

Relation between Word Order of Languages and the Entropy of Mitochondrial DNA Haplogroups Distribution of the Speakers' Population

言語の語順とその話者集団のミトコンドリア DNA ハプログループ分布の エントロピーとの関係

Terumasa EHARA
江原暉将

Ehara NLP Research Laboratory
江原自然言語処理研究室
<http://www.ne.jp/asahi/eharate/eharate/>

1 Introduction

We are investigating the relation between the word order of languages and the speakers' thought pattern. We have approached it through suicide rate and homicide rate (Ehara, 2015, 2017) and diversity of Y-chromosome haplogroups of the speakers' population (Ehara, 2018).

Y-chromosome DNA reflects paternal lineage and we measure its diversity with the entropy of Y-chromosome haplogroups (YHg) distribution. Higher entropy value means the population is rich in diversity and it relates to peaceful thought pattern. Lower entropy value means the population is poor in diversity and it relates to warlike thought pattern (Sakitani, 2008).

In this paper we use another metrics. It is the entropy of mitochondrial DNA haplogroups (MtHg) distribution of the speakers' population (shortly "MtHg entropy"). Mitochondrial DNA reflects maternal lineage.

Our conjecture is that head final language speakers' population has higher entropy value and head initial language speakers' population has lower entropy value. Ehara (2018) showed the conjecture is right for the YHg entropy.

2 Data

Data for the word order features are obtained from the WALS online database (Dryer, 2013). We use two dominant word order features:

- Order of Object (O) and Verb (V),
- Order of Adjective (A) and Noun (N).

WALS online provides O and V feature values for 1519 languages and A and N feature values for 1366 languages. Head initial languages tend to have VO and NA word order and head final languages tend to have OV and AN word order.

MtHg data of populations are obtained from Eupedia's page of "Distribution of European mitochondrial DNA (mtDNA) haplogroups by region in percentage" (Eupedia.com, 2018). From this page, we get the MtHg data from 80 populations in Europe and its surrounding regions. We estimate languages spoken by these populations. We can recognize 62 languages in these data. We merge population level data to language level data using the sample size of the population and each MtHg's relative frequency of the population. More concretely, the merging method is as follows. We use the set of all MtHg: $G = \{L, HV, H, HV0 + V, J, T1, T2, U2, U3, U4, U5, U, K, I, W, X, Other\}$. Let P be the set of populations, L be the set of languages and $\mathcal{L} : P \rightarrow L$ be a mapping which corresponds each population to its speaking language. For each $l \in L$ and $g \in G$, the probability (relative frequency) of g for l : $p(l, g)$ is calculated by

$$p(l, g) = \frac{\sum_{p \in \mathcal{L}^{-1}(l)} q(p, g) \times s(p)}{\sum_{p \in \mathcal{L}^{-1}(l)} s(p)}$$

where $q(p, g)$ is the probability (relative frequency) of g for the population p and $s(p)$ is the sample size of the population p .

The MtHg entropy for each language l : $H(l)$

is calculated by

$$H(l) = - \sum_{g \in G} p(l, g) \log_{10} p(l, g)$$

Appendix 1 shows the base data used in the research sorted by the MtHg entropy.

3 Analysis and results

We conduct t-test for O and V word order groups and A and N word order groups. We discard “no dominant order” data. Results are shown in Table 1. For O and V case, the mean value of the entropy of the OV language group is significantly ($p < 1\%$) higher than the mean value of the entropy of the VO language group. On the other hand, for A and N case, mean values are not different significantly ($p > 1\%$). This results shows our conjecture is true by half for MtHg case. For YHg case, Ehara (2018) showed for both O and V case and A and N case, the mean values of the entropy values are significantly different. The mean value of the YHg entropy of OV language group is higher than the mean value of the YHg entropy of VO language group. The mean value of the YHg entropy of AN language group is higher than the mean value of the YHg entropy of NA language group.

Table 1: Results of t-test for the MtHg entropy

	OV	VO		AN	NA
n	17	36	n	38	19
mean	0.980	0.897	mean	0.913	0.945
unbiased var.	0.007	0.012	unbiased var.	0.011	0.013
t	2.748		t	-1.070	
p	0.00828		p	0.2892	

Figure 1 shows the ranking of MtHg entropy values for each language from lower to higher. Many OV type languages have higher entropy values and many VO type languages have lower entropy values.

4 Correlation between MtHg entropy and YHg entropy

In this section, we compare the YHg entropy and the MtHg entropy. For 42 languages in the Appendix 1, YHg entropy can be obtained from Ehara (2018). They are listed in Appendix 2 with MtHg entropy and YHg entropy.

Scattering graph of MtHg entropy and YHg entropy is shown in Figure 2.

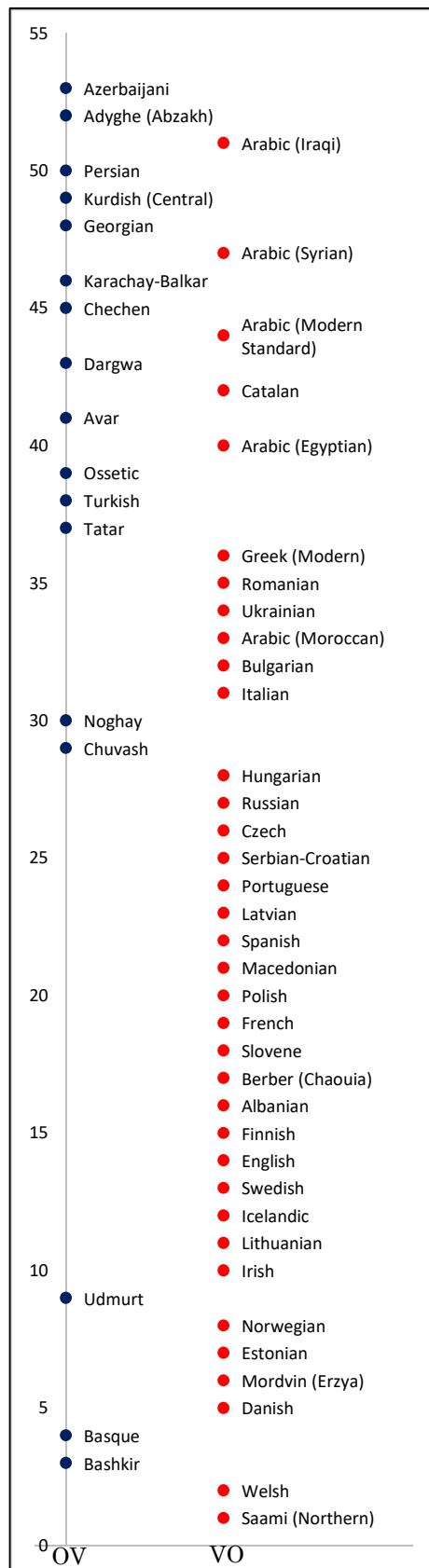


Figure 1: Ranking of MtHg entropy of each language (lower to higher)

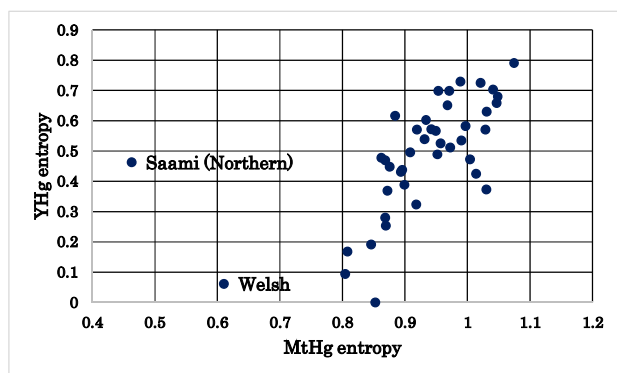


Figure 2: Scattering graph of MtHg entropy and YHg entropy for 42 languages

Correlation coefficient of MtHg entropy and YHg entropy is 0.597. Excluding the outliers (Saami (Northern) and Welsh), correlation coefficient is 0.742.

T-test results for MtHg entropy and YHg entropy by restricted 42 languages are shown in Table 2. MtHg entropy values divided by OV and VO are significantly different ($p < 5\%$). However, MtHg entropy values divided by AN and NA, YHg entropy values divided by OV and VO and YHg entropy values divided by AN and NA are not significantly different.

Table 2: Results of t-test for the MtHg entropy and the YHg entropy

(a) MtHg entropy

	OV	VO		AN	NA
n	14	24	n	29	13
mean	0.973	0.893	mean	0.922	0.924
standard dev.	0.091	0.122	standard dev.	0.111	0.121
t	2.154		t	-0.038	
p	0.03805		p	0.96985	

(b) YHg entropy

	OV	VO		AN	NA
n	14	24	n	29	13
mean	0.508	0.460	mean	0.495	0.455
standard dev.	0.223	0.179	standard dev.	0.169	0.232
t	0.736		t	0.631	
p	0.46621		p	0.53179	

5 Conclusion

Relation between word order (Object (O) / Verb (V) and Adjective (A) / Noun (N)) of languages and the entropy of Mitochondrial DNA haplogroups distribution of the speakers' population is examined. T-test results show OV word order language speakers' population tend to have higher entropy value than VO word order language speakers' population.

The languages used in this analysis are not

exhaustive. They are spoken in Europe and its surrounding regions. Analysis using languages all over the world is one of the remaining issues.

References

- Matthew S. Dryer. 2013. Order of Object and Verb and Order of Adjective and Noun, *In: Dryer, Matthew S. & Haspelmath, Martin (eds.) The World Atlas of Language Structures Online*. Leipzig: Max Planck Institute for Evolutionary Anthropology. (<http://wals.info/chapter/83> and 87, Accessed on 2015-3-23).
- Terumasa Ehara. 2015. Relation between Word Order Parameters and Suicide / Homicide Rates, *Journal of Yamanashi Eiwa College*, Vol.13, pages 9-29.
- Terumasa Ehara. 2017. Relation between the Word Order Characteristics and Suicide/Homicide Rates (6), *Proceedings of The 23th Annual Meeting of The Association for Natural Language Processing*, P4-4, pp.190-193.
- Terumasa Ehara. 2018. Relation between Word Order of Languages and the Entropy of Y-chromosome Haplogroup Distribution of the Speakers' Population, *Proceedings of The 24th Annual Meeting of The Association for Natural Language Processing*, P10-8, pp.1019-1022.
- Eupedia.com. 2018. Distribution of European mitochondrial DNA (mtDNA) haplogroups by region in percentage, *Eupedia*. (https://www.eupedia.com/europe/european_mtdna_haplogroups_frequency.shtml, Accessed on 2018-12-5).
- Mitsuru Sakitani. 2008. A long journey of Japanese populations revealed by DNA analysis, *Showado, Kyoto*, pages 37-39 (in Japanese).
- 崎谷満. 2008. DNAでたどる日本人10万年の旅, *昭和堂, 京都*, pages 37-39.

