



New Perspective on GWAS: East Asian Populations from the Viewpoint of Selection Pressure and Linear Algebra with AI

Masayuki Kanazawa*

Department of Human Sciences, ABO Center Inc., Tokyo, Japan

ABSTRACT

Genome Wide Association Studies (GWAS) are useful for comparing the characteristics of different human population groups. However, genomes can change rapidly over time when there is a strong selection pressure, such as a pandemic. The genetic information related to the immune system is thought to be very sensitive to such diseases. Therefore, it may be necessary to conduct not only the standard whole-genome GWAS but also a more detailed, chromosome-focused GWAS.

In this study, we compared chromosomes of immune system genes to those that are not thought to be related to the immune system, and analyzed GWAS results for SNPs in each chromosome to examine the differences. In order to keep the sample conditions as identical as possible, we limited the comparisons and the analyses to a few groups for which population movements were easy to interpret, and we also made sure the sample sizes were as close as possible. We selected a population of 403 East Asian people, consisting of 104 Japanese People in Tokyo (JPT), 103 Han Chinese People in Beijing (CHB) and 105 in Southern China (CHS), and 91 Korean People (KOR). PCA and Manhattan plot were used to analyze and compare the results.

Japanese, Chinese, and Korean populations formed distinctly different groups, with major differences observed. Validity of PCA and Manhattan plot was also discussed using Mahalanobis distance and AI.

Keywords: Genome; Chromosomes; Immune system; DNA sequencing

INTRODUCTION

Rapid advances in DNA sequencing technologies in recent years have led the International Human Genome Sequencing Consortium to announce the successful completion of the decoding of the human genome in 2003 [1]. Subsequently, a global whole-genome sequencing project (the 1000 Genomes Project) was initiated in 2008 with the goal of characterizing the genetic diversity of the world's population. As a result, genome information on more than 2,500 individuals worldwide is presently available as open data, including Japanese and Chinese individuals in East Asian populations [2]. More data are now available from other populations, including Koreans and Mongolians, which not been covered by the 1000 Genomes Project [3,4].

These data are being analysed mainly by linear algebraic methods, such as Genome Wide Association Studies (GWAS), and are beginning to reveal genomic characteristics in European populations

and details of human inter and intra-regional migrations from the past to the present [5,6]. Recently, studies have been conducted not only in Europe but also in East Asian populations such as Japanese, Chinese, and Korean populations [7-13].

Issues of previous studies

Most of the previous GWAS that compared human groups have focused on whole genomes or on genes associated with specific infectious diseases. However, not many studies have analyzed these associations at the chromosomal level, even though it is chromosomes where the genetic recombination is actually occurring. There are also several linear algebraic issues.

First, GWAS on whole genomes do not usually assume that genes change significantly over time due to a strong selection pressure; however, the genomes of the immune system, such as Human Leukocyte Antigen (HLA) system, are expected to change

Correspondence to: Masayuki Kanazawa, Department of Human Sciences, ABO Center Inc., Tokyo, Japan, E-mail: mkana@kfx.biglobe.ne.jp

Received: 24-Sep-2022, Manuscript No. RDT-22-18159; **Editor assigned:** 28-Sep-2022, PreQC No. RDT-22-18159 (PQ); **Reviewed:** 11-Oct-2022, PreQC No. RDT-22-18159; **Revised:** 18-Oct-2022, Manuscript No. RDT-22-18159 (R); **Published:** 25-Oct-2022, DOI: 10.35248/2329-6682.22.11.201.

Citation: Kanazawa M (2022) New Perspective on GWAS: East Asian Populations from the Viewpoint of Selection Pressure and Linear Algebra with AI. Gene Technol. 11:201

Copyright: © 2022 Kanazawa M. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

significantly in response to pandemics. Therefore, when the results of Principal Component Analysis (PCA) on Single Nucleotide Polymorphisms (SNPs) obtained by GWAS are used to estimate group associations back in time; it is not always guaranteed that the data are consistent from the past to the present. This may reduce the accuracy of the analysis results.

Second, in the previous studies, the volumes of sample sizes of the groups to be compared were often different; PCA on SNPs, which is often employed in GWAS, has the property of determining the scale (principal components) so as to ensure that the total variance of all samples is maximized. Therefore, the principal components that serve as the scale for comparison of all groups will be optimized for the group with the largest sample size, which may not necessarily guarantee total optimization. Thus, even if the same groups are compared, it is not uncommon among studies with different sample sizes to yield different results.

Aim of this study

To address the above issues, this study adopted the following alternative research approaches.

First, comparisons will be limited as much as possible to only a few groups for which population movement can be easily interpreted. Second, chromosomes with immune system genomes that are likely to change frequently will be selected, along with chromosomes that are not, and then the results of PCA of SNPs by GWAS for each chromosome will be compared.

Specifically, in the East Asian populations of Japanese, Chinese, and Korean subjects that are currently available as open data, chromosomes that contain genes closely related to infectious diseases and other chromosomes will be analyzed by PCA and Manhattan plot.

METHODOLOGY

Genome data used

Genome data from the aforementioned 1000 Genomes Project were used; they are publicly available as open data and include people in Tokyo, Japan (104 Japanese; JPT), Beijing, China (103 Han Chinese; CHB), and Southern China (105 Han Chinese; CHS) as East Asian populations [2]. For Koreans, we used data from the KPGP that is currently available as open data (91 Koreans; KOR) [3]. The number of people in each group is about the same (around 100), and the number of males and females is also about the same, which is expected to reduce PCA bias. Since these data are open data, anyone can use them freely.

The following advantages can also be expected from the data used in this study.

- 1) The geographical proximity and small number of groups make it relatively easy to verify the results of the analysis.
- 2) China has a long history, and events that may affect the genome (such as population movements, wars, and pandemics) are available in historical books and other written records, and can be compared with GWAS results.
- 3) Population sizes of these countries are relatively large (Japan 126 m, China 1,380 m, Korea 52 m) and group members are stable.

Target chromosomes

Based on the aforementioned reasons, the following chromosomes were selected for analysis, referring to the genome map created by the Japanese Ministry of Education, Culture, Sports, Science and Technology [14].

- 1) Chromosome 1, which is a common chromosome and is assumed to contain few genes related to the immune system.
- 2) Chromosome 6, which contains HLA genes, and is the center of the human immune system.
- 3) Chromosome 9, which contains the gene for ABO blood group, and has been studied in relation to many infectious diseases including COVID-19 [15,16].
- 4) Chromosome 12, which contains ALDH2 (Human Aldehyde Dehydrogenase 2), the "alcohol-sensitivity"; information encoded on this chromosome tends to be similar among people in Japan, southern China, and Korea, while also being considered to be distinct from other human groups [17].

Analysis methods

PCA was conducted for the above chromosomes, and the first and the second principal components were analyzed to calculate the Mahalanobis distances from JPT to CHB, CHS, and KOR. Similarly, analysis using Manhattan plot was conducted for each chromosome comparing the Japanese population to the others. The software used was plink v1.9 and 2.0 for PCA and R v4.1.3 for Manhattan plot. The alpha level was set at 0.05 and Bonferroni's correction was used.

Chromosome 6, which contains the genes of HLA, a major immune system, was divided into the first half, in which the genes of HLA are present, and the second half, in which they are absent. PCA was performed on each portion. The numbers of SNPs in the first half and the second half were set equal.

RESULTS

Chromosome 1

In the first and second principal components, where the differences were largest, the Japanese, Korean, and Chinese populations were clearly separated. The groups of Beijing (CHB) and southern China (CHS) also had their own characteristics, but as a whole they constituted one group (Figure 1). Manhattan plot showed no noticeable difference (Figure 2).

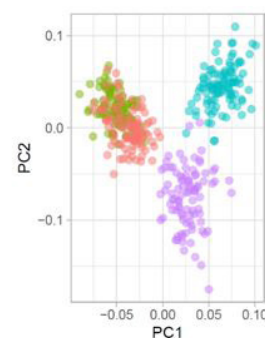
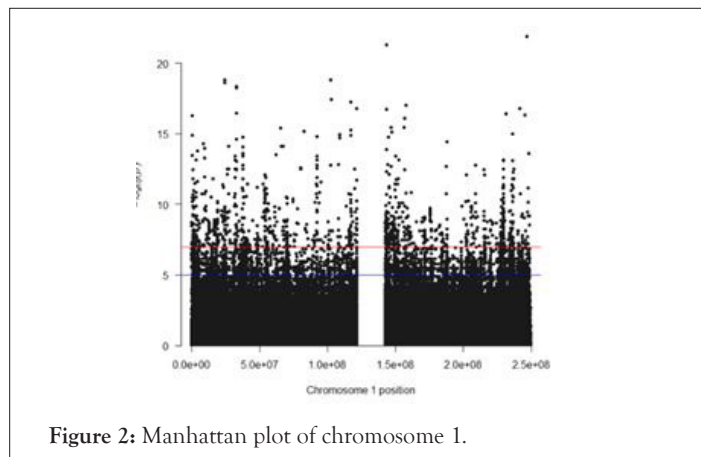
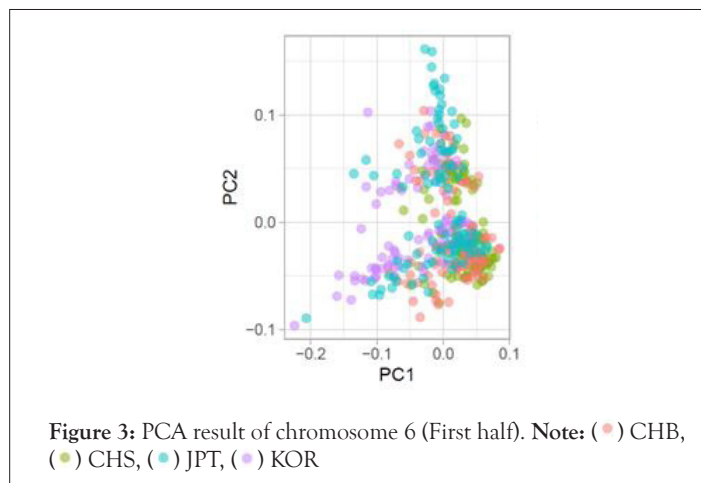


Figure 1: PCA result of chromosome 1. Note: (●) CHB, (●) CHS, (●) JPT, (●) KOR

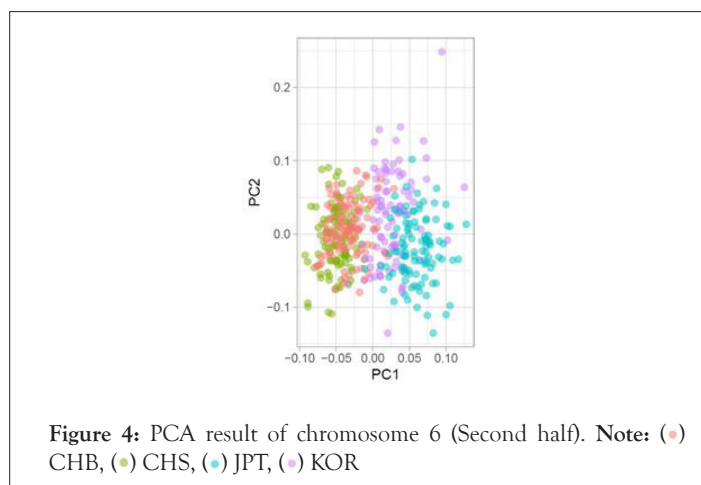


Chromosome 6

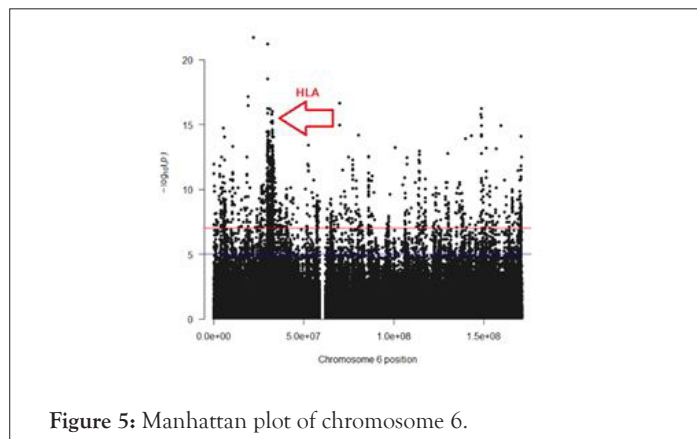
The results of the PCA of the first half of the study, where the HLA genes are located, were almost identical for Japanese, Chinese, and Korean populations, with individual differences more significant than group differences (Figure 3).



The second half of the chromosome showed the distinct characteristics of the Japanese, Chinese, and Korean populations, but the differences were smaller than those at chromosome 1 (Figure 4).

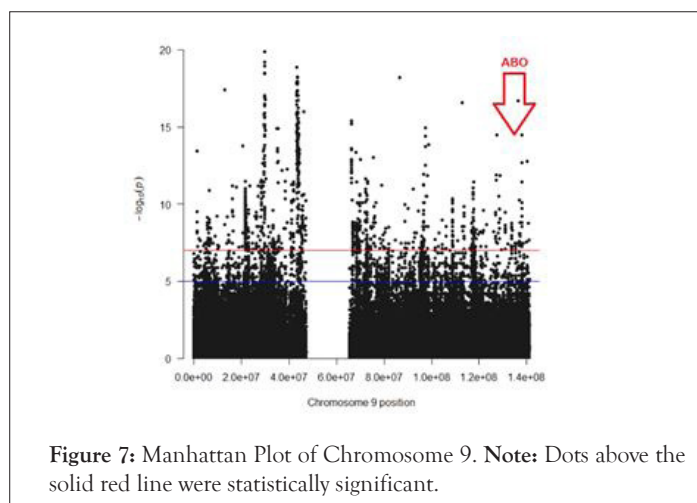
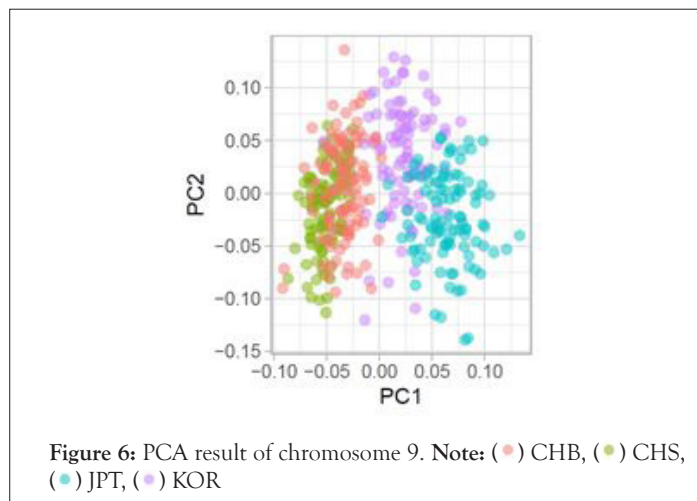


The Manhattan plot at chromosome 6 showed substantial differences in HLA positions, suggesting that the genes of HLA, a major immune system, had been significantly altered and mutated, resulting in extremely large differences between Japanese, Chinese, and Korean groups (Figure 5).



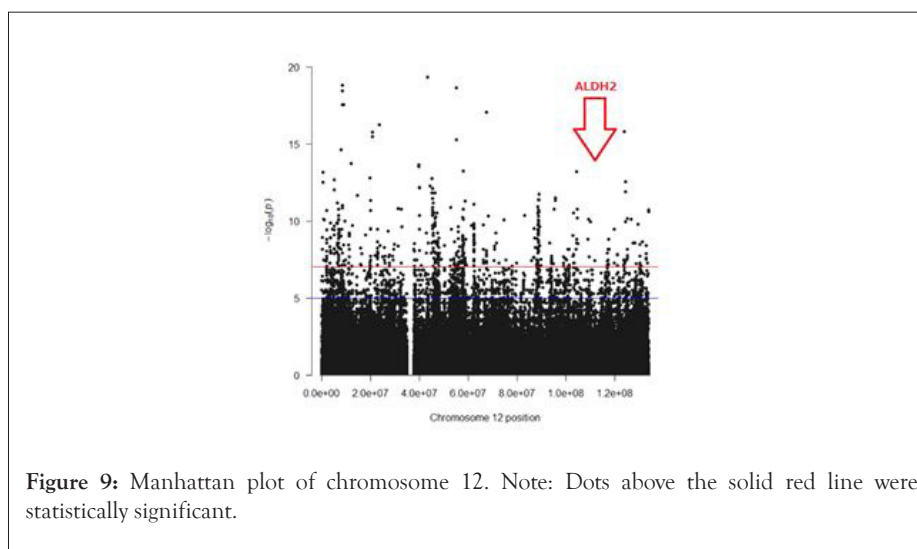
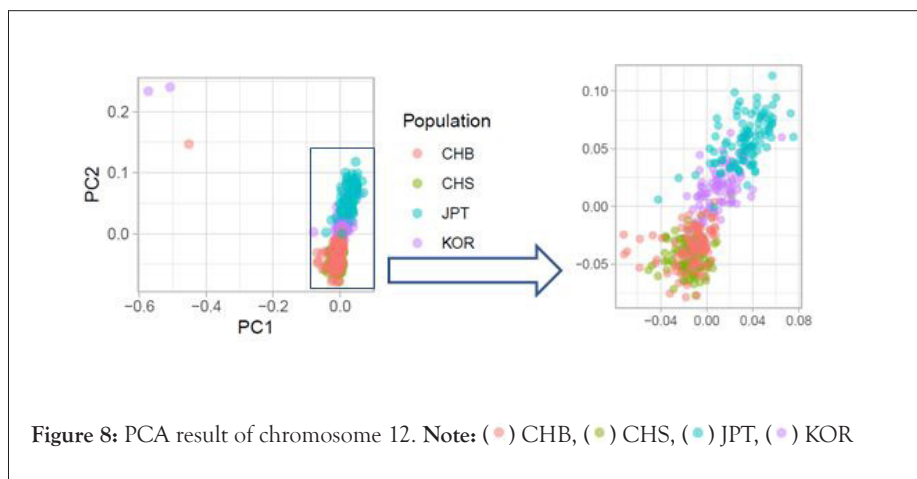
Chromosome 9

As with chromosome 1, the Japanese, Chinese, and Korean populations were separated (Figure 6), although the differences were smaller. Manhattan plot showed no noticeable difference in ABO position (Figure 7).



Chromosome 12

As with chromosome 1, the Japanese, Chinese, and Korean populations were separated. However, the first and second principal component figures showed that three individual subjects had values that were far from their groups (Figure 8). Manhattan plot showed no noticeable difference in ALDH2 position (Figure 9).



DISCUSSION

Estimation of selection pressure using linear algebra

Looking at the first and second principal components of PCAs, Mahalanobis distances, which indicate differences in human population groups, were considerably smaller in the first half of chromosome 6, which contains HLA, the core of the human immune system where natural selection pressure was strong, in comparison with other chromosomes (Table 1). In the other parts, such as the second half of chromosome 6, the Mahalanobis distance was not very different from that of chromosome 1. These values suggest that natural selection pressure affected the second half of chromosome 6 less than they affected the first half. In general, Mahalanobis distances from JPT (Tokyo, Japan) were also consistent with the physical distances from KOR (Korea), CHB (Beijing, China), and CHS (southern China), in that order.

As mentioned, the first half of chromosome 6, where the gene for HLA is located, was heavily altered and mutated in SNPs, resulting in extremely large differences among Japanese, Chinese, and Koreans. However, the DNA was very different; the PCA result

showed that Japanese, Chinese, and Koreans almost overlapped. This might mean very large individual differences.

Alternative approach using AI

In a different approach, we used AI, Microsoft Azure Machine Learning, to predict human population groups from the values of the first and the second principal components of PCAs. System parameters were set to defaults. The whole PCA data were divided into two parts, with 70% allocated for training and the remaining 30% for predicting. The result of the AI prediction is shown in Table 2. The larger the Mahalanobis distance the lower the accuracy. This result also suggests that natural selection pressure was strong.

The same dataset and AI were used to analyze the "genomic patterns" of each human population group. Okada et al. [17] suggested that the acetaldehyde decomposition gene (ALDH2) had the strongest natural selection pressure for Japanese people among all genomes. Since human chromosomes are diploid, each SNP has 0 to 2 mutations. Group differences in ALDH2 are shown in Table 3. The higher the number of mutations, the more sensitive to alcohol.

Table 1: Mahalanobis distances (1st and 2nd principal components) from JPT.

Population	Chromosome number				
	1	6 (1 st Half)	6 (2 nd Half)	9	12
CHB	6.722	0.673	3.859	4.333	5.140
CHS	7.365	0.953	4.464	4.946	5.477
KOR	5.753	0.717	1.759	2.307	2.508

Table 2: Result of human group prediction using AI.

Population	Chromosome number				
	1	6 (1 st Half)	6 (2 nd Half)	9	12
Accuracy	0.967	0.733	0.895	0.950	0.918
Algorithm	Gradient boosting	Random trees	Logistic Regression	SVM	SVM

Note: Each accuracy was the maximum value among several algorithms used.

Table 3: Number of ALDH2 mutations by human population group.

Population	Number of ALDH2 mutations			
	0	1	2	Total
JPT	59	40	5	104
	56.7%	38.5%	4.8%	100%
CHB	74	25	4	103
	71.8%	24.3%	3.9%	100%
CHS	57	39	9	105
	54.3%	37.1%	8.6%	100%
KOR	66	23	2	91
	72.5%	25.3%	2.2%	100%

We predicted human groups based on the pattern of mutation counts for all 100 SNPs, including ALDH2 and those before and after (POS 112211833-112263268). Similarly, the whole data were divided into two parts, with 70% allocated for training and the remaining 30% for predicting. The accuracy was 63.9% (algorithm: XG Boost Classifier), more than twice the 25% probability of a chance match to one of the four population groups (JPT, CHB, CHS and KOR). The result also suggests that natural selection pressure was strong.

Effect of selection pressure

From the results above, it can be inferred that natural selection by infectious diseases does not have a "positive" effect on specific genes, but rather a "negative" effect. In other words, genes change in the direction of larger diversity. This implies that individuals with genes that make them more susceptible to infectious diseases are more likely to go extinct, as there is a greater likelihood that they will decline rapidly since they will not be able to pass their genes on to the next generation. The large genetic variation probably means that it is advantageous for the survival of the species for the immune system to diversify, rather than the "survival of the fittest." This is similar to the situation with resistant bacteria or COVID-19; the latter is still mutating to escape vaccines.

It is estimated that Japan, being an island nation, has a relatively small population influx from outside compared to continental nations such as China and Korea. However, in light of the above, it may not always be a meaningful comparison to consider how much

of the ancient Japanese (Jomon) genes are inherited by modern Japanese people. Japan has been separated from the Asian continent since about 20,000 years ago, and from that point until about 3,000 years ago (the Jomon period), Japan was not an agricultural society. Later (the Yayoi period), with the arrival of paddy rice cultivation, the society transformed into an agricultural one, the staple food changed to rice, intense infectious diseases such as tuberculosis became prevalent, and the environment changed drastically [18]. Therefore, it is assumed that among the Jomon people, those individuals who could not cope with such changes died without children. Therefore, it is highly likely that the modern Japanese people who have the genes to cope with infectious diseases are considerably different from those of the Jomon people, who did not live in an agricultural society until 3,000 years ago.

Influence of paddy rice cultivation

The "low tolerance for alcohol" gene, which is said to have originated in the Yangtze River region, may be another example [19]. More than 6,000 years ago, many people began to gather and live near the Yangtze River floodplain, which was suitable for rice cultivation. At that time, because of a poor sanitary environment, food was often contaminated with harmful microbes and other substances that could cause infectious diseases. At such a time, alcoholic beverages made from rice were thought to be useful.

When people with a low tolerance for alcohol, or weak Acetaldehyde Decomposition Gene (ALDH2), drank alcoholic beverages, the level of acetaldehyde, a highly poisonous substance that cannot

be decomposed, would increase in his or her body. However, it appears that the poison might have also served as a drug that attacked harmful microbes. On the other hand, people without this weak gene had less acetaldehyde in their bodies and could not suppress those harmful microorganisms. Thus, people with the "low tolerance for alcohol" gene were more likely to survive and overcome infectious diseases.

In other words, it is possible that people in rice paddy farming areas felt selection pressure to develop a low tolerance for alcohol in order to protect themselves from infectious diseases. This "low tolerance for alcohol" gene, ALDH2, was eventually introduced to the Japanese archipelago along with the rice culture; over 40% of the present Japanese population has this gene. It is also thought that many ancient Japanese before rice paddy cultivation (Jomon people) did not have this gene.

Tuberculosis (TB) is another infectious disease thought to have arrived in Japan along with rice paddy cultivation [19]. People with blood types B and AB have the same type B antigen as Mycobacterium Tuberculosis, making it difficult for their immune system to function and making them susceptible to infection [20]. On the other hand, people with blood types A and O, which do not carry the type B antigen, are less susceptible to TB infection. In East Asia, paddy rice cultivation is prevalent in Japan, southern China, and Korea, where types A and O tend to be more common than types B and AB [21]. However, due to improved sanitary conditions in today's East Asia, it would be difficult to substantiate the above hypotheses.

Interpretation of PCA and Manhattan plot

According to the PCA and the Manhattan plot for each chromosome conducted in this study, the results seem to vary considerably depending on the selection pressure, human group selection methods, and sample sizes.

For example, on chromosome 6, which contains the HLA gene, the Mahalanobis distances for Japanese, Chinese, and Koreans are relatively close compared to other chromosomes according to PCA analysis (Table 1). However, Manhattan plot analysis results show that the differences between Japanese and Chinese or Japanese and Korean SNPs are larger than the other chromosomes (Figure 5), which were the exact opposite.

In addition, ALDH2, the "low tolerance for alcohol" gene commonly found in Japan, southern China, and Korea, is said to differ more from other human groups [14], but significant differences were not found in PCA or Manhattan plot (Figures 8 and 9), whereas differences were found in AI-based analysis (Table 3).

In light of the above, when making comparisons among human groups, sufficient attention should be paid not only to sample selection methods and sample sizes, but also to the linear algebraic nature of PCA and Manhattan plot. A comprehensive perspective will also be required when interpreting the results.

CONCLUSION

GWAS is useful for comparing characteristics of human groups. However, genomes may change rapidly over time when selection pressures, such as environmental changes, are strong. In particular, the immune system seems to be very sensitive to environmental changes. Therefore, there will be a need to perform GWAS not only on whole genomes, but also at the level of individual chromosomes when necessary.

Japanese, Chinese, and Korean people form distinctly different groups genetically. However, the sample size for this study is small (403 individuals), the targeted samples are limited to the East Asian population, and only a few basic methods were used for data analysis. Studies with a larger global dataset and methodological innovations will be needed to get one step closer to the truth.

REFERENCES

1. Human Genome Project Completed. National Human Genome Research Institute. 2003
2. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR. A global reference for human genetic variation. *Nature*. 2015;526(7571):68-74.
3. Jeon S, Bhak Y, Choi Y, Jeon Y, Kim S, Jang J et al. Korean genome project: 1094 Korean personal genomes with clinical information. *Science Advances*, 2020;6(22)peaz7835.
4. Bai H, Guo X, Narisu N, Lan T, Wu Q, Xing Y, et al. Whole-genome sequencing of 175 Mongolians uncovers population-specific genetic architecture and gene flow throughout North and East Asia. *Nature genetics*. 2018; 50(12):1696-704.
5. Tadaka S, Katsuo F, Ueki M, Kojima K, Makino S, Saito S, et al. 3.5 KJPNv2: An allele frequency panel of 3552 Japanese individuals including the X chromosome. *Human Genome Variation*. 2019;6(1):1-9.
6. Hofmann D, David Reich. *Who We Are and How We Got Here. Ancient DNA and the New Science of the Human Past* (Oxford: Oxford University Press, 2018, xxxi and 335pp, 28 illustr, pbk, ISBN 978-0-19-882126-7). *European Journal of Archaeology*. 2019;22(3):434-7.
7. Cao Y, Li L, Xu M, Feng Z, Sun X, Lu J et al. The ChinaMAP analytics of deep whole genome sequences in 10,588 individuals. *Cell research*. 2020;30(9):717-31.
8. Yoo SK, Kim CU, Kim HL, Kim S, Shin JY, et al. NARD: Whole-genome reference panel of 1779 Northeast Asians improves imputation accuracy of rare and low-frequency variants. *Genome medicine*. 2019; 11(1):1-0.
9. Robbeets M, Bouckaert R, Conte M, Saveljev A, Li T, An DI, et al. Triangulation supports agricultural spread of the Transeurasian languages. *Nature*. 2021;599(7886):616-621.
10. Jinan TA, Kanazawa-Kiriyama H, Saitou N. Human genetic diversity in the Japanese Archipelago: Dual structure and beyond. *Genes & Genetic Systems*. 2015;90(3):147-52.
11. Cooke NP, Mattiangeli V, Cassidy LM, Okazaki K, Stokes CA, Onbe, S, et al. Ancient genomics reveals tripartite origins of Japanese populations. *Science advances*. 2021;(7):38peabh2419.
12. Gelabert P, Blazyte A, Chang Y, Fernandes DM, Jeon S, Hong JG et al. Diverse northern Asian and Jomon-related genetic structure discovered among socially complex three kingdoms period Gaya region Koreans. *bioRxiv*. 2021.
13. Kanazawa M. New perspective on GWAS: East Asian populations from the viewpoint of selection pressure and Linear Algebra. *bioRxiv*. 2022.
14. Genome Map, Ministry of Education, Culture, Sports, Science and Technology Japan.
15. Hoiland RL, Fergusson NA, Mitra AR, Griesdale DE, Devine DV, Stukas S, et al. The association of ABO blood group with indices of disease severity and multiorgan dysfunction in COVID-19. *Blood advances*. 2020;4(20):4981-9.
16. Ellinghaus D, Degenhardt F, Bujanda L, Buti M, Albillos A, Invernizzi P, et al. Genomewide association study of severe COVID-19 with respiratory failure. *N Engl J Med*. 2020;383:1522-34.
17. Okada Y, Momozawa Y, Sakaue S, Kanai M, Ishigaki K, Akiyama M, et al. Deep whole-genome sequencing reveals recent selection signatures linked to evolution and disease risk of Japanese. *Nature communications*. 2018;9(1):1-0.

18. NHK (Japan Broadcasting Corporation). Evolutionary fate makes you want to drink: the unknown truth about alcoholic beverages.
19. Okazaki K, Takamuku H, Yonemoto S, Itahashi Y, Gakuhari T, Yoneda M, et al. A paleopathological approach to early human adaptation for wet-rice agriculture: The first case of Neolithic spinal tuberculosis at the Yangtze River Delta of China. *Int J Paleopathol.* 2019;1(24):236-44.
20. Oike Y, Kikuchi Y, Kushibiki H, Kudo M, Kobori K, Shintani K. Tuberculosis and ABO blood type. *Intern. Med.* 1954;42(11):835-838.
21. Mourant AE, Kopec AC, Domaniewska-Sobczak K. The distribution of the human blood groups and other polymorphisms (Monographs on Medical Genetics). 1976. Oxford: Oxford University Press.