

Semi-unsupervised quantization of strawberry shape diversity in modern germplasm

Mitchell J. Feldmann & Steven J. Knapp

Department of Plant Sciences, University of California – Davis, Davis CA, 95616

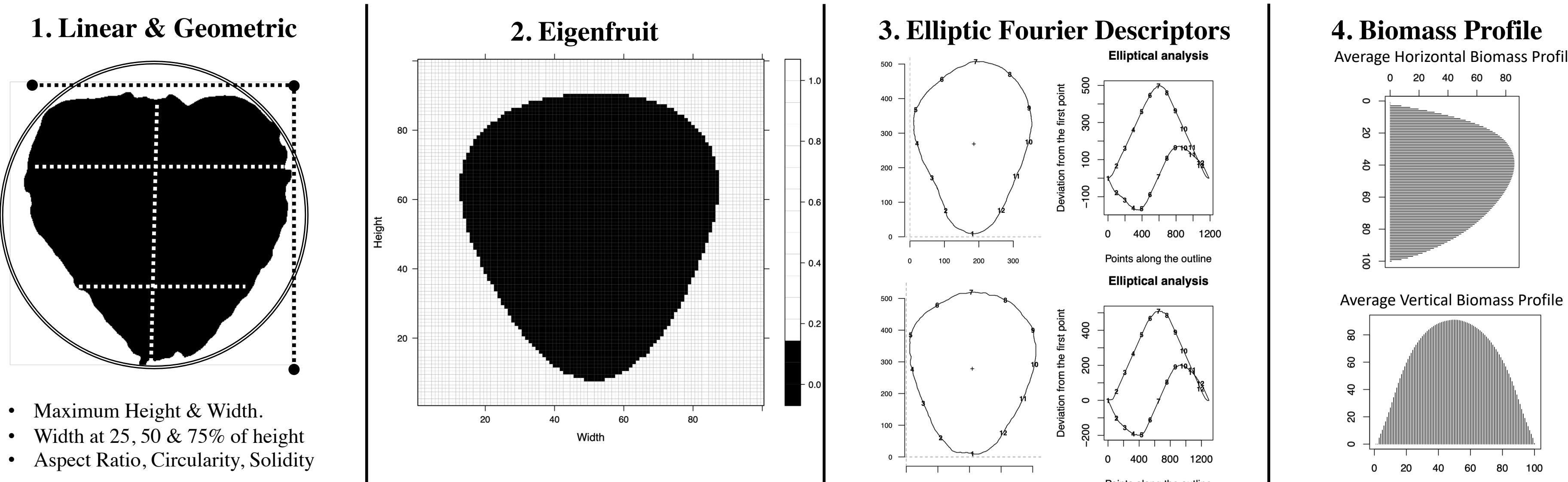


@MitchFeldmann
mjfeldmann@ucdavis.edu
mjfeldmann

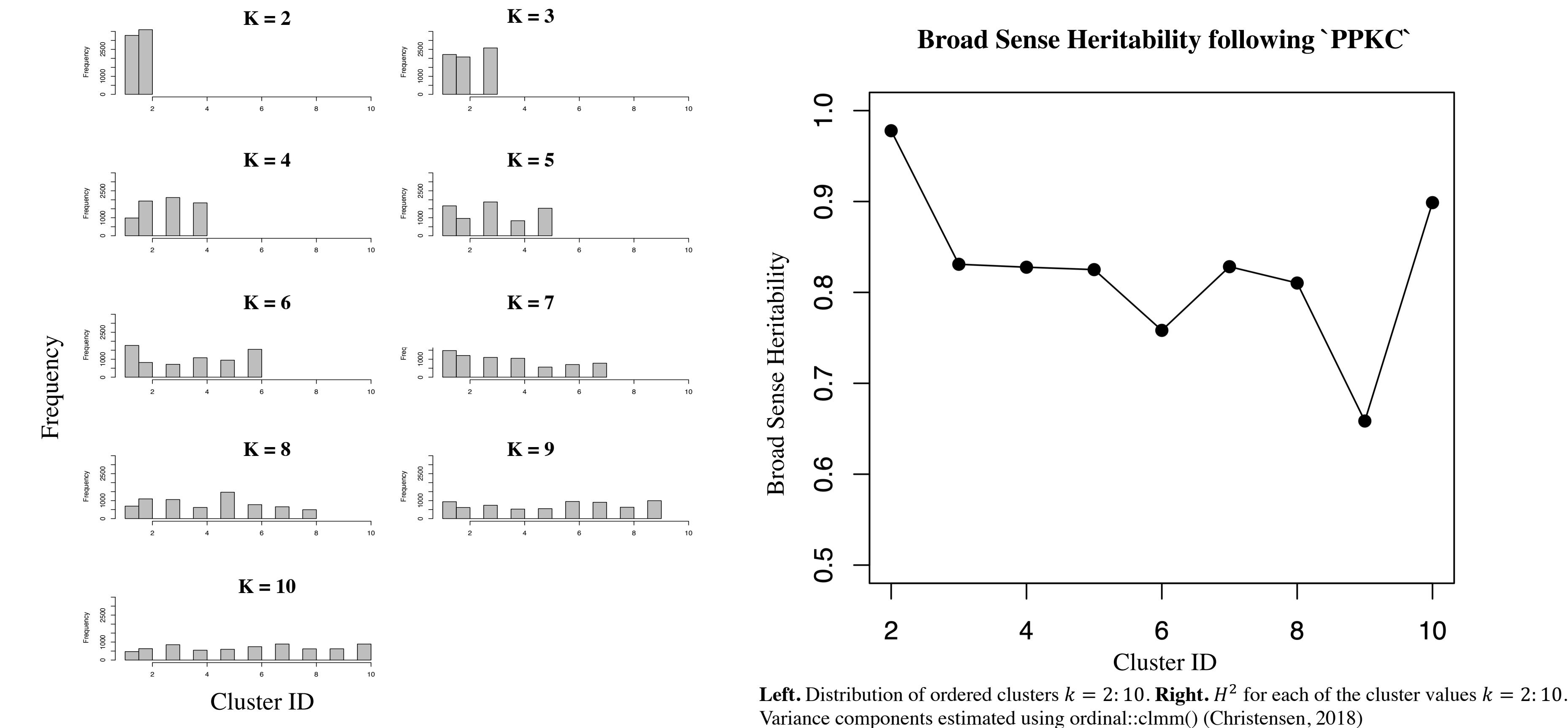
1. Summary

The shape of strawberry fruit is intimately linked, along with other visual attributes, to their marketability. To explore the existing morphospace in an unbiased method, K-means clustering was used to determine appropriate number of fruit categories at varying levels. A new method for transforming nominal descriptions onto an ordinal scale was applied to afford the use of robust statistical analyses common in quantitative genetics. Variance components estimated from multinomial mixed models indicated broad sense heritability (H^2) ranging from 0.65 for $k = 9$ and 0.98 for $k = 2$. Using a complex set of comprehensive and linear phenotypes, random forest regression was used to determine the most important quantitative phenotypes for classifying shapes. Of the 97 original quantitative descriptors only 15 were kept for classification. Mean accuracies from 10-fold cross validation ranged from 88-97% for discriminant models and 96-99% for support vector machines. These analyses suggest that there is a significant amount of genetic variation attributable to shape classifications and that few variables are needed for classification.

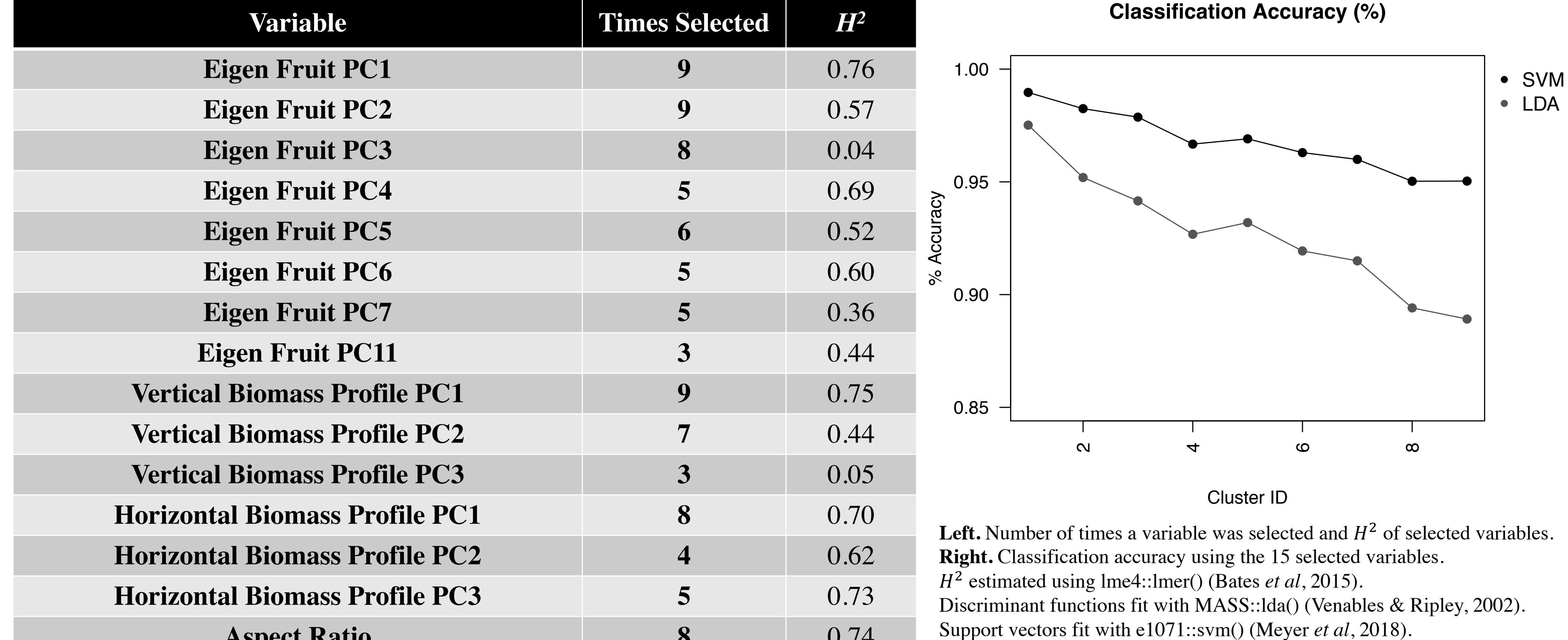
3. Quantitative Phenotypes



4. Distribution and H^2 of Ordered Clusters



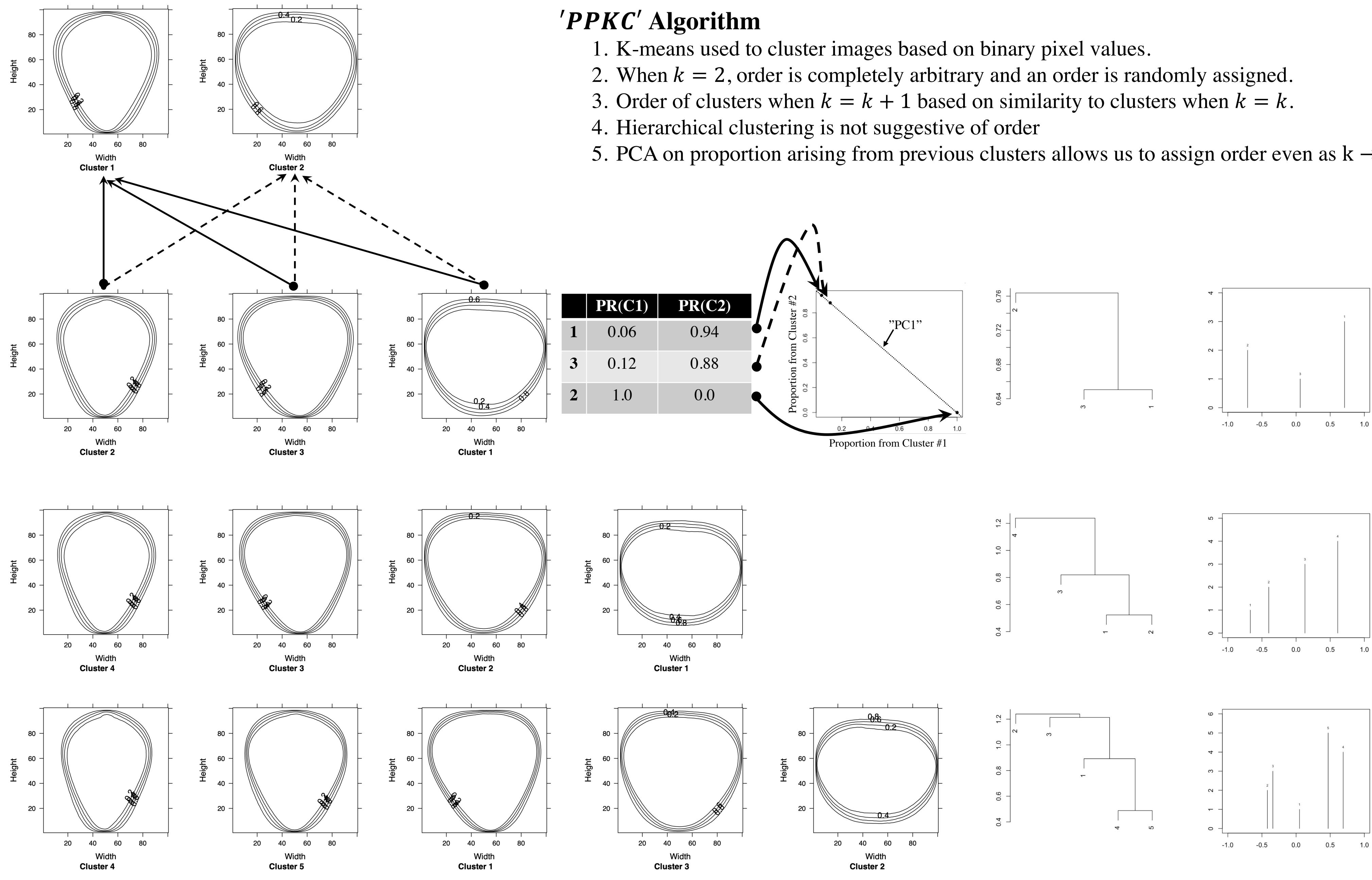
6. Heritability and Classification Accuracy



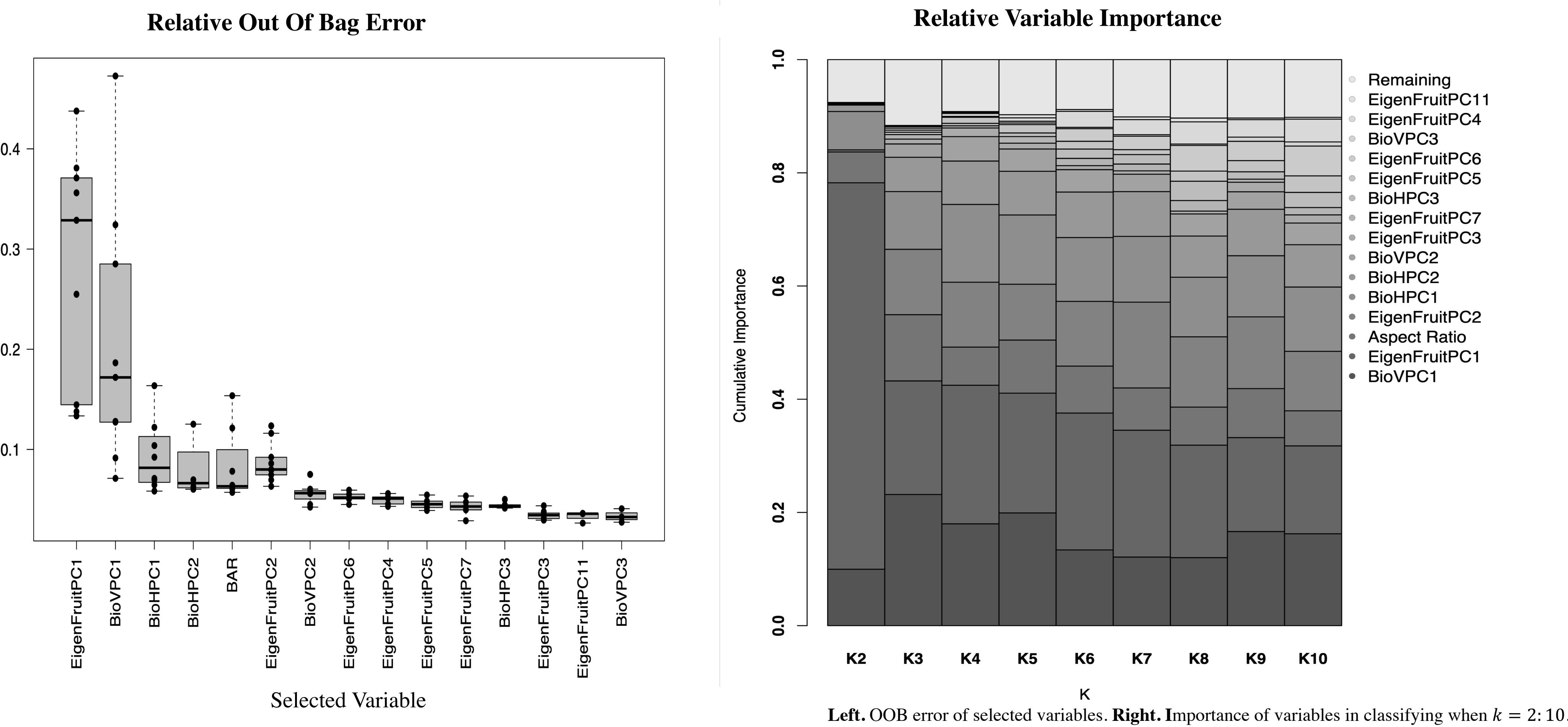
2. K-means Clustering & Principal Progression of K Clusters

'PPKC' Algorithm

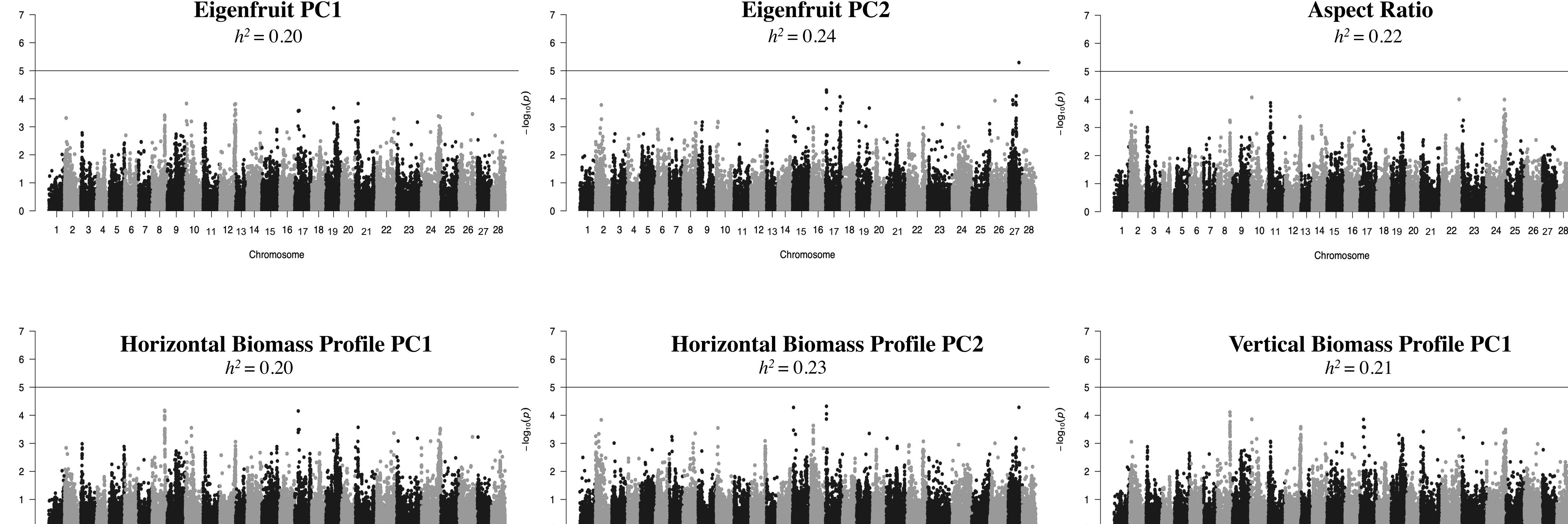
1. K-means used to cluster images based on binary pixel values.
2. When $k = 2$, order is completely arbitrary and an order is randomly assigned.
3. Order of clusters when $k = k + 1$ based on similarity to clusters when $k = k$.
4. Hierarchical clustering is not suggestive of order
5. PCA on proportion arising from previous clusters allows us to assign order even as $k \rightarrow \infty$



5. Variable Selection using Random Forests



7. Genome Wide Association Scan



Top left to bottom right: Manhattan plots for the 6 variables with the largest OOB error. PC1 and PC2 of EigenFruit, Aspect ratio, PC1 and PC2 of horizontal biomass profile, and PC1 of vertical biomass profile. GWAS performed using rrblup::gwas() (Endelman, 2011).