ㅓ	ㅕ	ㅗ	ㅛ	ㅜ	ㅠ	ㅡ	ㅣ
ㅏ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅑ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅓ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅕ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅗ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅛ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅜ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅠ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅡ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅣ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
P(x)	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

표 3. 조건적 발생확률 $p(z|y)$
table 3. conditional probability $p(z|y)$

X \ Y	ㅏ	ㅑ	ㅓ	ㅕ	ㅗ	ㅛ	ㅜ	ㅠ	ㅡ	ㅣ
ㅏ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅑ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅓ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅕ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅗ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅛ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅜ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅠ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅡ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
ㅣ	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
P(y)	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

