

WGCNA – Data Cleaning

```
1 rm(list=ls())
2
3 library(jpeg)
4 library(dplyr)
5 library(tidyr)
6 library(tibble)
7 library(stringr)
8 library(ggplot2)
9
10 library(foreach)
11 library(iterators)
12 library(parallel)
13 library(doParallel)
14
15 library(WGCNA)
16 # library(KEGGREST)
17 # library(biomaRt)
18
19 set.seed(1)
20
21 # Enable WGCNA threads to speed up calculations
22 enableWGCNAThreads()
23
24
```

```
35 #####
36 # Output folder
37 #####
38 output_path <- file.path("/home/ycth8/data/projects/05_30_2021_summer_WGCNA/Maize_proteomics_output/2021_06_10_data_cleaning")
39
40 if(!dir.exists(output_path)){
41   dir.create(output_path, showWarnings=FALSE, recursive=TRUE)
42   if(!dir.exists(output_path)){
43     quit(status=1)
44   }
45 }
46
47
```

```
48 #####
49 # Read in input file
50 #####
51
52 folder_path = file.path("/home/ycth8/data/projects/05_30_2021_summer_WGCNA/Maize_proteomics_data")
53
54 dat = read.csv(
55   file = file.path(folder_path, "3plus_2fold_genelist.csv"),
56   header = TRUE,
57   check.names = FALSE,
58   stringsAsFactors = FALSE
59 )
60
61 datTraits = read.csv(
62   file = file.path(folder_path, "PBAA_B73o2_Abs_Raw.csv"),
63   header = TRUE,
64   check.names = FALSE,
65   stringsAsFactors = FALSE
66 )
67
68
```

Original Genotype Input File

Accession	Description	O2_10_1	O2_10_2	O2_10_3	O2_14_1	O2_14_2	O2_14_3	O2_18_1	O2_18_2	O2_18_3	O2_22_1	O2_22_2
Zm00001d004855	Zm00001d004855	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d004960	Zm00001d004960	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d005090	Clathrin heavy chain	210	207	234	194	209	242	201	218	217	222	198
Zm00001d006433	Ras-related protein Rab-6A	23	26	23	18	20	23	15	19	20	19	18
Zm00001d006596	Elongation factor 1-alpha	120	129	127	123	141	147	136	138	123	153	145
Zm00001d006651	TUBA5, Tubulin alpha chain	165	134	153	122	114	111	106	107	114	109	97
Zm00001d007048	Zm00001d007048	36	31	28	32	22	34	33	31	31	29	27
Zm00001d007758	Zm00001d007758	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d007900	Eukaryotic translation initiation factor 3 subunit I	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d007936	Pyridoxin biosynthesis protein ER1	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d008651	Pyruvate decarboxylase	22	22	25	30	31	36	38	33	38	41	39
Zm00001d008739	Zm00001d008739	7	9	11	10	4	6	5	0.5	0.5	8	4
Zm00001d008743	Ubiquitinyl hydrolase 1	20	16	22	24	33	22	22	32	33	33	36
Zm00001d009311	Zm00001d009311	8	9	10	10	6	9	7	7	8	9	9
Zm00001d009652	Zm00001d009652	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d009653	Zm00001d009653	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d009708	Putative calcium-dependent protein kise family protein isoform ...	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d011183	THI1-1, Thiamine thiazole synthase chloroplastic [S...	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d014073	60S ribosomal protein L23	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
Zm00001d014844	0.5	19	21	17	40	53	54	59	63	61	54	60

Original Phenotype Input File

Genotype	DAP	Replicaion	Ala	Arg	Asx	Glx	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser
B73	10	1	168.34	26.01	147.24	118.77	185.03	25.52	52.02	65.28	49.08	26.01	25.52	93.25	82.45
B73	10	2	182.93	25.85	144.39	120.00	180.98	24.88	53.17	64.39	48.78	24.39	26.34	92.20	80.98
B73	10	3	190.85	26.14	131.24	120.78	191.90	27.71	49.15	61.18	50.72	24.05	26.67	88.37	84.18
B73	14	1	120.34	20.06	82.96	89.34	155.44	25.53	30.54	41.48	40.57	14.13	20.97	59.26	53.33
B73	14	2	140.35	23.86	94.04	89.82	172.63	24.80	48.19	66.90	44.91	18.71	25.26	56.61	65.96
B73	14	3	130.87	22.70	87.00	74.14	137.30	19.29	49.93	73.38	38.20	17.78	21.18	48.42	65.44
B73	18	1	126.78	23.86	84.21	85.15	161.87	24.33	58.95	94.97	43.51	18.71	26.20	60.82	66.43
B73	18	2	116.69	24.97	81.53	93.25	171.21	27.01	61.66	101.91	44.84	20.89	28.03	66.24	68.28
B73	18	3	114.58	23.74	74.32	88.26	160.52	27.87	54.19	90.32	43.35	17.03	27.87	59.87	60.90
B73	22	1	153.83	24.68	89.04	95.21	163.53	21.60	71.85	132.67	40.99	21.60	25.56	78.02	71.85
B73	22	2	135.80	20.23	80.90	82.35	132.91	16.98	69.71	123.88	33.95	22.39	21.31	68.98	65.73
B73	22	3	128.52	21.99	78.92	87.55	150.51	21.99	62.10	116.87	39.25	18.98	24.58	69.43	62.96
B73	26	1	135.53	23.53	85.18	93.18	162.82	24.47	60.71	106.82	38.59	17.88	26.35	66.35	62.59
B73	26	2	144.78	26.28	94.06	102.36	165.07	22.59	75.16	136.48	41.50	24.44	26.74	82.54	79.77
B73	26	3	199.47	29.47	107.89	117.37	193.68	26.32	95.79	168.42	45.79	25.26	31.05	96.32	87.37
B73	30	1	167.27	28.05	89.35	107.01	183.38	26.49	78.96	148.57	42.60	20.26	29.61	84.16	75.84
B73	30	2	143.66	25.13	80.42	88.38	148.27	20.52	71.62	121.88	38.53	20.94	23.87	68.27	69.53

```
69 #####
70 # Process the input file
71 #####
72
73 datExpr0 = dat %>%
74   pivot_longer(!c(Accession, Description), names_to = "Genotype", values_to = "Measurement") %>%
75   separate(Genotype, c("Genotype", "Repetition"), sep = "\\_(?=[^\\_]+$)", extra = "drop", fill = "right") %>%
76   group_by(Accession, Description, Genotype) %>%
77   summarise_at(vars(Measurement), list(Measurement = mean)) %>%
78   ungroup() %>%
79   select(-Description) %>%
80   pivot_wider(names_from = Accession, values_from = Measurement, values_fill = NA) %>%
81   as.data.frame(stringsAsFactors = FALSE)
82
83 # Copy values in column 1 into row names
84 rownames(datExpr0) <- datExpr0[,1]
85
86 # Remove column 1
87 datExpr0 <- datExpr0[,-1]
88
89 print(head(datExpr0))
90 print(tail(datExpr0))
91 print(dim(datExpr0))
92
```

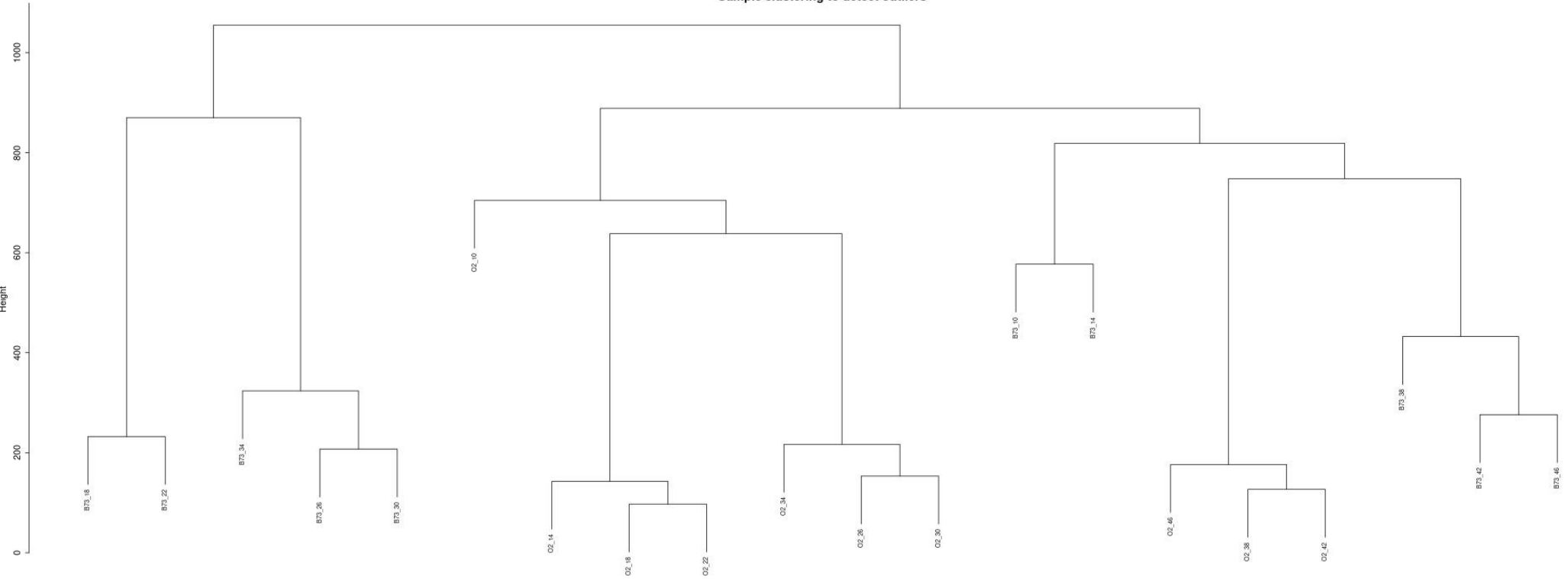
Processed Input File

	Zm00001d001895	Zm00001d002065	Zm00001d004837	Zm00001d004855	Zm00001d004960	Zm00001d005062	Zm00001d005090
B73_10	3.5	18	5.66666666666667	216.333333333333	70.3333333333333	0.5	0.5
B73_14	5.16666666666667	21.6666666666667	10.3333333333333	229	97.6666666666667	0.5	0.5
B73_18	9	22	14.3333333333333	206	137.666666666667	0.5	0.5
B73_22	11.6666666666667	15	16	209.666666666667	140.666666666667	0.5	0.5
B73_26	11.6666666666667	17.6666666666667	11	0.5	140.666666666667	0.5	216.333333333333
B73_30	9.33333333333333	11.6666666666667	10.6666666666667	0.5	149.666666666667	0.5	213.333333333333
B73_34	8.33333333333333	13.6666666666667	8.33333333333333	0.5	133.666666666667	0.5	179.333333333333
B73_38	8.66666666666667	10	10.3333333333333	0.5	136.333333333333	0.5	0.5
B73_42	6	11	11.3333333333333	0.5	111.333333333333	0.5	0.5
B73_46	6.33333333333333	13	0.5	0.5	114.666666666667	0.5	125
O2_10	7.33333333333333	8.66666666666667	0.5	0.5	0.5	0.5	217
O2_14	11.3333333333333	6.33333333333333	0.5	0.5	0.5	0.5	215
O2_18	15.3333333333333	5.66666666666667	0.5	0.5	0.5	0.5	212
O2_22	14.6666666666667	9	0.5	0.5	0.5	0.5	201.666666666667
O2_26	17	7	0.5	0.5	118.666666666667	16.3333333333333	205.666666666667
O2_30	21.3333333333333	3.83333333333333	0.5	0.5	107	18.3333333333333	213.666666666667
O2_34	18	4.66666666666667	0.5	0.5	101	19.6666666666667	181
O2_38	17.6666666666667	6	0.5	0.5	87.6666666666667	26.3333333333333	166.333333333333
O2_42	8.33333333333333	5.33333333333333	0.5	0.5	81.6666666666667	20.3333333333333	129
O2_46	6	6.66666666666667	0.5	0.5	79.3333333333333	15.6666666666667	89.3333333333333

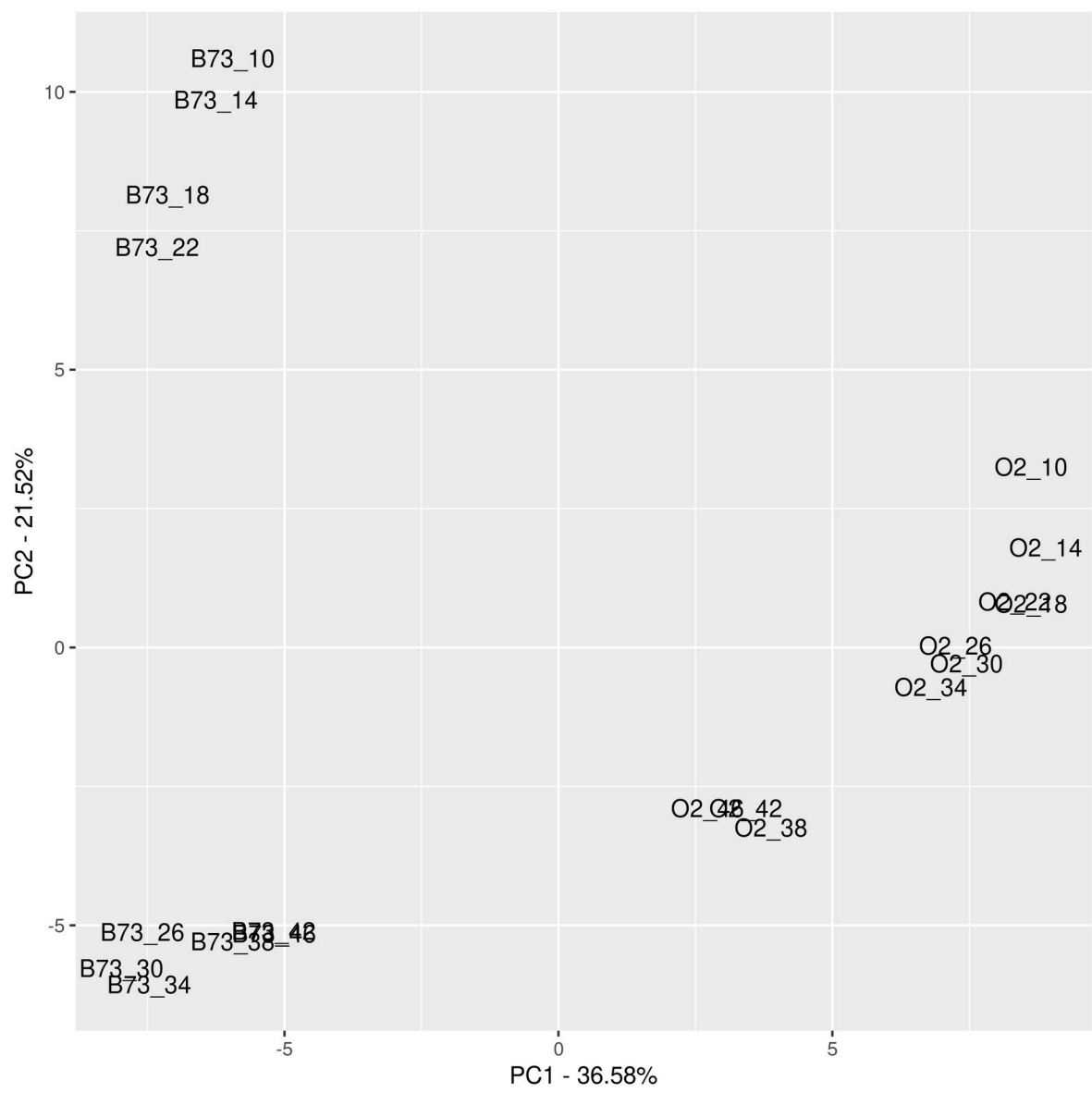

```
94 #####
95 # Quality check
96 #####
97
98 gsg = goodSamplesGenes(datExpr0, verbose = 3)
99
100 cat(rep("\n", 2))
101 print(gsg$allOK)
102
103 # If gsg$allOK is FALSE, remove genes and samples
104 if (!gsg$allOK)
105 {
106   # Optionally, print the gene and sample names that were removed:
107   if (sum(!gsg$goodGenes)>0)
108     printFlush(paste("Removing genes:", paste(names(datExpr0)[!gsg$goodGenes], collapse = ", ")))
109   if (sum(!gsg$goodSamples)>0)
110     printFlush(paste("Removing samples:", paste(rownames(datExpr0)[!gsg$goodSamples], collapse = ", ")))
111
112   # Remove the offending genes and samples from the data:
113   datExpr0 = datExpr0[gsg$goodSamples, gsg$goodGenes]
114 }
115
```

```
117 #####
118 # Cluster the data
119 #####
120
121 sampleTree = hclust(dist(datExpr0), method = "average")
122
123 # Plot tree
124 cat(rep("\n", 2))
125 jpeg(filename = file.path(output_path, "sampleTree.jpeg"), width = 1920, height = 720)
126 par(cex = 0.6)
127 par(mar = c(0,4,2,0))
128 plot(sampleTree, main = "Sample clustering to detect outliers", sub="", xlab="", cex.lab = 1.5, cex.axis = 1.5, cex.main = 2)
129 # abline(h = 15, col = "red")
130 dev.off()
131
132 # # Determine cluster under the line
133 # clust = cutreeStatic(sampleTree, cutHeight = 15.2, minSize = 10)
134 #
135 # print(table(clust))
136 #
137 # # Remove outlier sample
138 # keepSamples = (clust==1)
139 # datExpr = datExpr0[keepSamples, ]
140
141 datExpr = datExpr0
142 nGenes = ncol(datExpr)
143 nSamples = nrow(datExpr)
144
145 # Collect garbage
146 collectGarbage()
147
```

Sample clustering to detect outliers



```
149 #####
150 ## Performs principal components analysis (PCA)
151 #####
152 pca_df <- datExpr
153 pca <- prcomp(pca_df, center = TRUE, scale. = TRUE)
154
155 pca_variance <- pca$sdev^2
156 pca_variance_percentage <- round(pca_variance/sum(pca_variance, na.rm = TRUE)*100, digits=2)
157
158 x <- pca$x %>% as.data.frame(stringsAsFactors=TRUE)
159
160 ## Plot PC1 vs PC2 with GGplot2
161 p <- ggplot(x, aes(x=PC1, y=PC2)) +
162   geom_text(label=rownames(x)) +
163   labs(
164     x=paste0("PC1 - ", pca_variance_percentage[1], "%"),
165     y=paste0("PC2 - ", pca_variance_percentage[2], "%")
166   )
167
168 ggsave(
169   filename="PCA_Plot.jpg",
170   plot=p,
171   path=output_path
172 )
173
```



```
175 #####
176 # Save processed data
177 #####
178 write.csv(
179   x = datExpr,
180   file = file.path(output_path, "..", "datExpr.csv"),
181   na = "",
182   quote = FALSE
183 )
184
```

Original Phenotype Input File

Genotype	DAP	Replicaion	Ala	Arg	Asx	Glx	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser
B73	10	1	168.34	26.01	147.24	118.77	185.03	25.52	52.02	65.28	49.08	26.01	25.52	93.25	82.45
B73	10	2	182.93	25.85	144.39	120.00	180.98	24.88	53.17	64.39	48.78	24.39	26.34	92.20	80.98
B73	10	3	190.85	26.14	131.24	120.78	191.90	27.71	49.15	61.18	50.72	24.05	26.67	88.37	84.18
B73	14	1	120.34	20.06	82.96	89.34	155.44	25.53	30.54	41.48	40.57	14.13	20.97	59.26	53.33
B73	14	2	140.35	23.86	94.04	89.82	172.63	24.80	48.19	66.90	44.91	18.71	25.26	56.61	65.96
B73	14	3	130.87	22.70	87.00	74.14	137.30	19.29	49.93	73.38	38.20	17.78	21.18	48.42	65.44
B73	18	1	126.78	23.86	84.21	85.15	161.87	24.33	58.95	94.97	43.51	18.71	26.20	60.82	66.43
B73	18	2	116.69	24.97	81.53	93.25	171.21	27.01	61.66	101.91	44.84	20.89	28.03	66.24	68.28
B73	18	3	114.58	23.74	74.32	88.26	160.52	27.87	54.19	90.32	43.35	17.03	27.87	59.87	60.90
B73	22	1	153.83	24.68	89.04	95.21	163.53	21.60	71.85	132.67	40.99	21.60	25.56	78.02	71.85
B73	22	2	135.80	20.23	80.90	82.35	132.91	16.98	69.71	123.88	33.95	22.39	21.31	68.98	65.73
B73	22	3	128.52	21.99	78.92	87.55	150.51	21.99	62.10	116.87	39.25	18.98	24.58	69.43	62.96
B73	26	1	135.53	23.53	85.18	93.18	162.82	24.47	60.71	106.82	38.59	17.88	26.35	66.35	62.59
B73	26	2	144.78	26.28	94.06	102.36	165.07	22.59	75.16	136.48	41.50	24.44	26.74	82.54	79.77
B73	26	3	199.47	29.47	107.89	117.37	193.68	26.32	95.79	168.42	45.79	25.26	31.05	96.32	87.37
B73	30	1	167.27	28.05	89.35	107.01	183.38	26.49	78.96	148.57	42.60	20.26	29.61	84.16	75.84
B73	30	2	143.66	25.13	80.42	88.38	148.27	20.52	71.62	121.88	38.53	20.94	23.87	68.27	69.53

```
186 #####
187 # Process the trait data
188 #####
189
190 datTraits[datTraits[,1] == "o2", 1] <- "O2"
191
192 datTraits = datTraits %>%
193   pivot_longer(!c(Genotype, DAP, Replicaiton), names_to = "Trait", values_to = "Measurement") %>%
194   group_by(Genotype, DAP, Trait) %>%
195   summarise_at(vars(Measurement), list(Measurement = mean)) %>%
196   ungroup() %>%
197   unite("Genotype", Genotype:DAP, sep = "_") %>%
198   pivot_wider(names_from = Trait, values_from = Measurement, values_fill = NA) %>%
199   as.data.frame(stringsAsFactors = FALSE)
200
201 # Copy values in column 1 into row names
202 rownames(datTraits) <- datTraits[,1]
203
204 # Remove column 1
205 datTraits <- datTraits[,-1]
206
207 print(head(datTraits))
208 print(tail(datTraits))
209 print(dim(datTraits))
210
```


Processed Phenotype Input File

	Ala	Arg	Asx	Glx	Gly	His
B73_10	180.706666666667	26	140.956666666667	119.85	185.97	26.0366666666667
B73_14	130.52	22.2066666666667	88	84.4333333333333	155.123333333333	23.2066666666667
B73_18	119.35	24.19	80.02	88.8866666666667	164.533333333333	26.4033333333333
B73_22	139.383333333333	22.3	82.9533333333333	88.37	148.983333333333	20.19
B73_26	159.926666666667	26.4266666666667	95.71	104.303333333333	173.856666666667	24.46
B73_30	155.113333333333	25.53	85.2366666666667	95.8733333333333	155.266666666667	21.0266666666667
B73_34	149.816666666667	26.46	89.5666666666667	98.0433333333333	149.233333333333	18.22
B73_38	159.056666666667	28.43	87.9066666666667	105.13	153.66	19.6933333333333
B73_42	145.66	27.23	79.52	93.22	149.85	19.3933333333333
B73_46	155.34	29.1933333333333	90.4733333333333	100.916666666667	148.216666666667	18.0666666666667
O2_10	156.9	25.1133333333333	147.766666666667	104.1	153.053333333333	20.3233333333333
O2_14	140.636666666667	25.1533333333333	108.94	95.9066666666667	152.013333333333	21.4766666666667
O2_18	121.946666666667	25.52	98.1166666666667	96.2066666666667	152.603333333333	22.1333333333333
O2_22	129.033333333333	23.76	100.72	89.4333333333333	136.11	18.1866666666667
O2_26	126.86	25.1266666666667	106.406666666667	90.9833333333333	143.653333333333	19.9466666666667
O2_30	189.93	29.07	118.046666666667	94.48	160.173333333333	21.6966666666667
O2_34	153.553333333333	27.1066666666667	94.4066666666667	90.9266666666667	155.993333333333	21.6533333333333
O2_38	146.48	32.3133333333333	110.6	89.6533333333333	164.26	22.39
O2_42	151.676666666667	31.7	115.133333333333	96.94	167.67	21.38
O2_46	124.713333333333	35.4833333333333	112.5	89.8166666666667	162.35	22.2466666666667

```
212 #####
213 # Save processed trait data
214 #####
215 write.csv(
216   x = datTraits,
217   file = file.path(output_path, "..", "datTraits.csv"),
218   na = "",
219   quote = FALSE
220 )
```

```
223 #####
224 # Construct heatmap to visualize trait data
225 #####
226 # Re-cluster samples
227 sampleTree2 = hclust(dist(datExpr), method = "average")
228
229 # Convert traits to a color representation: white means low, red means high, grey means missing entry
230 traitColors = numbers2colors(datTraits, signed = FALSE)
231
232 cat(rep("\n", 2))
233 jpeg(filename = file.path(output_path, "tree_trait_heatmap.jpeg"), width = 1920, height = 720)
234 # Plot the sample dendrogram and the colors underneath.
235 plotDendroAndColors(sampleTree2, traitColors, groupLabels = names(datTraits), main = "Sample dendrogram and trait heatmap")
236 dev.off()
```
