Arthur Bernardeli
March 24, 2025
Exam 1 (Midterm) - AGRO 932

1    **Genomic Patterns in Distinct Soybean Breeding Pools: Selection and FST**

2    **Dynamics**

3    Soybean (*Glycine max* (L.) Merr.) is a globally important crop, cultivated for its

4    high-value seeds rich in both oil and protein. It serves as a critical component of food,

5    feed, and industrial markets, with significant production hubs across the Americas and

6    Asia (US Department of Agriculture Foreign Agricultural Service, 2017). In the United

7    States, soybean is a cornerstone of agricultural output and economic stability, occupying

8    millions of acres annually. Its dual utility, approximately 40% protein and 20% oil in seed

9    composition makes it especially valuable for both human consumption and livestock feed.

10   Given its essential role in global nutrition and supply chains, improving soybean traits

11   remains a primary objective for plant breeding programs worldwide.

12   Among the most economically and nutritionally significant traits in soybeans are

13   seed yield and seed protein content. However, these traits are negatively correlated,

14   presenting a substantial challenge for breeders. Genetic gains in seed yield have often

15   been accompanied by declines in protein concentration, and efforts to increase protein

16   content frequently come at the expense of yield potential. This antagonistic relationship

17   is rooted in complex genetic and physiological trade-offs, which are not yet fully

18   understood. As a result, breeding programs have often pursued these traits in parallel but

19   separate pipelines, selecting elite lines for either high yield or high protein. While this

20   strategy enables targeted improvements, it limits opportunities for simultaneous

21   enhancement of both traits and underscores the need to better understand the genetic

22   mechanisms underlying this trade-off.

23     Despite significant advances in genomic scan studies in soybean seed composition

24     traits, such as in Zhang et al. (2016, 2018), the genetic architecture and selection dynamics

25     that maintain or potentially resolve the yield-protein trade-off remain poorly defined.

26     There is limited insight into how sustained directional selection shapes allele frequencies

27     at trait-associated loci, and how such divergence may manifest across the genome in

28     breeding populations that share a common elite ancestry. Population genomics

29     approaches, such as the fixation index (FST), offer a powerful means of quantifying

30     genetic differentiation and identifying regions of the genome that have responded to trait-

31     specific selection pressures. Moreover, the integration of recombination and inter-group

32     crossing offers the potential to break unfavorable linkage blocks and generate superior

33     recombinants that transcend historical trait limitations.

34     This study aims to investigate these genomic and breeding dynamics using a

35     simulation-based approach. We simulated elite soybean populations subjected to

36     divergent selection for either high yield or high protein content. These simulations are

37     designed to mimic sequential selection and mating cycles commonly performed by

38     soybean breeding programs, enabling a controlled and realistic examination of selection

39     outcomes over time. The simulated populations initially share a common elite

40     background, from which two selection paths emerge. Each path undergoes repeated intra-

41     population selection, resulting in phenotypic advancement and genomic divergence.

42     Genotypic and phenotypic data are collected over multiple selection cycles, enabling the

43     analysis of genome-wide FST, chromosome-specific patterns of divergence, and changes

44     in within-group diversity and trait distributions.

45     The first aim of this study is to evaluate the degree of genetic differentiation

46     between the high-yield and high-protein populations following divergent selection. Using

47     genome-wide and chromosome-specific FST estimates, we assess how selection shapes

48    allele frequency distributions in each breeding pool. We hypothesize that loci associated

49    with yield or protein content will exhibit moderate to high differentiation values as

50    selection cycles advance.

51          The second aim is to determine how repeated cycles of intra-population selection

52    impact both genetic structure and phenotypic means within each group. We hypothesize

53    that selection will reduce within-group diversity by favoring specific haplotypes, thereby

54    increasing genetic homogeneity and driving further divergence at key loci across the

55    genome. This will be reflected in rising phenotypic means and increasing FST values

56    relative to the founder generation.

57          The third aim explores the consequences of crossing between the divergent

58    breeding pools, followed by selection in the recombinant progeny. This scenario is

59    designed to assess whether recombination can effectively break the negative correlation

60    between yield and protein, allowing for the development of superior lines that combine

61    favorable alleles from both parental groups.

62 **Material and Methods**

63 *Population Simulation and Structure*

64     To investigate genomic divergence and selection responses in soybean breeding

65 pools, we simulated one diverse founder population, each consisting of 20 high-yield lines

66 and 20 high-protein lines. These population will then represent breeding pools with

67 distinct selection histories, one selected for high seed yield and the other for high seed

68 protein content. Genotypic data were simulated using 6,000 evenly spaced single-

69 nucleotide polymorphism (SNP) markers, distributed across the 20 soybean

70 chromosomes (~300 markers per chromosome), to mimic the marker density commonly

71 used in breeding programs and genomic studies. Genotype data were imputed leveraging

72 a next-generation sequencing dataset based on the *Glycine max* reference genome

73 Williams 82 (*www.soybase.org*).

74     The high yield lines (group Y) was designed to reflect an elite germplasm

75 background with a historical focus on improving seed yield. This group had a high base

76 population mean for yield and a moderate mean for protein content, consistent with the

77 documented negative correlation between the two traits. The high protein lines (group P)

78 was simulated from a similar elite background but selected for increased seed protein

79 concentration. Due to the antagonistic relationship between the traits, this group exhibited

80 a reduced yield mean. Both populations shared common genetic ancestry, emulating real-

81 world breeding scenarios where selection diverges from a shared elite base.

*Simulation Platform and Trait Architecture*

83        Simulations were conducted using the AlphaSimR (Gaynor et al., 2021) package

84    in R (R Core Team, 2021), a widely used platform for modeling complex trait inheritance,

85    recombination, and selection in plant breeding, according to specific simulation

86    parameters. The simulated traits included seed yield and seed protein content, with a

87    genetic correlation of -0.45 imposed between them to reflect the biological constraint

88    observed in empirical breeding data. Heritability values were set at 0.45 for yield and 0.70

89    for protein.

90        Trait values were assigned at the founder level, generating variation within the

91    founder population for groups Y and P. This variation allowed for the initial assessment

92    of divergence and provided a foundation for within-population selection in subsequent

93    cycles.

94    *Experimental Procedure*

95        Step 1: Founder Population Simulation

96        The initial step involved simulating one founder population of 40 individuals,

97    each with assigned genotypes and trait values according to the parameters described

98    earlier. These populations served as the baseline for evaluating genomic differentiation

99    and trait performance.

100       Step 2: Intra-Group Recurrent Selection Cycles

101       Each simulated population underwent five generations of trait-specific selection.

102    In each cycle, the top 10% of individuals were selected based on phenotypic values for

103    their respective target trait, yield for the Y group and protein content for the P group.

104 Selected individuals were intermated within their group to produce progeny for the next

105 generation. Recombination and segregation were simulated during mating, and new

106 phenotypic values were assigned based on inherited genotypes and environmental effects.

107 This process was repeated for five cycles, allowing the accumulation of selection

108 responses and genetic changes across generations.

109     Step 3: FST, PCA, and phenotypic analysis

110     Following the founder simulation, genome-wide and chromosome-specific FST

111 values were calculated between the Y and P groups using SNP data. This allowed for the

112 quantification of genetic differentiation attributable to divergent selection histories. In

113 parallel, principal component analysis (PCA) was performed to visualize the overall

114 genetic structure and clustering of individuals by selection group. FST and PCA analyses

115 were repeated in each cycle to evaluate changes in population structure and genetic

116 clustering. Comparisons were made between each generation and its respective founder

117 population, as well as between the high-yield and high-protein populations in cycles 0

118 and 5. For soybeans, a similar procedure was performed by Yang et al. (2022), Silva et

119 al. (2025), and Andrijanić et al. (2023) to compute the fixation index in distinct soybean

120 pools. Phenotypic means of yield and protein content were recorded for each group at

121 every selection cycle to assess phenotypic gains and trade-offs over time.
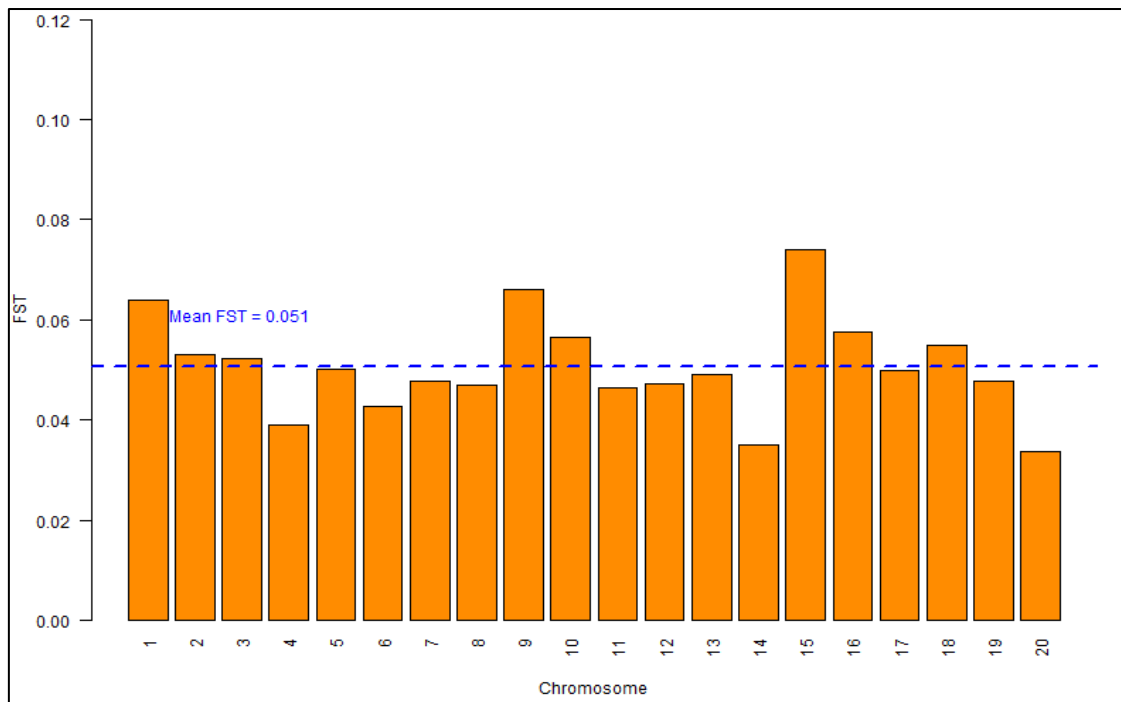
122 *Codes*

123     Codes are available in the Supplementary Material section.

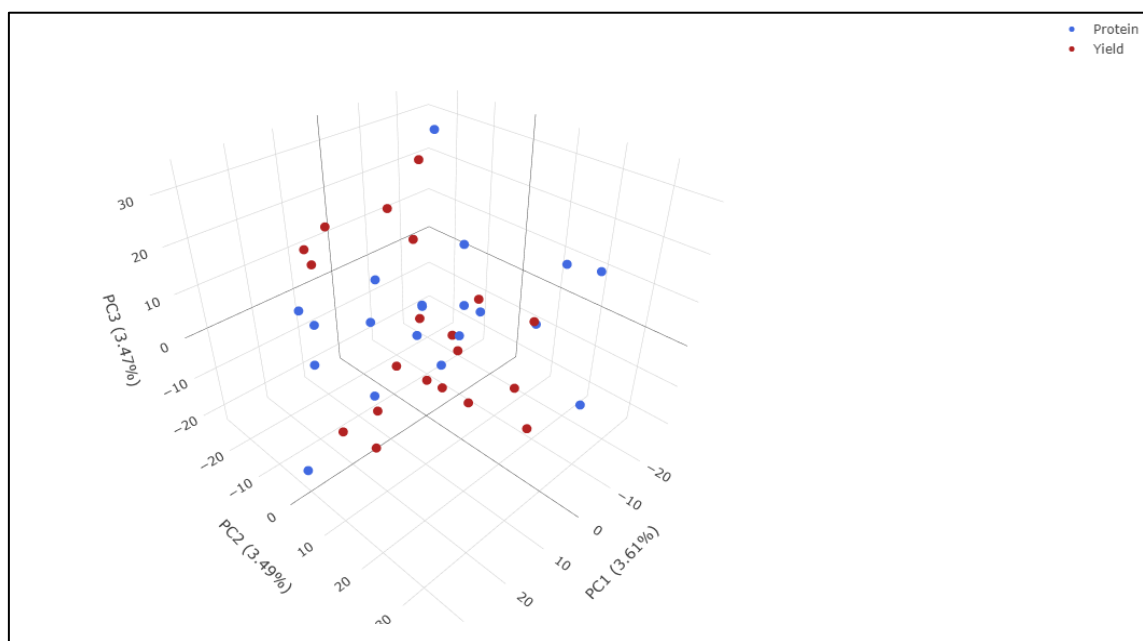124 **Results and Discussion**

125       This study investigated the genomic consequences and phenotypic benefits of

126 selection in two distinct soybean groups over five generations. Each group was initially

127 designed with identical genetic backgrounds but selected for contrasting traits, one for

128 high protein content (group P) and the other for high yield (group Y). Using fixation index

129 (FST) and principal component analysis (PCA), the genetic structure and differentiation

130 of the populations were tracked over selection cycles. Mean phenotypic responses were

131 monitored to evaluate the efficiency of selection. This approach enabled us to understand

132 the magnitude and direction of genomic change resulting from selection and how it

133 correlates with phenotypic performance.

134 *Initial Genetic Structure and Differentiation (Cycle 0)*

135       Before selection was applied, the two founder groups displayed minimal genetic

136 differentiation. Chromosome-wise FST values (Figure 1) were low, indicating a shared

137 genetic base between the two groups. Although these groups were designed for distinct

138 selection objectives (protein versus yield), the genetic structure had not yet diverged. PCA

139 results (Figure 2) confirmed this, with individuals from both groups overlapping in

140 multivariate space, showing no discernible clusters. This genetic similarity served as a

141 baseline, allowing the effects of subsequent selection to be clearly attributed to breeding

142 pressure rather than founder variation.

**Figure 1.** Chromosome-wise FST between individuals of the founder population (cycle 0).
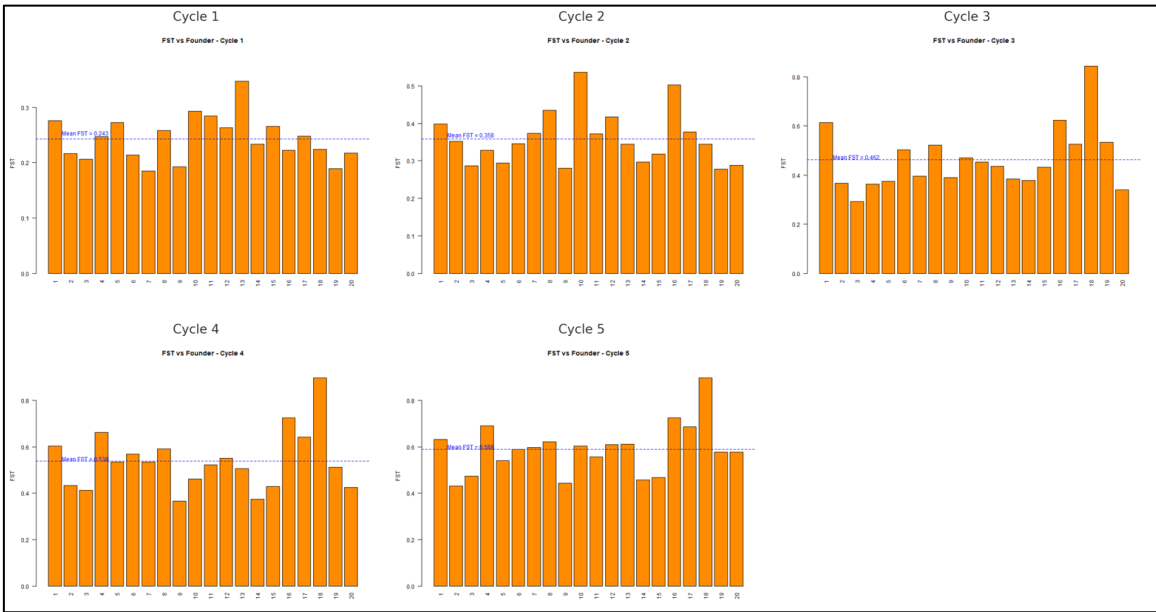


**Figure 2.** Principal component analyses plot (PCA) between individuals of the founder population (cycle 0).
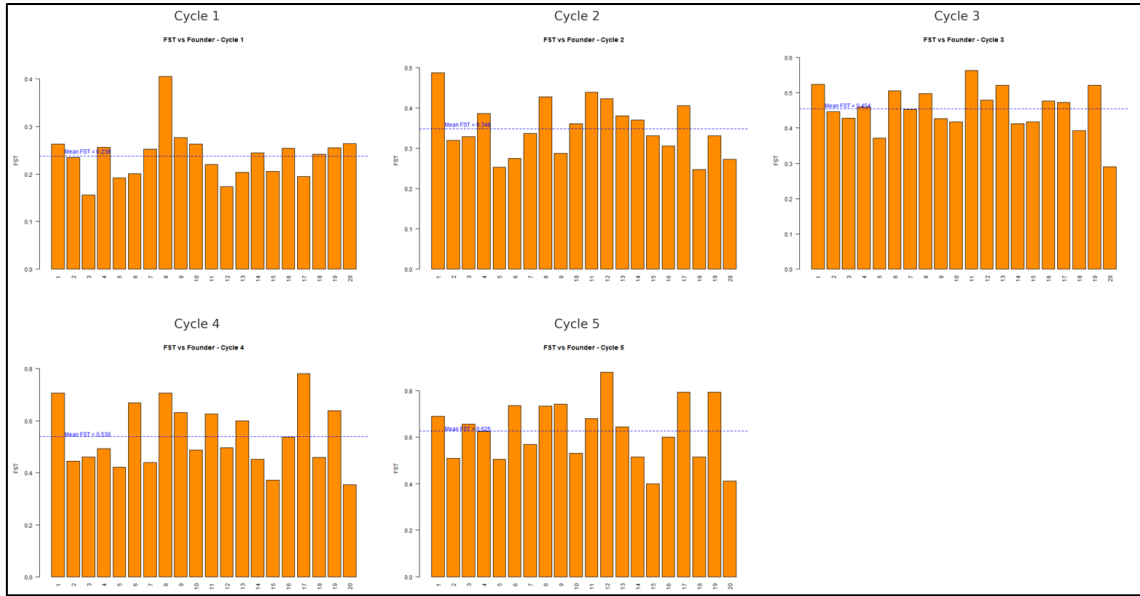
*FST Across Selection Cycles*

144       Following selection, an increase in FST values was observed in both populations

145   relative to their founder state. In group P, FST rose from 0.243 in Cycle 1 to 0.589 in

146   Cycle 5 (Figure 3), while in group Y, it increased from 0.238 to 0.626 over the same

147   period (Figure 4). These trends reflect selection pressure, favoring alleles associated with

148   the respective target traits. The steady rise in FST over cycles is indicative of both reduced

149   within-population genetic diversity and increased divergence from the ancestral gene

150   pool. This suggests that only a subset of alleles, likely those conferring favorable

151   phenotypic effects, were retained through recombination and selection.



**Figure 3.** FST between protein population and founder population for each selection cycle.

**Figure 4.** FST between yield population and founder population for each selection cycle.
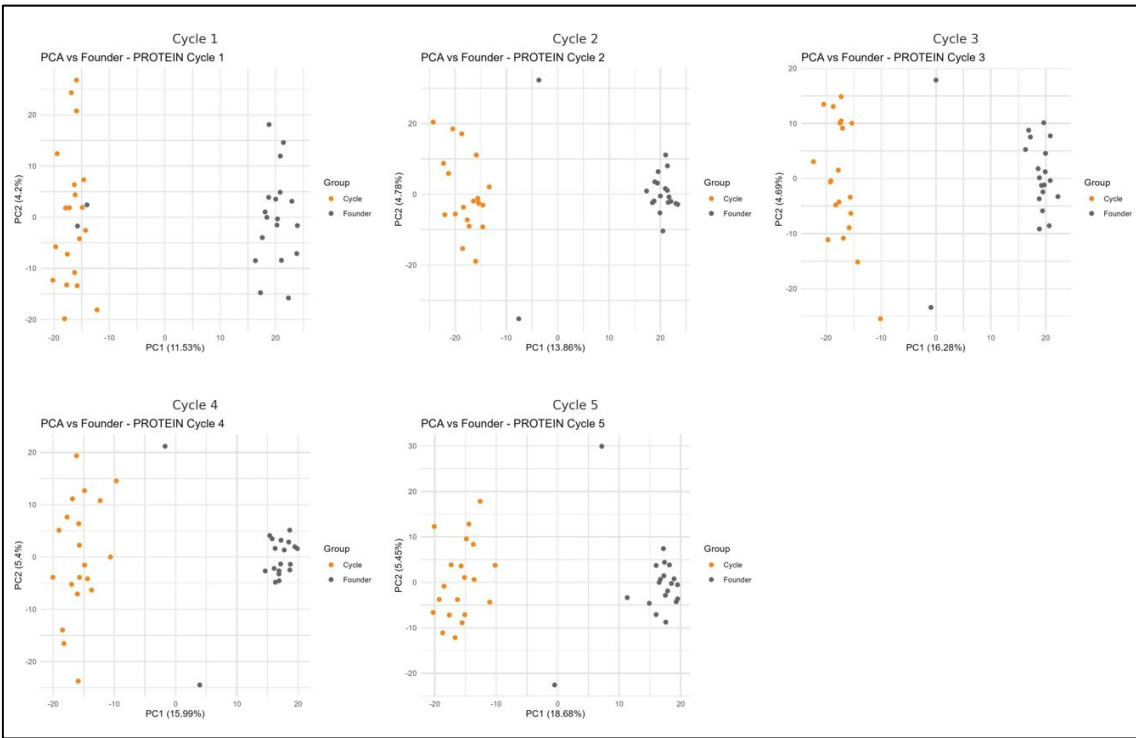
152        To further investigate the impact of trait-based selection, we directly compared

153    groups P and Y. Table 1 shows the interpopulation FST for each cycle, which increased

154    from 0.277 in Cycle 1 to 0.622 by Cycle 5. This growing divergence emphasizes how

155    selection for contrasting phenotypes reshapes the genome in divergent directions. As

156    selection advanced, each population fixed or retained different allelic combinations that

157    improved their respective trait values, which cumulatively enhanced genetic separation.

158    This is particularly relevant in applied breeding programs where the trade-off between

159    protein and yield must be managed strategically.

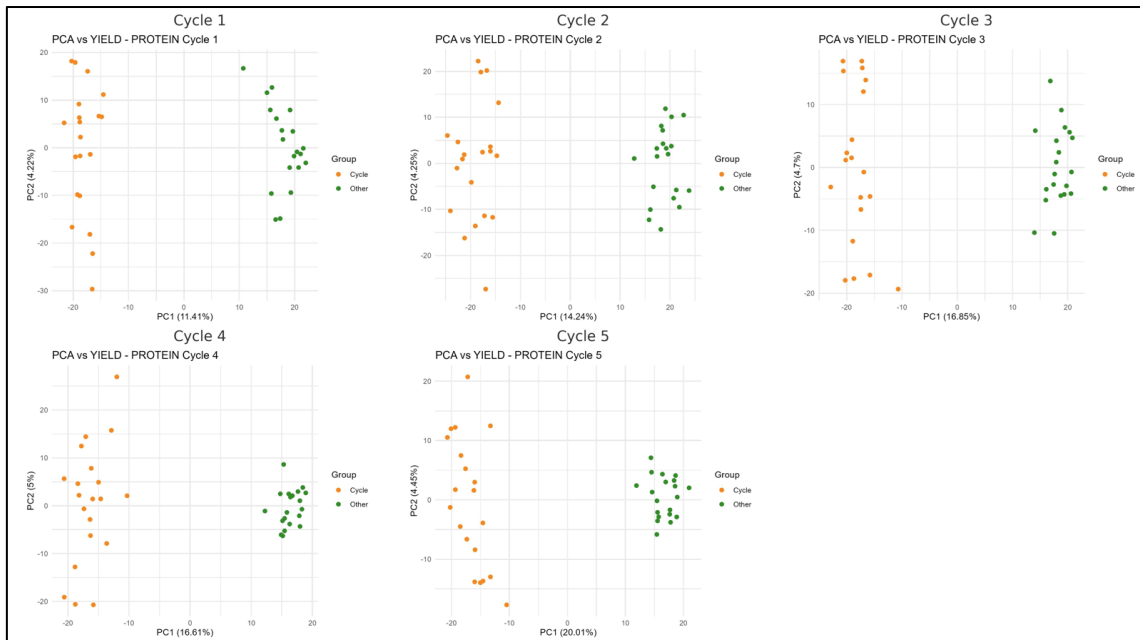**Table 1.** FST between groups P and Y for each selection cycle.

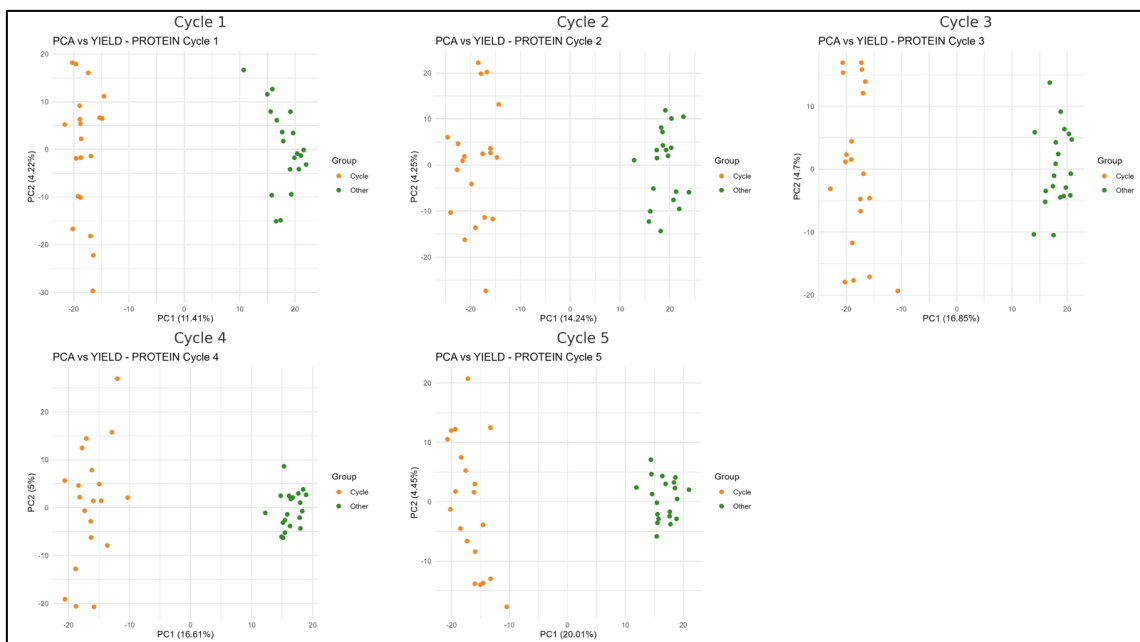| Cycle | FST |
|---------|-------|
| Cycle 1 | 0.277 |
| Cycle 2 | 0.389 |
| Cycle 3 | 0.494 |
| Cycle 4 | 0.566 |
| Cycle 5 | 0.622 |

*PCA Across Selection Cycles*

161        The genomic divergence described by FST was also clearly visualized using PCA.

162    In the early cycles, the overlap between groups was still visible, but from Cycle 1 onward,

163    distinct clusters emerged for both protein- and yield-selected lines (Figures 5 and 6). In

164    Cycle 5, PCA plots demonstrated complete separation, confirming that directional

165    selection had generated distinct genetic trajectories. This pattern was reinforced when

166    comparing the two groups at each cycle (Figure 7). These plots are consistent with our

167    hypothesis of loss of shared alleles and accumulation of beneficial mutations or

168    combinations specific to the selection regime.



**Figure 5.** PCA between group P and founder population for each selection cycle.

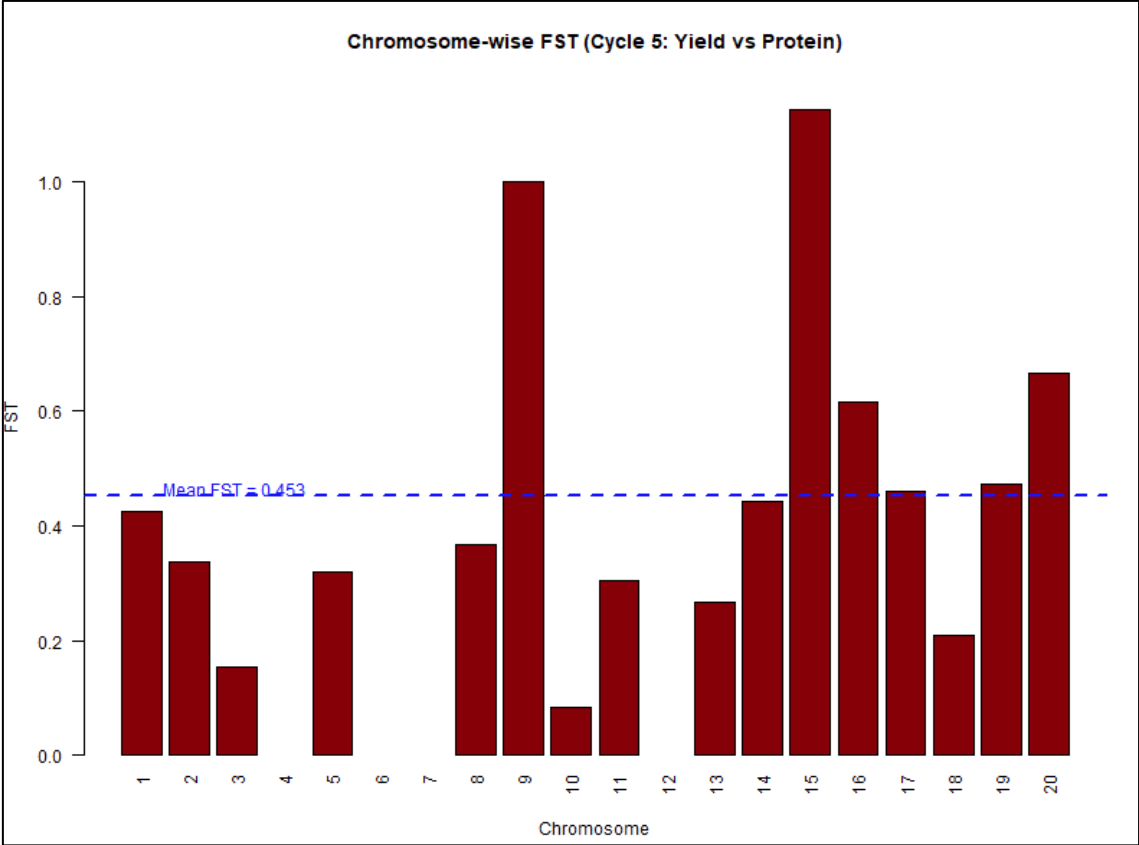**Figure 6.** PCA between Y group and founder population for each selection cycle.



**Figure 7.** PCA between P and Y groups for each selection cycle.
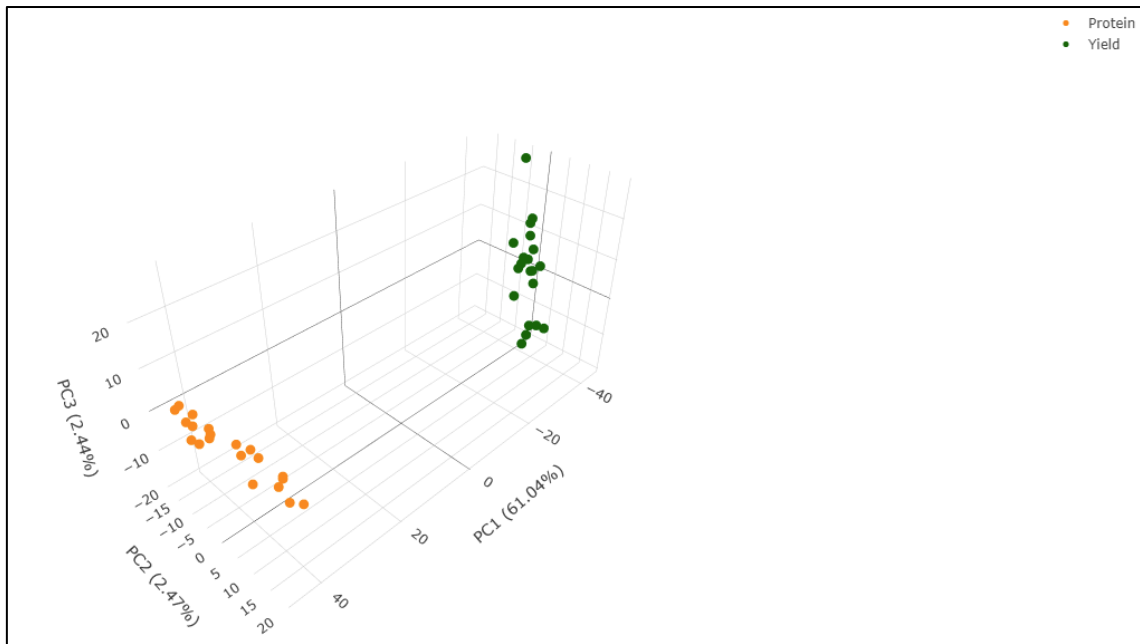
169    *FST and PCA in Cycle 5*

170         Chromosome-specific FST analysis between the P and Y groups in Cycle 5

171    (Figure 8) revealed variable patterns of divergence. Some chromosomes showed stronger

172    differentiation, potentially harboring loci with major effects for protein or yield traits.

173    This heterogeneity is expected given the polygenic nature of both traits. It also suggests

174   that despite whole-genome selection, specific genomic regions were disproportionately

175   impacted. PCA of Cycle 5 (Figure 9) supports this observation, with tight clustering

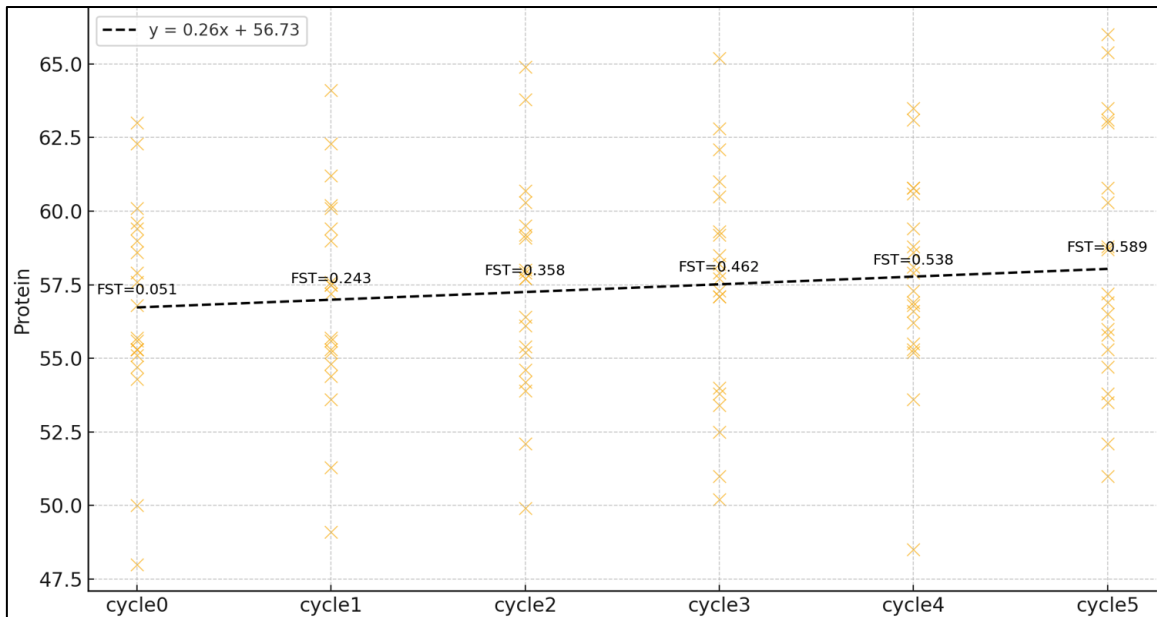176   within populations and broad separation between them.



**Figure 8.** Chromosome-wise FST between yield and protein population in cycle.
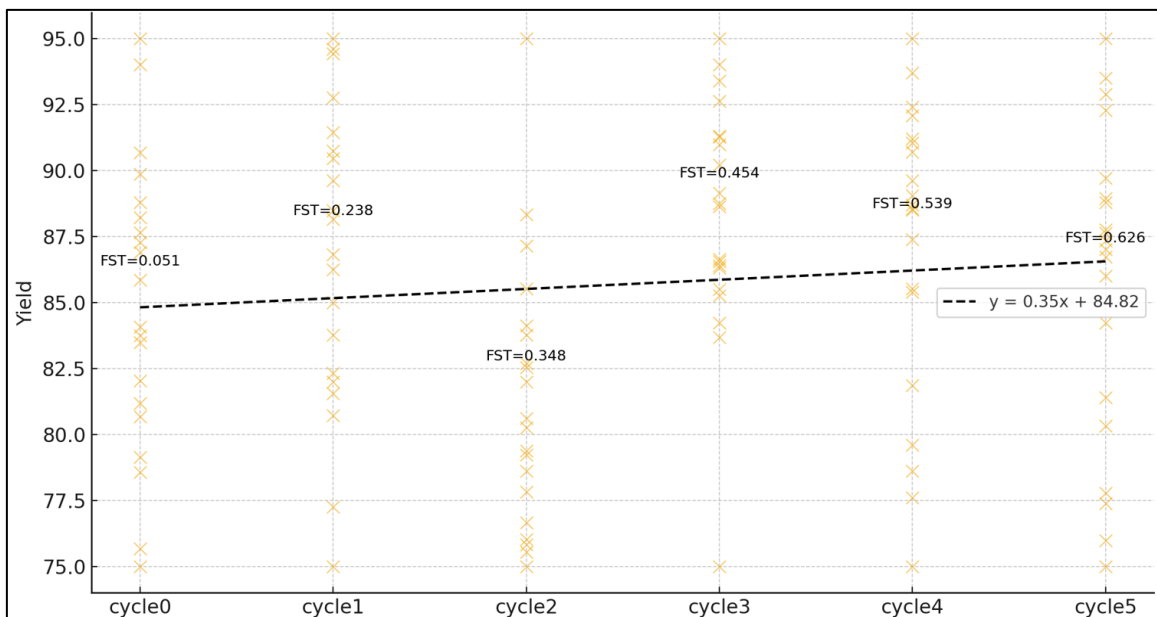
**Figure 9.** PCA between yield and protein population in cycle.

*FST Values and Changes in Phenotypes Across Cycles*

177

178    The consistent increase in FST coincided with favorable shifts in mean phenotypic

179    values. For group P, average protein content rose across selection cycles (Figure 10).

180    Likewise, group P exhibited rising yield values (Figure 11), highlighting selection

181    success. It also reflects the narrowing of genetic diversity due to cycles of selection.

**Figure 10.** Changes in FST and protein content mean for each selection cycle for group P. Protein content (y-axis) is represented on a percentage basis.



**Figure 11**. Changes in FST and yield content mean for each selection cycle for group Y. Yield (y-axis) is represented on a bushels per acre basis.

182    This study demonstrates the power and consequences of selection in breeding

183    populations by combining population genomics with selection scenarios. It contributes to

184    our understanding of how elite breeding pools diverge and interact, ultimately informing

185    strategies to improve both yield and protein content in soybeans. Even starting from

186    genetically similar founders, consistent selection for distinct traits produced rapid

187 genomic divergence and substantial phenotypic gains. These outcomes reinforce the

188 utility of genomic tools such as FST and PCA in monitoring genetic change. For breeders,

189 the results underscore the trade-offs inherent in selection: while targeted gains are

190 possible, they often come at the cost of genetic diversity. In practice, this necessitates

191 striking a balance between progress and the maintenance of long-term variability.

192 Moreover, the divergence observed here provides a rationale for developing pools based

193 on protein and yield traits.

194 **Future directions**

195 The next phase of this study will assess the impact of FST in a new population

196 generated from progenies resulting from crosses between the high-yield and high-protein

197 lines in cycle 5. We expect that the mean FST will decrease relative to the founder

198 population, as the cycle 6 population is expected to redistribute genetic variability from

199 both parental pools across the genome. For the same reason and given the negative

200 correlation between yield and protein content, it is also expected that the mean phenotypic

201 values for both traits will decline compared to those observed in Cycle 5.

202 **References**

203 Andrijanić, Z., Nazzicari, N., Šarčević, H., Sudarić, A., Annicchiarico, P., & Pejić, I.

204 (2023). Genetic diversity and population structure of European soybean

205 germplasm revealed by single nucleotide polymorphism. *Plants*, *12*(9), 1837.

206 Gaynor RC, Gorjanc G, Hickey JM (2021). "AlphaSimR: an R package for breeding

207 program simulations." G3 Gene|Genomes|Genetics, 11(jkaa07).

208 https://doi.org/10.1093/g3journal/jkaa017.

209     R Core Team (2021). R: A language and environment for statistical computing. R
210         Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-
211         project.org/.

212     Silva, A. C. D., Gregorio da Silva, D. C., Ferreira, E. G. C., Abdelnoor, R. V., Borém, A.,
213         Arias, C. A., ... & Marcelino-Guimarães, F. C. (2025). Genetic diversity,
214         population structure in a historical panel of Brazilian soybean cultivars. *PloS*
215         *one*, *20*(1), e0313151.

216     US Department of Agriculture Foreign Agricultural Service. (2017). Oilseeds: World
217         Production Markets and Trade Reports (Washington, DC: FAS).

218     Yang, C., Yan, J., Jiang, S., Li, X., Min, H., Wang, X., & Hao, D. (2022). Resequencing
219         250 soybean accessions: new insights into genes associated with agronomic traits
220         and genetic networks. *Genomics, Proteomics & Bioinformatics*, *20*(1), 29-41.

221     Zhang, J., Song, Q., Cregan, P. B., & Jiang, G. L. (2016). Genome-wide association study,
222         genomic prediction and marker-assisted selection for seed weight in soybean
223         (Glycine max). Theoretical and Applied Genetics, 129, 117-130.

224     Zhang, J., Wang, X., Lu, Y., Bhusal, S. J., Song, Q., Cregan, P. B., ... & Jiang, G. L.
225         (2018). Genome-wide scan for seed composition provides insights into soybean
226         quality improvement and the impacts of domestication and breeding. Molecular
227         Plant, 11(3), 460-472.