



## Research Article

# Principal component analysis and grouping of sorghum (*Sorghum bicolor* L. Moench) gene pool for genetic diversity

D. Kavithamani<sup>1</sup>, A. Yuvaraja<sup>2</sup> and B. Selvi<sup>1</sup>

<sup>1</sup>Department of Millets, CPBG, TNAU, Coimbatore

<sup>2</sup>Agricultural College and Research Institute, Madurai

E-Mail: kavitharice@gmail.com

(Received: 18 Mar 2019; Revised: 04 Nov 2019; Accepted: 14 Nov 2019)

### Abstract

Principal component and hierarchical cluster analyses were carried out with eight quantitative traits in 100 germplasm accessions of sorghum (*Sorghum bicolor* L. Moench). Principal factor analysis identified three principal components which explained about 66.35% variability. PC 1 explained maximum variability of 31.05 % of total variation for morphological traits and PC 2 loaded with 19.66 % of total phenotypic variability and PC 3 had contributed 15.64 % of the total variation. PCA revealed that 100 seed weight, Plant height, leaf blade length and width were important traits depending upon their loading variables on a common principal axis. Sorghum germplasm accessions viz., IS 6312, IS 6238, MS 7885, AS 4242, AS 5763, IS 7270, AS 4669 and AS 7076 were identified as superior genotypes based on principal component trait analysis. Hierarchical cluster analysis emphasized on 100 accessions and grouped into eight clusters. Maximum of 32 genotypes were grouped in cluster III and showed more divergence of plant height related traits. Cluster number V and VIII had only one genotype each and both were to be superior for 100 g weight and length of panicle branches. Both PCA and clustering registered same level of variability between genotypes. Hence these superior accessions may further be utilized in breeding programmes for evolving sorghum varieties with high yield.

### Key words

Sorghum, Principal component analysis, Clustering, Genetic diversity

### Introduction

Sorghum (*Sorghum bicolor* (L.) Moench) is the fifth most important cereal crop worldwide after wheat, rice, maize and barley. It forms the most important dryland cereal crop for the semi-arid tropics together with maize and pearl millet. It is grown in at least 86 countries, in an area of 38 million hectares and with annual grain production of about 58 million tonnes. The average productivity reaches 1.5 tonnes per hectare (FAOSTAT, 2014). Sorghum is one of the most important crops grown for food and feed. Sorghum is a dietary staple for most of developing countries. Sorghum is the single most important cereal in the lowland areas because of its extreme resistance to water stress. *Sorghum bicolor* contains both cultivated and wild relative races, and it provides a substantial amount of genetic diversity for traits of agronomic importance to develop the crop's different variety of interest for plant breeders (Hart *et al.*, 2001 and Ali *et al.*, 2007). Sorghum has one of the largest crop germplasm collections, consisting of more than 42,000 accessions worldwide (Kassahun TESFAYE, 2017). Genetic diversity is one of the important cornerstones of crop improvement. Diversity provides the raw materials from which desirable or favourable alleles for improved agronomic traits of interest can

be selected. Subsequently, breeding for improved varieties are facilitated through the incorporation of new alleles into well adapted or elite lines. The assembly through collections, preservation and distribution of germplasm representing large diversity for each crop species then, is of utmost significance for plant breeding. A critical and challenging step towards utilization of conserved germplasm is the characterization of existing genetic diversity within the collection (Casa *et al.*, 2008 and Uphadyaya *et al.*, 2009). The largest diversity of the crop germplasm provides greater opportunities for improvement regarding its environmental adaptability and acquiring better agronomic traits from the crop species. Identifying and selecting the best varieties meeting specific local food and industrial requirements from this great biodiversity is of high importance for the food security assurance of any given country. Having a good knowledge of the genetic diversity of a crop often enables the plant geneticist to select the desirable family for the breeding programme and gene introgression from distantly related germplasm. The more variable genotypes / accessions can be crossed to produce better varieties that can tolerate a range of environmental changes to abiotic and biotic stresses. Therefore, a

better understanding of the genetic diversity in sorghum crop species will facilitate the further improvement of this cereal crop concerning its genetic architecture. Statistical process of categorization is generally by multivariate methods as it has wide use in summarizing and describing the innate discrepancy among the genotypes. Principal Component Analysis (PCA), Cluster analysis and discriminate analysis are the important multivariate analysis methods (Oyelola, 2004). Cluster analysis is concerned with classifying earlier unclassified materials, whereas PCA can be used to find out the resemblance between the variables and classify the genotypes (Leonard and Peter, 2009). PCA may be used to disclose the patterns and eradicate redundancy in data sets as variations regularly arise in crop species for yield and grain quality (Maji and Shaibu, 2012). Large datasets are increasingly common and are often difficult to interpret. Principal component analysis (PCA) is a technique for reducing the dimensionality of such datasets, increasing interpretability but at the same time minimizing information loss. It does so by creating new uncorrelated variables that successively maximize variance. Finding such new variables, the principal components, reduces to solving an eigenvalue/eigenvector problem, and the new variables are defined by the dataset at hand, not *a priori*, hence making PCA an adaptive data analysis technique (Jolliffe IT and Cadima J. 2016). The main objective of this research was to determine the range of variation among sorghum accessions in general and to classify them into clusters based on their similarity features regarding the quantitative characters traits under study and also to generate data on their performance for plant breeders for further evaluation of the crop in particular.

### Materials and Methods

A total of one hundred sorghum germplasm accessions received from Dr. Ramaiah gene bank located at department of Plant Genetic Resources, Centre for Plant Breeding and Genetics, Tamil Nadu Agricultural University, Coimbatore for DUS trait characterization and seed multiplication. Crop was raised during Summer 2018 at Department of Millets, TNAU, Coimbatore. The soil for the experiment is black with uniform topography and free from waterlogged conditions. Each genotype was grown in 4 rows with a spacing of 60 cm between rows and 15 cm between plants. Standard agronomic package and practices were followed to raise a healthy crop. The analysis is based on the determination of the eight quantitative traits *viz.*, time of panicle emergence (TPE), plant height upto base of flag leaf (PH. base), total plant height (PH.Totl), stem diameter (STD), length of leaf

blade (LL), width of leaf blade (LW), length of branches in panicle (PBL) and 100 grain weight (GWT). Analysis of variance using descriptive statistics such as mean, standard deviation and coefficient of variation for all the eight traits were calculated. Clustering of genotypes into similar groups was performed using Ward's hierarchical algorithm based on squared Euclidean distances. To identify the patterns of morphological variation, (PCA) was conducted. The observations recorded on eight traits were statistically analysed by using Statistical Tool for agricultural Research (STAR) package for analysis.

### Results and Discussion

The First order Statistical measures *viz.*, maximum, minimum, sum, mean, Standard Deviation (SD) and Coefficient of Variation (CV) for the measured traits are shown in Table 1. The largest variation was observed for length of branches in panicle with CV of 49.55% followed by plant height upto base of flag leaf *ie.*, 32.86%, total plant height recorded the variation of 30.55%, stem diameter recorded 23.44%, 100 grain weight 22.11% variation, width of leaf blade observed 19.61%, followed by length of leaf blade registered 13.85% variation among different morphological traits studied. The minimum level of variation was observed by the trait time of panicle emergence (CV- 10.60%).

The genotype AS 4242 has taken the minimum of 45 days for days to 50% flowering/ time of panicle emergence. The longest duration of 70 days had taken by IS 7196 for time of panicle emergence. The overall mean for time of panicle emergence in sorghum germplasm accessions was 56 days. The mean value for plant height upto base of flag leaf was 166.24 cm with the minimum and maximum value of 54.60 cm and 276.20 cm were recorded by the genotypes IS 6238 and AS 2699 respectively. For greater coefficient of variation (30.55%) was observed for total plant height *ie.*, shortest height of 89.30 cm was observed by AS 5763 and tall growth was shown by IS 8895 (322.60 cm) and average total plant height for 100 germplasm were 205.45 cm. The stem diameter ranged variation from 0.90 cm to 2.40 cm. The lowest stem thickness was measured by five genotypes *viz.*, IS 6236, AS 557/1, AS 1041, AS 2752 and highest stem diameter was recorded by AS 4104 and IS 7270. The mean value of 1.28 cm was observed for stem diameter.

The genotype MS 7885 had recorded with 95 cm leaf blade length and four accessions namely IS 6218, IS 6236, AS 4153 and AS 4242 showed shortest leaf blade length of 50 cm each. The leaf blade width ranged from narrow level of 4.2 cm (AS 557/1) to wider level of 11.10 cm (AS 4669).

The mean value for leaf blade width was 7.04 cm. Likewise length of branches in the panicle ranged from 2.10 cm to 21.50 cm. Among the traits studied length of branches in the panicle had the highest coefficient of variation. The genotype AS 2752 observed for least panicle branch length and AS 7076 showed highest length. The seed yield trait 100 grain weight recorded the overall mean of 2.22 g. The genotype IS 6312 was recorded the highest 100 grain weight of 3.8 g and genotype IS 38132 had the lowest weight of 1.3 g. Among the traits studied more than 20 per cent coefficient of variation had recorded by length of branches in panicle, plant height upto base of flag leaf, total plant height, stem diameter, length of branches in panicle and 100 grain weight.

The principal component analysis plays an extremely useful tool for studying large number of data and is desired to extract most significant data from those data points. PCA is conducted in a sequence of steps, with subjective decisions being made at many of these steps. The number of components extracted is equal to the number of variables being analysed. The first component can be expected to account for a large amount of the total variance. Each succeeding component accounts for progressively smaller amounts of variance. Data were considered in each component with Eigen values more than 1 as per the suggestions given by Brejda *et al.* (2000), which determines as a minimum 10 % of the variation. Superior Eigen values are considered as best attributes in principal components. PCA has shown the genetic diversity of the larger population or gene pool. Principal component analysis showed that out of eight components derived first three explained most of the total variations present in the gene pool. The first three principal components with Eigen value > 1 contributed about 66.35 % of the total variability among the 100 sorghum germplasm accessions evaluated for different morphological traits. The remaining five components contributed only 33.65% towards the total morphological diversity among the accessions studied. These results are presented in Table 2.

The PC 1 contributed maximum variability of 31.05% followed by PC 2 shows total phenotypic variability of 19.66% and PC 3 had contributed 15.64 % of the total variation. The important morphological traits in PC 1 were due to variations among the accessions mainly for 100 grain weight had positive factor loading value. The remaining traits *viz.*, time of panicle emergence, plant height upto base of flag leaf and total plant height had contributed negatively. PC 2 was related to diversity among sorghum genotypes due to time of panicle emergence, stem diameter, length of leaf

blade and width of leaf blade. Similarly, PC 3 expressed positive loading values for variations among genotypes resulted from plant height upto base of flag leaf, length of leaf blade and 100 grain weight and negative contributions were shown by length of branches in panicle.

Scree plot explained the percentage of variation associated with each principal component obtained by drawing a graph between Eigen values and principal component numbers. PC1 showed 31.05% variability with the Eigen value of 2.48. The Eigen values are gradually declined from PC1 to PC8. The Eigen values for remaining principal components were 1.57, 1.25, 0.86, 0.69, 0.51, 0.49 and 0.11 for PC2, PC3, PC4, PC5, PC6, PC 7 and PC 8 respectively (Fig.1). The distribution of sorghum germplasm accessions accounted by different variables from component 1, component 2 and component 3 (accessions arranged by their plot number) were distributed in different groups, which clearly showed genetic diversity among sorghum accessions (Fig. 2 to 4). Kassahun TESFAYE, 2017 utilized 117 Sorghum accessions and two standard checks for genetic diversity study. The study found that first four principal components (PCs) with eigenvalues greater than 1 explained about 71.9% of the total variation among accessions for all traits evaluated. The first principal component (PC1) obtained from the study was 26.9% and responsible loading factors were Leaf number at maturity and Plant height. Ahmed *et al.*, 2015 estimated the genetic divergence among 127 sorghum genotypes and reported 83.38 % total variation was observed for first six components and remaining two components were responsible for 46.11% variation only.

For considering a minimum threshold Eigen value of one, the first three components accounted for a cumulative of about 66.35% of the whole phenotypic diversity observed among the germplasm accessions. These findings are in accordance with the earlier findings of Nachimuthu *et al.*, 2014; Atul Kumar *et al.*, 2017 and Abraha *et al.*, 2015. Moreover, the principal component analysis also showed that the variation in the germplasm accessions cannot be explained based on few characters. Tesfamichaei *et al.*, 2015 explained the order of diminishing importance, the explanation of greater proportion of the entire phenotypic diversity involved were panicle traits, leaf traits, yield related traits and plant phenology. The same trend was observed in this investigation *viz.*, 100 grain weight, time of panicle emergence, plant height upto base of flag leaf, stem diameter, length of leaf blade and width of leaf blade. These results confirmed the previous results that also described the importance of these traits in

contributing towards the overall diversity of the sorghum germplasm land races (Ayana and Bekele, 1999). Principal component 1 explained maximum variability from total variation for morphological traits and PC 2 loaded with 19.66 % of total phenotypic variability and PC 3 had contributed 15.64 % of the total variation.

By using Agglomerative clustering method 100 germplasm accessions of sorghum grouped into eight distinct clusters (Table 3). The minimum distance between the clusters varies between 0.5 to 8.0 shown in Figure 5. Among the clusters, cluster III had the highest genotypic members (32) followed by cluster II had 23 genotypes and cluster VI held 17 genotypes. But cluster V and VI had only one genotype in each. Grouping of genotypes into different clusters confirmed the presence of variation among genotypes. Maji and Shaibu, 2012 reported that germplasm evaluation and characterization is a routine endeavour for plant breeders, and application of PCA tool, cluster and multivariate statistical analysis provide a useful means for estimating morphological diversity within and between germplasm collections. These tools are useful for the evaluation of potential breeding value and used to detect significant differences between germplasm and magnitude of deviation among crop species.

Results of the study revealed that there is a large quantity of variability in sorghum germplasm. PCA identified only few characters played important role in classifying the variation present in the germplasm. The results of the PCA revealed that the 66.35% of the total variability was explained by the first three principal components. The same trend was observed in grouping of genotypes by clustering analysis. Combination of both PCA and clustering provides the absolute results of variability existed in the gene pool (Fig.5). In this investigation also same findings were observed. Among the eight clusters Cluster number V (IS 6312) and VIII (AS 7076) had only one genotype each and both were to be superior for 100 g weight (IS 6312- 3.8g; AS 7076-3.2g) and length of panicle branches (IS 6312- 11 cm; AS 7076-21.5 cm) respectively.

The other important components time of panicle emergence, plant height upto base of flag leaf, total plant height, stem diameter, length of leaf blade and width of leaf blade were found out by PCA. In clustering method also showed highest variation for all these traits and grouped the genotypes into different clusters. Selection of parents based on these findings will produce more of genetic diversity in the crop evolution. Highest level of variability existing in the genotypes and important traits will open the scope for additional

enhancement of the cultivars in crop improvement programmes in sorghum.

Principal component analysis was utilized to examine the variation and to estimate the relative contribution of various traits for total variability. The present study identified seven best performing genotypes viz., IS 6312, IS 7830, IS 8120, AS 3479, AS 4242, AS 5476 and AS 7076 for considering the grain yield traits. These germplasm accessions may be utilized in the recombination breeding programmes to develop high yielding new sorghum varieties.

### Acknowledgement

The author would like to thank the Director CPBG, TNAU, Coimbatore, Dr. N. Nadarajan (Former Director, IIPR, Kanpur), Dr. N. Manivannan (Professor & Head, NPRC, Vamban), Dr. K. Iyanar (Assoc. Professor, Dept. of Millets) and Dr. Arunachalam, (Asst. Professor, AC & RI, Madurai) for organizing the valuable training workshop on “Systematic Plant Breeding Approaches – Transforming Field Data to Publications”.

### References

- Ahamed, K. U., Akhter, B., Islam, M.R., Alam, M.K., and Hossain, M.M. 2015. An assessment of genetic diversity in sorghum (*Sorghum bicolor* L. Moench) germplasm. *Bull. Inst. Trop. Agr., Kyushu Univ.*, **38**: 47-54.
- Ali, M.A., Jabran, K., Awan, S.I., Abbas, A., Zulkiffal, M., Acet, T., Farooq, J. and Rehman, A. 2011. Morpho-Physiological Diversity and Its Implications for Improving Drought Tolerance in Grain Sorghum at Different Growth Stages. *Australian Journal of Crop Science.*, **5**: 311-320.
- Ayana, A. and Bekele, E., 1999. Multivariate Analysis of Sorghum (*Sorghum bicolor* (L.) Moench) Germplasm from Ethiopia and Eritrea. *Genet Resource Crop Evolution.*, **46**: 273-284.
- Brejda, J.J., Moorman, T.B., Karlen, D.L., Dao, T.H. 2000. Identification of regional soil quality factors and indicators. I. Central and Southern High- Plains. *Soil Sci. Soc. Am. J.*, **64**: 2115-2124.
- Casa, A.M., Pressoir, G., Brown, P.J., Mitchell S.E., Rooney, W.L. 2008. Community resources and strategies for association mapping in Sorghum. *Crop Science.*, **48**: 30-40.
- Food and Agriculture Organization Crop Production Statistics (FAOSTAT). 2014. World Sorghum Production and Utilization. FAO, Rome.



- Hart, G., Schertz, K., Peng, Y., Syed, N. 2001. Genetic mapping of *Sorghum bicolor* (L.) Moench QTLs that control variation in tillering and other morphological character. *Theoretical and Applied Genetics.*, **103**: 1232–1242.
- Jolliffe IT and Cadima J. 2016. Principal component analysis: a review and recent developments. *Phil. Trans. R. Soc. A* **374**: 20150202.
- Kassahun TESFAYE., 2017. Genetic diversity study of sorghum (*Sorghum bicolor* (L.) Moench) genotypes, Ethiopia. *Acta Universitatis Sapientiae Agriculture and Environment.*, **9**: 44-54.
- Leonard, K., and Peter, R.J. 2009. Finding Groups in Data: An Introduction to Cluster Analysis. John Wiley & Sons. Vol. **344**.
- Maji, A.T. and Shaibu, A. A. 2012. Application of principal component analysis for rice germplasm characterization and evaluation. *J. Plant Breed. Crop Sci.*, **4**: 87-93.
- Nachimuthu, V. V., Robin, S., Sudhakar, D., Raveendran, M., Rajeswari, S., and Manonmani, S. 2014. Evaluation of rice genetic diversity and variability in a population panel by principal component analysis. *Indian J Sci Technol.*, **10**: 1555-1562.
- Oyelola BA. 2004. *The Nigerian Statistical Association preconference workshop*; 2004 Sep 20–21; University of Ibadan.
- Tesfamichael Abraha, Stephen Mwangi Githiri, Remmy Kasili, Woldeamlak Araia, Aggrey Bernard Nyende. 2015. Genetic Variation among Sorghum (*Sorghum bicolor* L. Moench) Landraces from Eritrea under Post-Flowering Drought Stress Conditions. *American Journal of Plant Sciences.*, **6**: 1410-1424.
- Upadhyaya, H.D, Reddy L.J, Dwivedi SL, Gowda CLL, Singh., S. 2009. Phenotypic diversity in cold-tolerant peanut (*Arachis hypogaea* L.) germplasm. *Euphytica.*, **165**: 279-291.



**Table 1. Characteristic means and variations for 100 Sorghum germplasm accessions**

Variable	Sum	Mean	SD	CV	Minimum Value	Maximum Value
Time of Panicle emergence (TPE)	5633	56.33	5.97	10.60	45.00	70.00
Plant height upto base of flag leaf (PH.base)	16624.3	166.24	54.62	32.86	54.60	276.20
Total Plant height (PH.Totl)	20545.2	205.45	62.76	30.55	89.30	322.60
Stem Diameter (STD)	128	1.28	0.30	23.44	0.90	2.40
Length of Leaf Blade (LL)	6664	66.64	9.23	13.85	50.00	95.00
Width of Leaf Blade (LW)	703.6	7.04	1.38	19.61	4.20	11.10
Length of Branches in Panicle (PBL)	777	7.77	3.85	49.55	2.10	21.50
100 grain weight (GWT)	221.6	2.22	0.49	22.11	1.30	3.80

**Table 2. Principal component analysis for different morphological traits recorded in *Sorghum bicolor* L. Moench**

Principal Components	PC 1	PC 2	PC 3	PC 4	PC 5	PC 6	PC 7	PC 8
Eigen values	2.4843	1.5727	1.2510	0.8676	0.6981	0.5169	0.4956	0.1155
% of total Variance	0.3105	0.1966	0.1564	0.1084	0.0873	0.0646	0.0618	0.0144
Cumulative Variance %	0.3105	0.5071	0.6635	0.7720	0.8590	0.9236	0.9856	1.0000
<b>Factor loading by different morphological traits</b>								
Time of Panicle emergence (TPE)	<b>-0.3250</b>	<b>0.4255</b>	-0.0838	0.0109	<b>0.7031</b>	-0.1036	<b>-0.4315</b>	-0.1237
Plant height upto base of flag leaf (PH.base)	<b>-0.5103</b>	-0.1974	<b>0.3118</b>	0.0498	-0.2866	0.2823	-0.2239	<b>-0.6242</b>
Total Plant height (PH.Totl)	<b>-0.5883</b>	-0.0184	0.0892	0.1741	-0.1568	0.1882	-0.0995	<b>0.7380</b>
Stem Diameter (STD)	-0.0285	<b>0.6043</b>	-0.2557	<b>0.4103</b>	-0.1192	<b>0.4396</b>	<b>0.4173</b>	-0.1564
Length of Leaf Blade (LL)	-0.2435	<b>0.3219</b>	<b>0.5784</b>	-0.2244	0.0264	<b>-0.4124</b>	<b>0.5299</b>	-0.0205
Width of Leaf Blade (LW)	0.2319	<b>0.4505</b>	0.1937	<b>-0.6124</b>	-0.2639	0.3671	-0.3350	0.1231
Length of Branches in Panicle (PBL)	-0.2848	0.2087	<b>-0.5550</b>	-0.1869	<b>-0.4865</b>	<b>-0.5176</b>	-0.1123	-0.0967
100 grain weight (GWT)	<b>0.3047</b>	0.2537	<b>0.3773</b>	<b>0.5817</b>	-0.2788	-0.3288	<b>-0.4172</b>	0.0317

**Table 3. Distribution of 100 germplasm accessions into eight clusters**

Clusters	Cluster size	Name of the genotypes in cluster
I	3	IS 6218 (1), IS 6238 (4), IS 9783 (78)
II	23	IS 6232 (2), IS 6304 (8), IS 9985 (18), IS 18473 (26), IS 18475 (27), IS 38132 (28), AS 270/1 (32), AS 321 (33), AS 557/1 (34), AS 1041 (35), AS 2699 (36), AS 2752 (37), AS 2902 (39), AS 4104 (45), AS 4599 (52), AS 4647 (53), AS 6327 (65), AS 6616 (67), AS 6620 (68), IS 9645 (76), IS 9912 (79), AS 7749 (85), MS 7885 (100)
III	32	IS 6236 (3), IS 6297 (6), IS 6303 (7), IS 8120 (13), IS 15130 (23), AS 3289 (40), AS 3641 (42), AS 3998 (43), AS 4003 (44), AS 4153 (46), AS 4226 (47), AS 4242 (48), AS 4243 (50), AS 5546 (57), AS 5763 (58), AS 5787 (59), AS 5794 (60), AS 5812 (61), AS 6002 (62), AS 6184 (63), AS 6342 (66), AS 6642 (69), AS 6653 (70), IS 1126 (72), IS 9780 (77), IS 7695 (83), MS 7818 (88), MS 7806 (89), MS 7819 (90), MS 7862 (96), MS 7863 (97), MS 7881 (99)
IV	10	IS 6241 (5), IS 7270 (12), AS 4249 (49), AS 4294 (51), AS 4669 (54), MS 7904 (55), AS 5476 (56), IS 7830 (74), MS 7894 (86), MS 7816 (87)
V	1	IS 6312 (9)
VI	17	IS 6334 (10), IS 9143 (16), IS 10027 (19), IS 15221 (25), AS 2797 (38), AS 3479 (41), AS 6191 (64), IS 7196 (73), IS 9098 (75), MS 178 (80), MS 342 (81), MS 5826 (82), AS 7738 (84), MS 7820 (91), MS 7837 (92), MS 7851(95), MS 7868 (98)
VII	13	IS 6549(11), IS 8296(14), IS 8845(15), IS 9407(17), IS 11024(20), IS 11056(21), IS 14385(22), IS 38060(24), IS 68621(29), IS 70034 (30), AS 1 (31), AS 9841 (93), MS 7843 (94)
VIII	1	AS 7076 (71)

\* Parenthesis denotes serial number for drawing dendrogram.

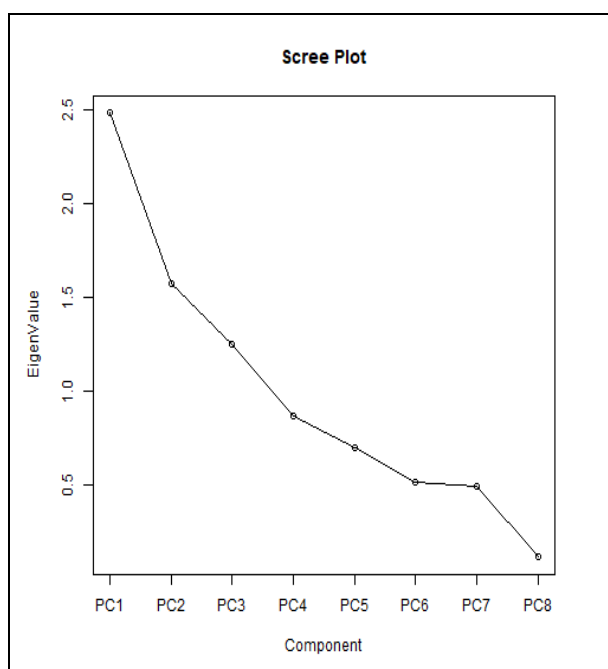


Fig. 1. Scree plot for Eigen values and principal components

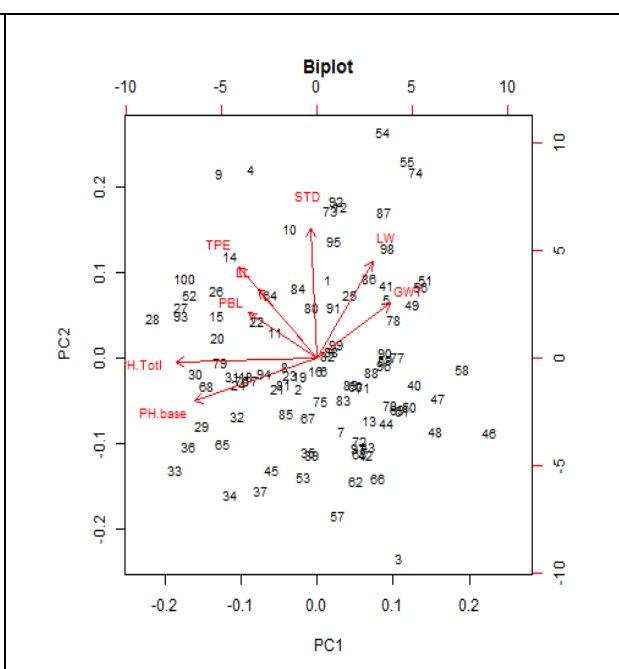


Fig. 2. Distribution of genotypes across first two components based on PCA

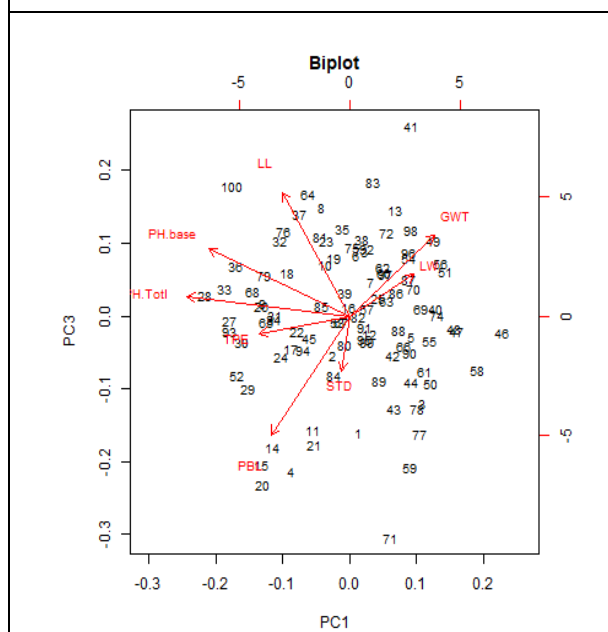


Fig. 3. Distribution of genotypes across first and third components based on PCA

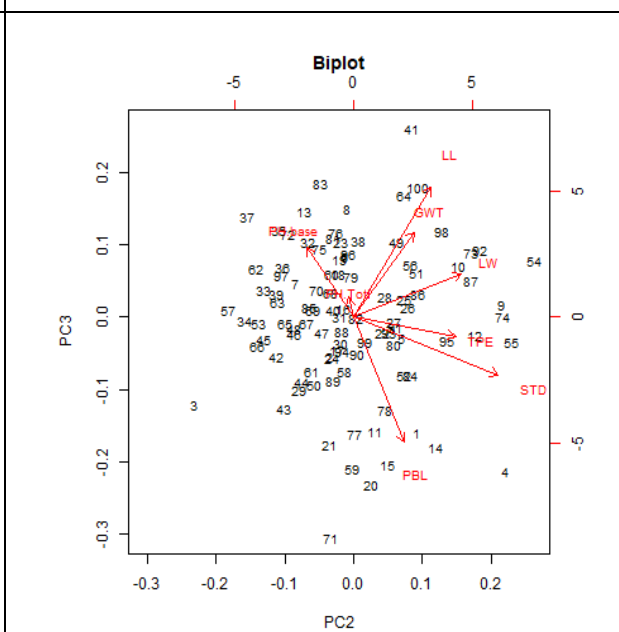
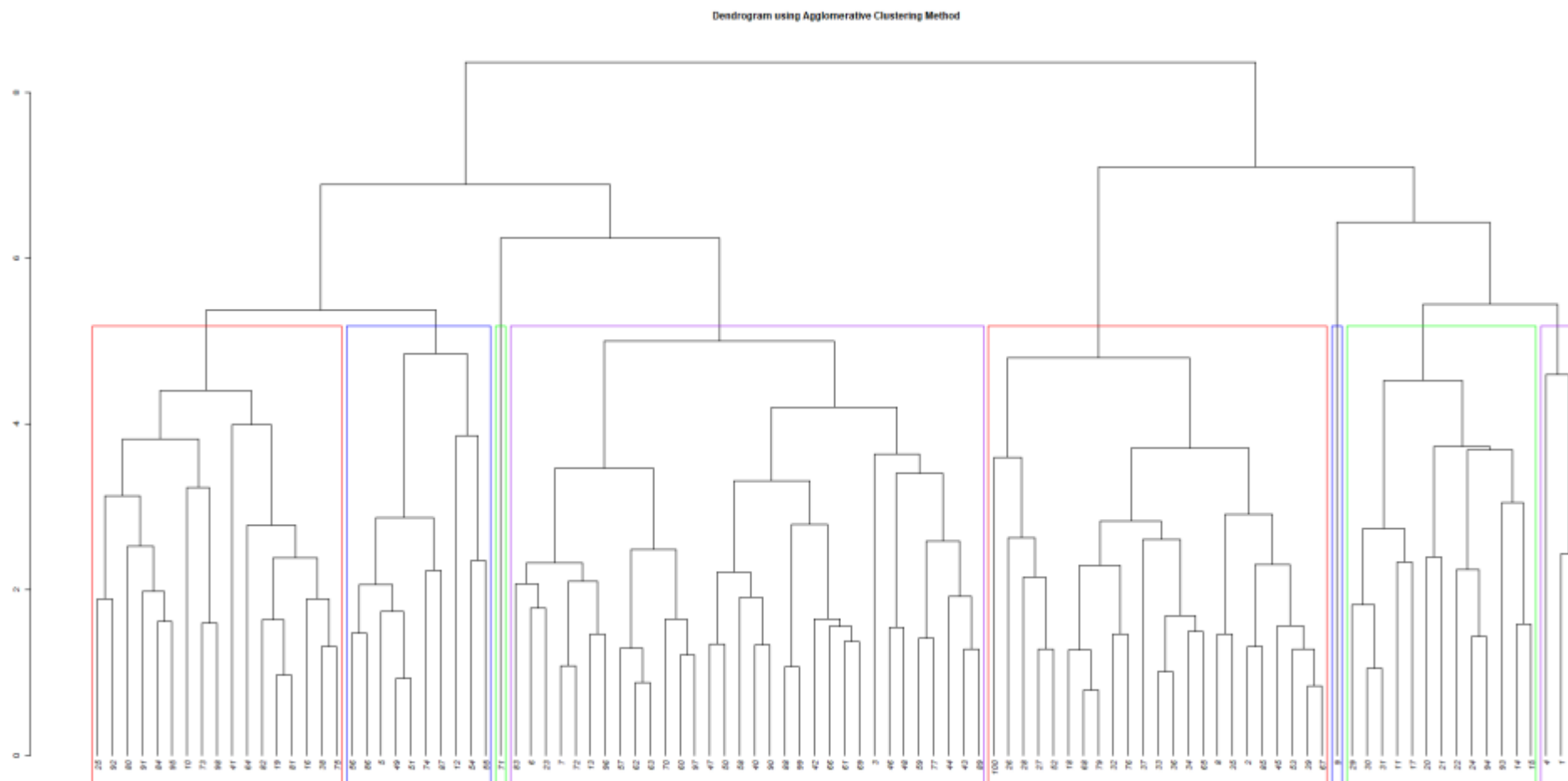


Fig. 4. Distribution of genotypes across second and third components based on PCA





**Fig. 5. Phenotypic dendrogram generated by Agglomerative Clustering method**