# Searching for best lower dimensional visualization angles for high dimensional RNA-Seq data

**Wanli Zhang** [Corresp., 1] , **Yanming Di** [1]

[1] Statistics, Oregon State University, Corvallis, Oregon, United States

Corresponding Author: Wanli Zhang
Email address: zhangwa@stat.oregonstate.edu

The accumulation of RNA-Seq gene expression data in recent years has resulted in large and complex data sets of high dimensions. Exploratory analysis, including data mining and visualization, reveals hidden patterns and potential outliers in such data, but is often challenged by the high dimensional nature of the data. The scatterplot matrix is a commonly used tool for visualizing multivariate data, and allows us to view multiple bivariate relationships simultaneously. However, the scatterplot matrix becomes less effective for high dimensional data because the number of bivariate displays increases quadratically with data dimensionality. In this study, we introduce a selection criterion for each bivariate scatterplot and design/implement an algorithm that automatically scan and rank all possible scatterplots, with the goal of identifying the plots in which separation between two pre-defined groups is maximized. By applying our method to a multi-experiment *Arabidopsis* RNA-Seq data set, we were able to successfully pinpoint the visualization angles where genes from two biological pathways are the most separated, as well as identify potential outliers.

# Searching for best lower dimensional visualization angles for high dimensional RNA-Seq data

Wanli Zhang[1][¶] and Yanming Di[1][¶]

[1]Department of Statistics, Oregon State University, Corvallis, Oregon, United States of America

¶ These authors contributed equally to this work.

## Abstract

The accumulation of RNA-Seq gene expression data in recent years has resulted in large and complex data sets of high dimensions. Exploratory analysis, including data mining and visualization, reveals hidden patterns and potential outliers in such data, but is often challenged by the high dimensional nature of the data. The scatterplot matrix is a commonly used tool for visualizing multivariate data, and allows us to view multiple bivariate relationships simultaneously. However, the scatterplot matrix becomes less effective for high dimensional data because the number of bivariate displays increases quadratically with data dimensionality. In this study, we introduce a selection criterion for each bivariate scatterplot and design/implement an algorithm that automatically scan and rank all possible scatterplots, with the goal of identifying the plots in which separation between two pre-defined groups is maximized. By applying our method to a multi-experiment *Arabidopsis* RNA-Seq data set, we were able to successfully pinpoint the visualization angles where genes from two biological pathways are the most separated, as well as identify potential outliers.

# 1.   INTRODUCTION

High throughput RNA sequencing (RNA-Seq) has been widely adopted for quantifying relative gene expression in comparative transcriptome analysis. In recent years, the increasing number of RNA-seq studies on the model plant *Arabidopsis thaliana* have resulted in an ever-accumulating amount of data from multiple RNA-Seq experiments. In this article, we will develop tools for the exploration and visualization of such multi-experiment data.

For examining treatment effects of individual genes under multiple conditions and across multiple experiments, a vector summarizing the differential expression (DE) results under different treatment conditions seems adequate. To visualize the DE profile under different treatments, a line plot can be used. However, since genes work interactively in all biological processes, it is of interest to examine expression patterns of groups of genes, through which the genes' biological context can be better understood. In light of this, researchers often would like to both identify the general trend and pinpoint individual aberrations in the expression profile of genes belonging to the same biological pathway,

38 as well as compare the profiles between multiple pathways.

39  When multiple genes are being examined together, the line plots are less effective for

40 visualizing DE or expression profiles: The lines often cross each other, making it difficult

41 to identify the grouping and understand the behavior of individual genes. One common

42 alternative visualization method is the scatterplot, which shows expression level under two

43 treatment conditions at a time. Scatterplots are effective in showing clustering patterns

44 and outliers, greatly assisting with data exploration (Elmqvist et al., 2008). For high

45 dimensional data, one has the option of using the scatterplot matrix, in which each panel

46 is the scatterplot for the corresponding pairs of features. However, manual scanning

47 of all possible pairwise scatterplots can be arduous or even fruitless at times, because

48 the number of possible visualization angles increases quadratically with respect to data

49 dimensionality ($p$ choose 2 possible angles).

50  In this paper, we propose to automatically search for the best low dimensional visual-

51 ization angles (2-, 3-, or 4-dimensional) based on a context-sensitive, numeric measure of

52 importance, thereby reducing the amount of effort invested in scatterplot scanning. In our

53 study, we hope to explore the patterns and differences in gene expression profile between

54 two phytohormone signaling pathways, and therefore, we would like the top ranked scat-

55 terplots to contain as much information as possible on pathway classification. We thus

56 define such an importance measure for the dimension subsets such that the scatterplots

57 will show the largest separation between different pre-defined groups in the data set.

58  For this study, we will look for feature subsets upon which the pathways ethylene (ET)

59 and jasmonate (JA) are the most separated, and quantify the between-group separation by

60 calculating the repeated cross-validation (RCV) error of misclassification using MclustDA

61 (Fraley and Raftery, 2002), a model-based classification method. In Figure 1, we show

62 one of the top ranked 2-subset feature combinations that give the greatest separation

63 between two pathways, as well as subset giving the smallest separation. Comparing the

64 two scatterplots, we can observe that the two pathway groups in Figure 1(a) are more

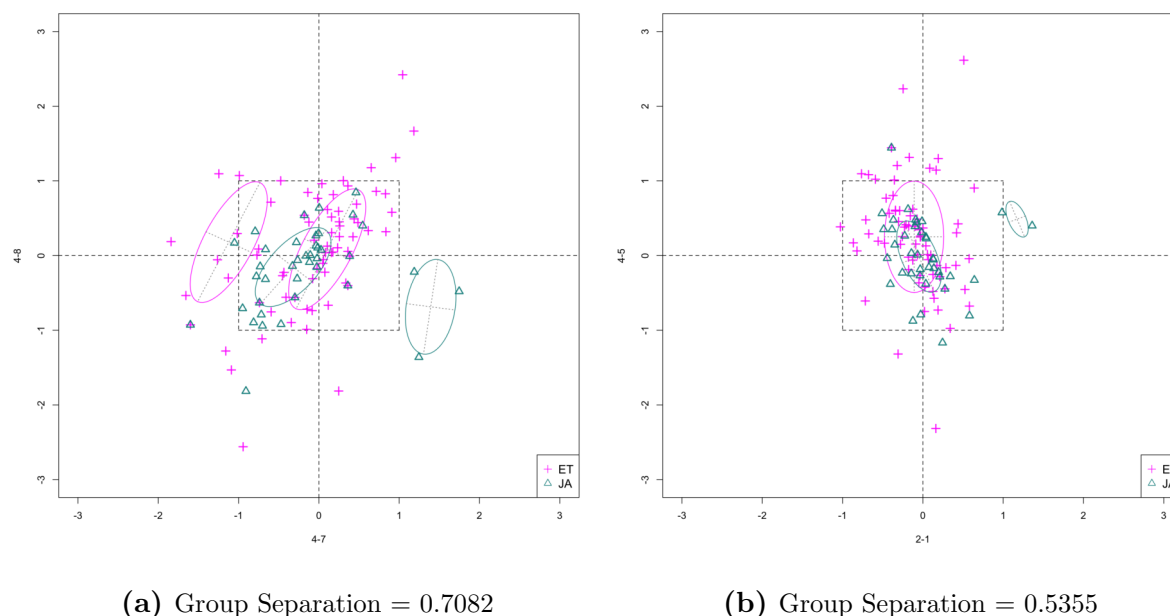65 visually distinguishable than those in Figure 1.

**(a)** Group Separation = 0.7082      **(b)** Group Separation = 0.5355

**Figure 1:** Scatterplots of 2-dim feature subsets reflecting maximum and minimum group separations. Dashed-line square marks $\pm 1$ range from the origin. Different classes distinguished with color. Ellipses correspond to component mean and covariance fitted by MclustDA. Treatment $i$-$j$ represents the $j$th treatment in experiment $i$.

66      The rest of the paper is formatted as follows: Section 2 outlines the collection and

67   processing of the data and information on the experiments and biological pathways. The

68   statistical methods are described in Section 3. In Section 4, we list the results obtained by

69   applying our method to the collected data. Finally, we state our conclusion and discuss

70   limitations and possibilities for future work in Section 5. Additional proofs and graphs

71   are included in the Appendix.

## 2.   DATA DESCRIPTION AND PROCESSING

### 2.1   Collecting experimental data

72   In this study, we use the data collected and processed by Bin Zhuo (Zhuo et al., 2016). All

73   5 datasets were acquired from the National Center for Biotechnology Information (NCBI)

74   website `www.ncbi.nlm.nih.gov` and processed through a customized assembly pipeline

75   to obtain a matrix of counts for genes in samples. All datasets originate from RNA-Seq

76   experiments on the model plant *Arabidopsis thaliana*, with treatment conditions (includ-

77 ing genetic variants) varying between experiments. All experiments were conducted on

78 the leaf tissue. The number of treatments/factor levels also vary among the experiments.

79 The GEO (Gene Expression Omnibus) accession numbers (which can be used to directly

80 search for the experiment/dataset information) are available as part of the meta-data,

81 and the assembly pipeline is described in detail in Zhuo et al. (2016). We have included

82 the basic information on the experiments in Table 1.

| ID | GEO accession # | Title |
|----|-----------------|-------|
| 1 | GSE36626 | Dynamic Deposition of the Histone H3.3 Variant Accompanies Developmental Remodeling of Arabidopsis Transcriptome (mRNA-Seq) |
| 2 | GSE39463 | Time-course RNA-seq analysis of the barley MLA1 immune receptor-mediated response to barley powdery mildew fungus Bgh in Arabidopsis thaliana |
| 3 | GSE48235 | Four distinct types of dehydration stress memory genes in Arabidopsis thaliana |
| 4 | GSE51304 | Non-CG methylation patterns shape the epigenetic landscape in Arabidopsis |
| 5 | GSE54677 | Transcriptional gene silencing by Arabidopsis Microrchidia homologues involves the formation of heteromers |

**Table 1:** Experiment information

83 For each of the 5 experiments, a negative binomial regression model is fitted to the

84 normalized counts, where the normalization factors are computed using the 104 genes

85 shared by top 1000 most stably expressed genes in three tissue groups (Zhuo et al., 2016).

86 After removing the columns corresponding to the baseline expression levels of the control

87 groups, the resulting matrix summarizes the log (base 2) fold changes under different

88 treatments across the 5 experiments: Each column represents the log fold changes of gene

89 expression between one treatment group (or a gene knockout mutant) and the control

90 group (or wildtype) in one of the experiments.

## 2.2 Finding pathway genes

For this study, we focus our attention on the signaling pathways of two phytohormones: ethylene (ET) and jasmonic acid (JA). As a plant hormone, ethylene is commercially important due to its regulation on fruit ripening (Lin et al., 2009). JA acts as a key cellular signal involved in the activation of immune responses to most insect herbivores and necrotrophic microorganisms (Ballaré, 2010).

For each pathway, we first use AmiGO 2 (`http://amigo.geneontology.org/amigo/landing`) to search for the list of genes involved, and then identify the subset of genes in our data set that belong to the pathway through cross-reference. Genes with a fold change of $< 2$ under all treatment conditions are filtered out. The name, GO accession number, and the number of genes in each pathway are listed in Table 2.

| ID | Pathway name | GO accession # | # Genes |
|----|--------------|----------------|---------|
| ET | Ethylene-activated signaling pathway | GO: 0009873 | 86 |
| JA | Jasmonic acid mediated signaling pathway | GO: 0009867 | 48 |

**Table 2:** Pathway information

In Figure 2, we display the expression profile of genes that belong to each pathway group. Under certain individual treatment-control contrasts (e.g. 2-3, 4-3, 5-1), there exist observable similarities between the distribution of expression levels, while it is more difficult to tell under other treatments.
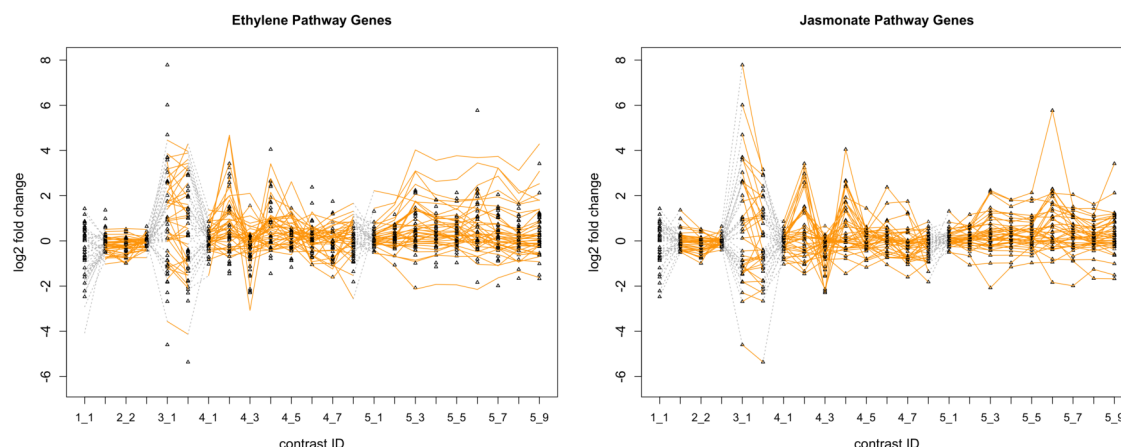
**Figure 2:** Gene expression profile plot for pathways ET and JA. Treatments from the same experiment are joined by orange lines. Different experiments are joined by grey dashed lines. Feature $i$-$j$ represents the $j$th treatment-control contrast in experiment $i$.

# 3. METHOD

## 3.1 Mixture discriminant analysis via MclustDA

105 In this section, we will start by introducing a classification method named MclustDA, and
106 then define a measure for group separation using cross-validation results with MclustDA.
107 Finally, we lay out our strategy for reducing data dimensionality with the ultimate goal
108 of simplifying navigation of scatterplots.

109

110 **MclustDA model**    In discriminant analysis (DA), known classifications of some ob-
111 servations are used to classify others. The number of classes, $G$, is assumed to be known.
112 For probabilistic DA methods, it is assumed that observations in class $k$ follow a class
113 specific probability distribution $f_k(\cdot)$. Let $\tau_k$ represent the proportion of observations in
114 class $k$. According to Bayes's theorem, it follows that

$$P(\boldsymbol{y} \in \text{class } j) = \frac{\tau_j f_j(\boldsymbol{y})}{\sum_{k=1}^{G} \tau_k f_k(\boldsymbol{y})},$$

115 where observation $\boldsymbol{y}$ is assigned to the most probable class.
116    Commonly used DA methods, including Fisher's linear discriminant analysis (LDA)
117 and quadratic discriminant analysis (QDA), assume a multivariate normal density for
118 each class:

$$f_k(\boldsymbol{y}) = \phi(\boldsymbol{y}|\mu_k, \Sigma_k).$$

119 The method is called LDA if the covariance matrices for all classes coincide ($\Sigma_k = \Sigma$
120 for $k = 1, ..., G$), and is called QDA if the class covariances are allowed to vary.

121 MclustDA (Fraley and Raftery, 2002), an extension and generalization to LDA and
122 QDA, models each class density as a mixture of multivariate normals. The density for
123 class $j$ is as follows:

$$f_j(\boldsymbol{y}|\theta_k) = \sum_{k=1}^{G_j} \tau_{jk}\phi(\boldsymbol{y}|\mu_{jk}, \Sigma_{jk}),$$

124 where $G_j$ is the number of components for class $j$, $\{\tau_{jk}\}$ are mixing proportions for com-
125 ponents in class $j$, and $\theta_k$ is the vector of parameters for the normal mixture. Component
126 covariances $\Sigma_{jk}$ are allowed to vary both within and between classes.

127 Parameters within each class are separately estimated by maximum likelihood via the
128 EM algorithm (Dempster et al., 1977), which is equivalent to fitting a Mclust (Fraley and
129 Raftery, 2002) model for each class. And just like Mclust, MclustDA performs model
130 selection within each class for the number of mixture components as well as covariance
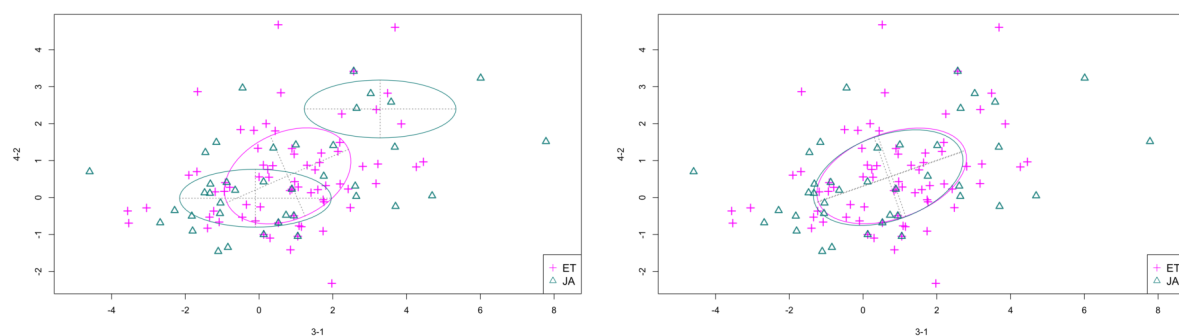131 matrix parameterizations with Bayesian information criterion (Schwarz, 1978).

132

133 **Comparison with LDA**    In our study, MclustDA is chosen over LDA/QDA as the
134 classifier due to its greater flexibility in describing the data. In RNA-Seq analysis, we
135 typically assume that the majority of genes are not differentially expressed, and therefore
136 we expect to see a cluster of points around the origin. Since MclustDA proposes to fit
137 more than one normal component to each class, it readily captures the cluster of non-DE
138 genes as well as any abnormalities that might be of interest.

139 In Figure 3, we fitted a MclustDA model and a LDA model on dimensions [3-1, 4-
140 2] of our data, separately. In MclustDA fit, each class is described with a mixture of
141 two bivariate normal components, with the ellipses representing fitted covariance matrix
142 estimates. For details in how the ellipses are constructed, see Appendix A.

143 Class JA is fitted with a component centered near the origin, representing genes
144 with low expression levels under both treatments, as well as a component centered at

$_{145}$ $(2.276, 1.663)$ that encompasses relatively active genes. Class ET is represented by a

$_{146}$ single normal component centered at $(0.537, 0.406)$.

$_{147}$     In comparison, due to model assumptions, LDA fitted a bivariate normal density to

$_{148}$ each class with covariances being equal, and in this case, the estimated centers almost

$_{149}$ coincide with each other. The fitted normal densities are only able to capture the general

$_{150}$ shape and orientation of each class, while MclustDA provides us with a more detailed

$_{151}$ anatomy of geometric and distributional properties in each class.



**(a)** MclustDA fit with 1 and 2 components in each class    **(b)** LDA model assuming equal covariance matrix for each class

**Figure 3:** Comparison of MclustDA and LDA fit of the same data. Fitted components and points from different classes are distinguished with color. Ellipses correspond to component covariances.

## 3.2   Quantification of group separation

$_{152}$ Our definition of group separation measure is motivated by the relationship between

$_{153}$ visualized separation and misclassification probability (from a model-based classifier).

$_{154}$     Suppose we wish to separate two populations $\pi_1$ and $\pi_2$. Let $X = [X_1, ..., X_p]$ denote

$_{155}$ the $p$-dimensional measurement vector of an observation. We assume that densities $f_1(x)$

$_{156}$ and $f_2(x)$ describe the variability of the two populations. Let $p_1$ and $p_2$ denote prior

$_{157}$ probability of each population. Define $c(1|2)$ and $c(2|1)$ as costs of misclassifying an

$_{158}$ object from class $2(1)$ as class $1(2)$. Here we let $c(1|2) = c(2|1) = 1$ to simplify the

$_{159}$ formulation. Let $\Omega$ denote the entire sample space, and $\Omega = R_1 \cup R_2$, where $R_1$ is the set

$_{160}$ of values of $x$ for which we classify objects into $\pi_1$, and $R_2 = \Omega - R_1$.

$_{161}$     The probability of misclassifying an object from $\pi_1$ as $\pi_2$ is:

$$P(2|1) = P(X \in R_2|\pi_1) = \int_{R_2} f_1(x)\mathrm{d}x,$$

162 and similarly, we have

$$P(1|2) = P(X \in R_1|\pi_2) = \int_{R_1} f_2(x)\mathrm{d}x.$$

163 By definition, we can calculate the probability of misclassifying any object:

$$P(\text{misclassified as } \pi_1) = P(X \in R_1|\pi_2)P(\pi_2) = P(1|2)p_2,$$

$$P(\text{misclassified as } \pi_2) = P(X \in R_2|\pi_1)P(\pi_1) = P(2|1)p_1.$$

164 The Total Probability of Misclassification (TPM) is defined as the probability of either
165 misclassifying a $\pi_1$ object or misclassifying a $\pi_2$ object, i.e.

$$\text{TPM} = p_1 P(2|1) + p_2 P(1|2). \tag{1}$$

166 Suppose $Y = \{Y_1, ..., Y_{N_1}\} \sim \pi_1$ and $Z = \{Z_1, ..., Z_{N_2}\} \sim \pi_2$ are two i.i.d samples from
167 the two populations. Assume that a classification system has been trained and tested on
168 this data set, and results in the following confusion matrix in Table 3:

**Predicted Class**

|  |  | $\pi_1$ | $\pi_2$ | **total** |
|---|---|---|---|---|
| **Actual Class** | $\pi_1$ | $n_{1|1}$ | $n_{2|1}$ | $N_1$ |
|  | $\pi_2$ | $n_{1|2}$ | $n_{2|2}$ | $N_2$ |
|  | **total** | $N_1'$ | $N_2'$ |  |

**Table 3:** Confusion matrix

169 Then the misclassification error rate (MER), i.e. probability of misclassifying any
170 object, is given by:

$$\text{MER} = \frac{n_{1|2} + n_{2|1}}{N_1 + N_2} = \frac{n_{1|2}}{N_2} \cdot \frac{N_2}{N_1 + N_2} + \frac{n_{2|1}}{N_1} \cdot \frac{N_1}{N_1 + N_2}. \tag{2}$$

171  Under the assumption that each object is independently classified, the number of mis-

172  classified $\pi_1$ objects, $N_{2|1}$, follows a Binomial distribution with parameters $(N_1, P(2|1))$.

173  Likewise, the number of misclassified $\pi_2$ objects, $N_{1|2}$, follows a Binomial distribution

174  with parameters $(N_2, P(1|2))$. The maximum likelihood (ML) estimators for $P(2|1)$ and

175  $P(1|2)$ can be easily computed:

$$\widehat{P(2|1)} = \frac{n_{2|1}}{N_1}; \quad \widehat{P(1|2)} = \frac{n_{1|2}}{N_2}.$$

176  Now, if we set $p_1 = N_1/(N_1 + N_2)$ and $p_2 = N_2/(N_1 + N_2)$ as prior probabilities for $\pi_1$

177  and $\pi_2$, then under independence assumption, it follows that

$$\text{MER} = p_1 \widehat{P(2|1)} + p_2 \widehat{P(1|2)},$$

178  that is, MER is a maximum likelihood, and hence consistent, estimate of TPM.

179  In practice, however, the MER tends to underestimate TPM because the same data

180  has been used for both training and testing. In this study, we use cross-validation to

181  address this issue.

182

183  **Repeated stratified cross-validation**    One of the most commonly used method to

184  estimate the expected error rate is cross-validation (CV). For a $K$-fold CV, the original

185  data is randomly split into $K$ equally sized subsamples, of which $K - 1$ (training set) are

186  used to train a classifier and the remaining one (validation set) is used to test the trained

187  classifier. For a binary classification problem, the misclassification error rate (MER), as

188  defined in (2), is typically computed using the validation set as a performance measure for

189  the classifier. The training-validation process is iterated over all $K$ folds, each time using

190  a different subsample as validation set, and the resulting $K$ MER values are averaged. In

191  stratified cross-validation, the folds are selected so that they contain approximately the

192  same proportion of classes as the original data. It has been shown in previous studies

193  that stratified CV tends to perform uniformly better than CV, in terms of both bias and

194  variance (Kohavi, 1995).

195  Due to the randomness in partitioning the sample into $K$ folds, we have introduced

196  variation into the $K$-fold CV estimator. One way to reduce this variation is to repeat the

197  whole cross-validation process multiple times using different pseudorandom allocations

198 of instances to training and validation folds for each repetition (Kim, 2009), and report

199 the average of CV estimators across all repetitions. This method is often referred to

200 as the repeated cross-validation (RCV). For improved repeatability of results, common

201 seeding has been recommended in earlier studies (Powers and Atyabi, 2012). In our

202 implementation, we set a fixed random number seed for each repetition of CV.

203 Let $C \times K$-CV denote a $K$-fold CV with $C$ repetitions. There has been much discus-

204 sion on the optimal choice of $C$ and $K$ (Kohavi, 1995; Kim, 2009; Powers and Atyabi,

205 2012). Increasing $C$ tends to decrease the variance of the RCV estimator, but at the same

206 time increases the computational time. The choice of $K$ takes into account the tradeoff

207 between bias and variance of the CV estimator (of the expected error rate). For small

208 $K$, less data is used to train the classifier and therefore the error estimate tends to be

209 biased. For large $K$, the estimator becomes less biased due to more data being used in

210 training, but its variance is inflated due to higher correlation between different training

211 folds. Kohavi (1995) recommends using a stratified 10-fold CV with multiple runs, and

212 we chose $C = 10$ considering the amount of computation required as well as the specs of

213 our hardware.

214

215 **Quantify group separation**    We define the group separation index (GSI) as

$$\text{GSI} = 1 - \hat{\epsilon}_{\text{rcv}}, \tag{3}$$

216 where $\hat{\epsilon}_{\text{rcv}}$ denotes the repeated stratified CV estimator of the total misclassification prob-

217 ability using MclustDA as the classifier.

218 Intuitively, for a chosen feature subset, a small CV error indicates that the data can be

219 more easily classified when projected onto these dimensions, which, in our expectation,

220 can be reflected in the graphical representation of the data by showing that different

221 classes can be more easily distinguished through simple visualization.

## 3.3 Feature subset selection via GSI ranking

222 In this section, we describe the data in each pathway with a low dimensional representation

223 for easier interpretation by selecting a parsimonious subset of features (treatments) that

224 contain as much information on pathway classification/separation as possible. In other

225 words, we hope to find the dimensions to project the data onto such that the separation

226 between two pathways is as large as possible. We use GSI, as defined in (3), to measure

227 the separation between two pathway groups.

228     In order to find the optimal feature subset in terms of group separation, we designed

229 and implemented the following algorithm:

230     **Step 1**: Determine the number of features $M$ to keep. Choose $M$ from $\{2, 3, 4\}$.

231     **Step 2**: List all $M$-subsets of features exhaustively. Call this collection of subsets

232     $\mathcal{F}_M$.

233     **Step 3**: For each member of $\mathcal{F}_M$, subset the data accordingly. Calculate and

234     record a $10 \times 10$ stratified CV error rate (and equivalently, GSI) with MclustDA as

235     classifier on each subsetted data. For each fold of CV, use misclassification error

236     rate as measure of fit.

237        – CV model fitting: First fit a MclustDA model to the entire subsetted data,

238          setting maximum number of components as 2. Then use the same fitted model

239          (number of components, parameterization) for every fold of CV.

240     **Step 4**: Rank the feature subsets in $\mathcal{F}_M$ according to their GSI values. Feature

241     subsets with higher GSI values are ranked higher.

242     **Step 5**: Repeat above steps for other values of $M$.

243     For the purpose of finding "good" angles for data visualization, we will examine the

244 scatterplots and scatterplot matrices generated by top-ranked feature subsets. The results

245 will be discussed in Section 4.

246

247 **Random number seed**    To ensure reproducibility of our results, for each of 2-, 3-

248 and 4-subset selection process, we followed the following protocol to set random number

249 seeds:

250     **Step 1**: Choose a list of 50 random number seeds. Partition the list into 5 batches

251     of 10 seeds.

252     **Step 2**: For each feature subset, run 10-fold stratified CV for 50 times, each time

253     using a different seed from the list.

254 **Step 3**: Average results within each of 5 batches of 10 random seeds to obtain

255 $10 \times 10$ stratified CV result. For instance, average of seeds $1\sim10$ results serves as

256 first run of $10 \times 10$ RCV; average of seeds $11\sim20$ serves as second run, etc.

## 3.4 Dimension reduction via PCA

257 Principal component analysis (PCA) maps the data onto a lower dimensional space in such

258 a way that the variance of the data in the low-dimensional representation is maximized.

259 As a dimension reduction technique, usually only the first few principal components (PCs)

260 are used. Despite its popularity in the field of data visualization, the formulation of PCA

261 does not involve any class information in the data, which implies that the projected

262 directions corresponding to the largest variance may not contain the best separability

263 information.

264 To verify this observation, using the expression data from all 5 experiments, we calcu-

265 lated its principal components, and treat them as the new (projected) features. Then for

266 the first 2, 3 and 4 PCs, respectively, we calculated the group separation index for each

267 case using $10 \times 10$-CV with MclustDA and compare the results with ours.

# 4. RESULTS

## 4.1 Repeated cross-validation with MclustDA

268 With the secondary purpose of testing the stability of repeated CV, we executed multiple

269 runs for each of the 2-, 3-, and 4-subset feature selection procedures. The top ranked

270 feature subsets as well as their corresponding GSI values are presented in Tables $4\sim6$.

**Table 4:** Top ranked 2-subsets from multiple runs of $10 \times 10$ RCV. Ties are marked with asterisk ($*$). Combinations appearing in all 5 runs are highlighted with distinguishing colors.

**(a)** Run 1

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [4-7, 4-8] | 0.698 |
| 2 | [3-1, 3-2] | 0.694 |
| 3 | [4-7, 5-2] | 0.688 |
| 4 | [3-1, 5-2] | 0.681 |
| 5 | [2-1, 4-7] | 0.680 |

**(b)** Run 2

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [4-7, 4-8] | 0.703 |
| 2 | [3-1, 3-2] | 0.696 |
| 3 | [4-7, 5-2] | 0.690 |
| 4 | [3-1, 5-8] | 0.682 |
| 5 | [3-1, 5-2] | 0.681 |

**(c)** Run 3

| Rank | Subset | GSI |
|------|--------|-----|
| 1* | [3-1, 3-2] | 0.697 |
| 2* | [4-7, 5-2] | 0.697 |
| 3 | [4-7, 4-8] | 0.689 |
| 4 | [3-1, 5-2] | 0.685 |
| 5 | [2-1, 4-7] | 0.675 |

**(d)** Run 4

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [3-1, 3-2] | 0.704 |
| 2 | [4-7, 4-8] | 0.695 |
| 3 | [3-1, 5-2] | 0.687 |
| 4 | [3-1, 5-8] | 0.684 |
| 5 | [4-7, 5-2] | 0.681 |

**(e)** Run 5

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [4-7, 4-8] | 0.708 |
| 2 | [3-1, 3-2] | 0.702 |
| 3 | [4-7, 5-2] | 0.689 |
| 4 | [4-4, 5-2] | 0.688 |
| 5 | [3-1, 5-2] | 0.683 |

**Table 5:** Top ranked 3-subsets from multiple runs of $10 \times 10$ RCV. Ties are marked with asterisk ($*$). Combinations appearing in all 5 runs are highlighted with distinguishing colors.

**(a)** Run 1

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [3-2, 4-7, 5-2] | 0.724 |
| 2 | [3-2, 5-2, 5-6] | 0.720 |
| 3 | [1-1, 3-1, 5-2] | 0.718 |
| 4 | [3-1, 5-2, 5-6] | 0.710 |
| 5 | [2-2, 2-3, 3-1] | 0.707 |

**(b)** Run 2

| Rank | Subset | GSI |
|------|--------|-----|
| 1* | [1-1, 3-1, 5-2] | 0.717 |
| 2* | [3-2, 4-7, 5-2] | 0.717 |
| 3 | [3-2, 5-2, 5-6] | 0.713 |
| 4 | [1-1, 2-2, 3-1] | 0.708 |
| 5 | [2-2, 2-3, 3-1] | 0.708 |

**(c)** Run 3

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [3-2, 5-2, 5-6] | 0.725 |
| 2 | [1-1, 3-1, 5-2] | 0.717 |
| 3* | [3-1, 5-2, 5-9] | 0.715 |
| 4* | [3-2, 4-7, 5-2] | 0.715 |
| 5 | [3-1, 4-7, 5-2] | 0.714 |

**(d)** Run 4

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [3-2, 5-2, 5-6] | 0.736 |
| 2 | [1-1, 3-1, 5-2] | 0.728 |
| 3 | [2-2, 2-3, 3-1] | 0.713 |
| 4* | [3-1, 4-7, 5-2] | 0.712 |
| 5* | [3-2, 4-7, 5-2] | 0.712 |

**(e)** Run 5

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [1-1, 3-1, 5-2] | 0.726 |
| 2 | [3-2, 4-7, 5-2] | 0.724 |
| 3 | [3-2, 5-2, 5-6] | 0.720 |
| 4 | [2-3, 3-1, 5-2] | 0.715 |
| 5 | [2-2, 2-3, 3-1] | 0.710 |

**Table 6:** Top ranked 4-subsets from multiple runs of $10 \times 10$ RCV. Ties are marked with asterisk ($*$). Combinations appearing in all 5 runs are highlighted with distinguishing colors.

**(a)** Run 1

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [1-1, 4-2, 4-4, 5-2] | 0.730 |
| 2 | [4-4, 4-5, 4-7, 5-8] | 0.728 |
| 3 | [3-1, 4-4, 4-5, 5-2] | 0.727 |
| 4 | [4-1, 4-7, 4-8, 5-1] | 0.725 |
| 5 | [3-2, 4-5, 4-7, 5-7] | 0.724 |

**(b)** Run 2

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [1-1, 4-2, 4-4, 5-2] | 0.745 |
| 2 | [3-1, 4-4, 4-5, 5-2] | 0.741 |
| 3 | [2-3, 3-1, 4-4, 4-5] | 0.734 |
| 4 | [2-3, 3-1, 4-7, 5-2] | 0.725 |
| 5 | [3-2, 4-5, 4-7, 5-7] | 0.724 |

**(c)** Run 3

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [1-1, 4-2, 4-4, 5-2] | 0.735 |
| 2 | [4-4, 4-5, 4-7, 5-8] | 0.731 |
| 3* | [3-2, 4-5, 4-7, 5-7] | 0.726 |
| 4* | [4-1, 4-7, 4-8, 5-1] | 0.726 |
| 5 | [3-1, 4-4, 4-5, 5-2] | 0.725 |

**(d)** Run 4

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [1-1, 4-2, 4-4, 5-2] | 0.739 |
| 2* | [3-1, 4-4, 4-5, 5-2] | 0.731 |
| 3* | [3-2, 4-5, 4-7, 5-7] | 0.731 |
| 4 | [4-4, 4-5, 4-7, 5-8] | 0.723 |
| 5 | [3-1, 3-2, 4-7, 5-2] | 0.722 |

**(e)** Run 5

| Rank | Subset | GSI |
|------|--------|-----|
| 1 | [1-1, 4-2, 4-4, 5-2] | 0.740 |
| 2 | [3-1, 4-4, 4-5, 5-2] | 0.735 |
| 3* | [3-2, 4-5, 4-7, 5-7] | 0.730 |
| 4* | [2-2, 2-3, 3-1, 4-7] | 0.730 |
| 5 | [3-2, 4-6, 4-7, 5-6] | 0.723 |

### 4.1.1 Stability of RCV model selection results

Although the top ranked feature subsets sometimes differ between multiple RCV runs, we are still able to observe high degree of overlap between the results:

273    For 4-subset (Table 6), [1-1, 4-2, 4-4, 5-2], [3-1, 4-4, 4-5, 5-2] and [3-2, 4-5, 4-7, 5-7]

274 are among top ranked feature combinations in all 5 runs.

275    For 3-subset (Table 5), feature combinations [3-2, 4-7, 5-2], [3-2, 5-2, 5-6] and [1-1,

276 3-1, 5-2] are ranked top for all 5 runs.

277    For 2-subset (Table 4), [4-7, 4-8], [3-1, 3-2], [4-7, 5-2] and [3-1, 5-2] are among top

278 ranked feature combinations for all runs.

### 4.1.2   Top Ranked Scatterplot: Same Experiment

279 In Figure 4, we show the scatterplot of the data projected onto dimensions [4-7, 4-8], one

280 of the top ranked 2-subset feature combinations. These two features originate from the

281 same experiment.

**Figure 4:** Scatterplot of data projected on dimensions 4-7 and 4-8. Pathways are distinguished with color. Ellipses represent estimated covariances fitted by MclustDA. Potential outliers highlighted and labeled with their names. Dashed-line square is $\pm \log(2)$ range from the origin.

**Experiment 4**     Since both features originate from the same experiment, we will focus on the context of this experiment and first present some background information. The purpose of Experiment 4 is to characterize non-CG methylation and its interaction with histone methylation in *Arabidopsis thaliana*(Stroud et al., 2014). Non-CG methylation is a category of DNA methylation, where methyl groups are added to the DNA molecule, altering its chemical structure and thereby changing its activity. DNA methylation is usually catalyzed by DNA methyltransferases (MTases), which transfer and covalently bind methyl groups to DNA. In *Arabidopsis*, the principal DNA MTases include chromomethy-

<sup>290</sup> lase (CMT) and domains rearranged MTase (DRM) proteins, in particular CMT3 and

<sup>291</sup> DRM2. Expression of DRM1 is scarcely detected, while the function of CMT2 has not

<sup>292</sup> been studied as well as that of CMT3.

<sup>293</sup>     Histone methylation is a process by which methyl groups are transferred to amino acids

<sup>294</sup> of histone proteins. Histone methylation can either increase or decrease gene transcription,

<sup>295</sup> depending on which amino acids are methylated and the degree of methylation. The

<sup>296</sup> methylation process is most commonly observed on lysine residues (K) of histone tails H3

<sup>297</sup> and H4, among which H3K9 (lysine residue at 9th position on H3) serves as a common site

<sup>298</sup> for gene inactivation. Lysine methylation requires a specific MTase, usually containing an

<sup>299</sup> evolutionarily conserved SET domain. In *Arabidopsis*, Su(var)3-9 homologue 4 (SUVH

<sup>300</sup> 4), SUVH 5 and SUVH 6 are the major H3K9 MTases.

<sup>301</sup>     Feature 4-7 corresponds to the *drm1 drm2 cmt2 cmt3* quadruple gene knockout mu-

<sup>302</sup> tant, created by crossing *cmt2* to *cmt3* and *drm1 drm2* double mutants. It was found

<sup>303</sup> that non-CG methylation was eliminated in such mutants, indicating that DRM1, DRM2,

<sup>304</sup> CMT2 and CMT3 proteins are collectively responsible for all non-CG methylation in *Ara-*

<sup>305</sup> *bidopsis*. Feature 4-8 corresponds to the *suvh4 suvh5 suvh6* triple mutant. The control

<sup>306</sup> group of this experiment corresponds to wildtype *Arabidopsis*. Table 7 summarizes the

<sup>307</sup> above information.

| Feature ID | Sample GEO accession # | Description |
|---|---|---|
| 4-0 (control) | GSM1242374, GSM1242375 | Wildtype |
| 4-7 | GSM1242388, GSM1242389 | *drm1 drm2 cmt2 cmt3* quadruple mutant |
| 4-8 | GSM1242390, GSM1242391 | *suvh4 suvh5 suvh6* triple mutant |

**Table 7:** Feature information

<sup>308</sup> **Outliers**     Potential outliers from JA pathway, as highlighted and labeled in the scat-

<sup>309</sup> terplot, fall into the fourth quadrant, which implicates that these genes are up-regulated

<sup>310</sup> under 4-7 (DNA methylation) but down-regulated under 4-8 (histone methylation). In-

<sup>311</sup> formation on these genes is collected from TAIR (Berardini et al., 2015) and displayed in

<sup>312</sup> Table 8. One interesting discovery we made was that one of the outliers, **AT3G56400**,

<sup>313</sup> functions as a repressor of JA-regulated genes. In other words, its gene product inhibits

<sup>314</sup> the expression of other genes related to JA regulation.

| Gene name | Description |
|---|---|
| AT5G44210 | encodes a member of the ERF (ethylene response factor) subfamily B-1 of ERF/AP2 transcription factor family (ATERF-9) |
| AT2G44840 | Same function as AT5G44210; Cell-to-cell mobile mRNA |
| AT5G26170 | WRKY Transcription Factor, Group II-c; Involved in jasmonic acid inducible defense responses. |
| AT3G56400 | WRKY Transcription Factor, Group III; Repressor of JA-regulated genes; Activator of SA-dependent defense genes. |
| AT1G28400 | GATA zinc finger protein |

**Table 8:** Outlier information

³¹⁵ **Pattern differences**     The first thing we can observe from the scatterplot is that a

³¹⁶ majority of genes are expressed at a low level (with fold change $< 2$) under both treatment

³¹⁷ conditions, as demonstrated by the clustered points inside $\pm 1$ square. Although most

³¹⁸ genes are expressed at a relatively low level, we are still able to identify the difference

³¹⁹ between the two pathways. If a differential expression (DE) analysis is performed and

³²⁰ only DE genes are included in our model, it will be less likely for us to spot the same

³²¹ structural difference as before because we would lose much group level information by

³²² filtering out non-DE genes.

³²³     Secondly, not considering the outliers, genes belonging to the JA pathway are mostly

³²⁴ concentrated around the origin as well as in quadrant III, meaning that most JA genes

³²⁵ are down-regulated under both treatments. The expression pattern of ET pathway genes,

³²⁶ however, is more diverse than that of JA genes. These genes populate all four quadrants

³²⁷ of the coordinate system, with the highest density in quadrant I followed by quadrant II

³²⁸ and III. That is, a majority of ET genes are up-regulated under both treatments, while

³²⁹ most of the others are down-regulated under 4-7.

### 4.1.3   Top Ranked Scatterplot: Different Experiments

In Figure 5, we show the scatterplot of another top ranked feature combination, [3-1, 5-2], which come from two different experiments.
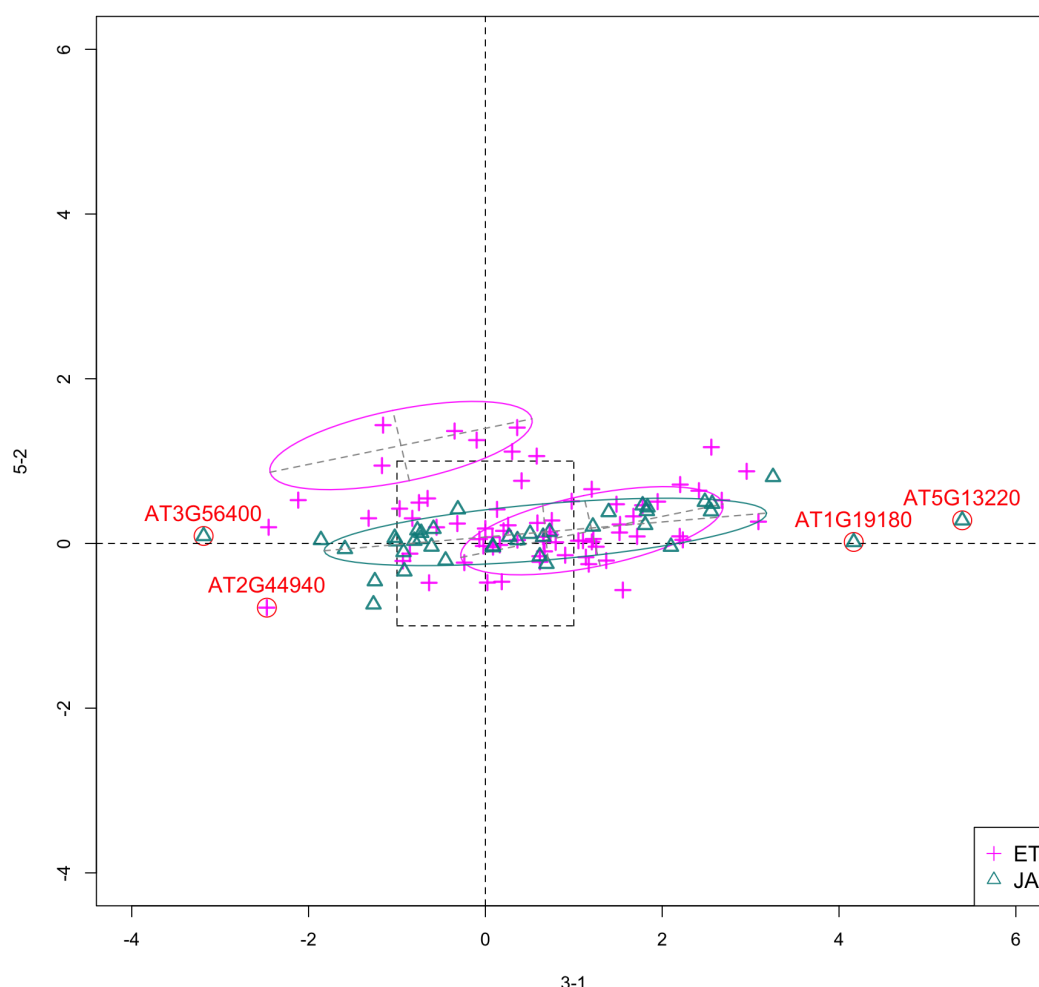


**Figure 5:** Scatterplot of data projected on dimensions 3-1 and 5-2. Pathways are distinguished with color. Ellipses represent estimated covariances fitted by MclustDA. Potential outliers highlighted and labeled with their names. Dashed-line square is $\pm \log(2)$ range from the origin.

**Experiment 3**   The focus of this study is the response of *Arabidopsis* to multiple consecutive dehydration stresses (Ding et al., 2013). Based on the observation that pre-exposure to abiotic stresses (including dehydration) may alter plants subsequent responses by improving resistance to future exposures, the researchers hypothesized the existence of "memory genes": genes that provide altered response to subsequent stresses (Ding et al.,

337   2012).

338      A RNA-Seq study is performed to determine the transcriptional responses of *Arabidop-*
339   *sis* plants that have experienced multiple exposures to dehydration stress and compare
340   them with the transcriptional behavior of plants encountering the stress for the first time.
341   The dehydration treatments are applied in the following fashion:

342   (1) Plants were removed from soil and air-dried for 2h. Call this exposure Stress 1 (S1).

343   (2) Plants were then rehydrated for 22h by being placed in humid chambers with their
344        roots in a few drops of water. Call this step Recovery 1 (R1).

345   (3) Air-dry R1 plants for 2h. This is Stress 2 (S2), followed by R2, which is the same
346        as R1.

347   (4) Air-dry R2 plants for 2h. This is Stress 3 (S3).

348      RNA-Seq analyses were then performed on leave tissues from pre-stressed/watered
349   plants (control), S1 plants and S3 plants. For each treatment group, plants from two
350   independent biological samples were used. In our data, feature 3-1 corresponds to S1, or
351   first drought stress. See Table 9 for a summary.

352

353   **Experiment 5**      In this study, the researchers examine the functional relationship be-
354   tween members of the *Arabidopsis* microrchidia (AtMORC) ATPase family (Moissiard
355   et al., 2014), which have been shown to be involved in transposon repression and gene
356   silencing. Three of seven MORC homologs were examined: AtMORC1, AtMORC2 and
357   AtMORC6. RNA-Seq experiment using single and double mutants indicates that At-
358   MORC1 and AtMORC2 act redundantly in gene silencing. Wildtype *Arabidopsis* was
359   used as control group. Treatment groups include both single and double mutant lines:
360   *atmorc2-1*, *atmorc2-4*, *atmorc1-2*, *atmorc1-5*, and *atmorc1-2 atmorc2-1*, in which two in-
361   dividual alleles were used for *atmorc1* and *atmorc2*. In our data, feature 5-2 corresponds
362   to the single mutant line *atmorc2-1*. Table 9 includes summary information on this ex-
363   periment.

| Feature ID | Sample GEO accession # | Description |
|---|---|---|
| 3-0 (control) | GSM1173202, GSM1173203 | Watered condition |
| 3-1 | GSM1173204, GSM1173205 | First drought stress |
| 5-0 (control) | GSM1321694, GSM1321704 | Wildtype |
| 5-2 | GSM1321696, GSM1321706 | *atmorc2-1* mutant |

**Table 9:** Feature information for experiments 3 and 5

₃₆₄ **Outliers**     In Figure 5, we highlighted a few observations considered as outlying, and

₃₆₅ as before, looked up their information using TAIR. A brief description for each outlier is

₃₆₆ included in Table 10. Gene **AT3G56400** is again identified as an outlier, mainly because

₃₆₇ of its highly negative expression level under treatment 3-1, while the near-zero expression

₃₆₈ level under 5-2 indicates its inactivity under this treatment. Gene **AT5G13220** has the

₃₆₉ highest expression level under 3-1 among all JA genes, and at the same time not as active

₃₇₀ under 5-2. This gene is interesting because it functions as a repressor of JA signaling,

₃₇₁ and its high expression level could be an implication for repression of JA singaling for

₃₇₂ *Arabidopsis* plants going through first drought stress (3-1).

| Gene name | Description |
|---|---|
| AT3G56400 | WRKY Transcription Factor, Group III; Repressor of JA-regulated genes; Activator of SA-dependent defense genes. |
| AT1G19180 | a.k.a. JAZ1 Nuclear-localized protein involved in JA signaling; JAZ1 transcript levels rise in response to a jasmonate stimulus. |
| AT5G13220 | a.k.a. JAS1, JAZ10 Repressor of JA signaling |
| AT2G44940 | Integrase-type DNA-binding superfamily protein |

**Table 10:** Outlier information for 3-1 and 5-2

₃₇₃ **Pattern differences**     From the scatterplot, the first thing we can observe is that quite

₃₇₄ a few genes from both pathways are up- or down-regulated under treatment 3-1, while

375 genes are expressed at an overall low level under 5-2. Nevertheless, a few genes from ET

376 group show overexpression pattern under 5-2. JA pathway genes populate quadrants I,

377 II and III, while ET pathway genes are mainly located in quadrants I, II and IV. Overall,

378 under 5-2, ET genes tends to be more active than JA genes.

379    In the previous two sections we singled out two of the top ranked scatterplots for

380 discussion. Interested readers are directed to the appendix for additional scatterplots for

381 top ranked feature subsets (Figures 6∼13).

## 4.2   GSI for PC transformed data

382 In Table 11, we report the GSI for PC transformed data, as well the maximum GSI

383 achieved by subsets of the original data. The proportion of total variation explained is

384 66.5% for first 2 PCs, 78.2% for first 3, and 85.6% for first 4. Through comparison,

385 we observe that using PCs as new features does not necessarily maximize the separation

386 between the distinct groups in the data, therefore confirming our statement in Section 3.4.

**Table 11:** Separation index for PC transformed data and maximum GSI for original data

| # of PCs | GSI | Features | max GSI achieved |
|:---:|:---:|:---:|:---:|
| 2 | 0.638 | 2 | 0.708 |
| 3 | 0.642 | 3 | 0.736 |
| 4 | 0.639 | 4 | 0.745 |

# 5.   CONCLUSION

387 **Conclusion**    In this article, we defined a numeric measure for the separation between

388 different groups of data, and used said measure to perform low dimensional feature sub-

389 set selection in order to find the most interesting angles to visualize high dimensional

390 data. By applying our method to a multi-experiment RNA-Seq data on *Arabidopsis* leave

391 tissues, we found that the top ranked feature subsets did demonstrate some interesting

392 differences in expression patterns between two biological pathways, which shows that our

393 method can be a potentially powerful tool in the exploratory analysis of such high dimen-

394  sional integrated/assembled data from various sources.

395

396  **Significance**      Firstly, our method yields well documented results. We enumerated the

397  group separation index for every low dimensional feature subset, and constructed the

398  scatterplots/scatterplot matrices for each case. If scientists know beforehand which fea-

399  tures are of interest, they will be able to directly access the corresponding entry in our

400  result. Secondly, through the application of mixture discriminant analysis, we were able

401  to summarize the expression pattern of groups of genes using a mixture of only a handful

402  of normal components. Furthermore, using the fitted MclustDA ellipses as visual aid, we

403  were able to clearly show the geometric structure of each group and make comparisons.

404  Finally, as seen in Figure 4, through visualization of the unfiltered data, we are able to

405  identify difference in expression patterns of non-DE genes between two biological path-

406  ways.

407

408  **Limitations** & **Future Work**      A limitation of our method is the difficulty of scaling

409  our feature selection method to data of higher dimensions. The first concern is the heavy

410  computational burden required for RCV. In our implementation, although we used parallel

411  computing to speed up computation as much as possible, the actual running times for 3-

412  and 4-dimensional subset are not quite satisfactory (Table 12), mainly due to the large

413  number of possible subsets. However, in practice, the 2-subset results are usually more

414  interpretable and visually appealing than its higher dimensional counterparts. Therefore,

415  we recommend doing only 2-dim feature subset selection for exploratory purposes.

**Table 12:** Average running time for 10-fold cross-validation for all feature subsets, aver-
aged over 50 runs with different random number seeds.

| Subset dim. | # of subsets | Avg. runtime (s) |
|:---:|:---:|:---|
| 2 | 253 | 65.04 |
| 3 | 1771 | 512.61 |
| 4 | 8855 | 2241.43 |

416      Another reason is that the scatterplot matrix becomes less informative when the num-

417  ber of displayed dimensions exceeds 4. Even in our study, scatterplot matrices of di-

mensions 3 and 4 cannot fully reflect geometric properties of the data. For 3-d and 4-d angles, the scatterplot matrix only shows projections to all axial dimensions, which doesn't precisely convey the amount of separation between two classes, computed using all 3 or 4 dimensions. It is difficult to visualize the geometric and topological differences by only looking at individual panels of scatterplots. To more effectively visualize higher dimensional feature subsets, we can consider using interactive visualization tools, such as GGobi Swayne et al. (2003) and R Shiny Chang et al. (2017). Both tools allow users to identify the same point in all panels of a scatterplot matrix, significantly increasing its visual expressiveness.

**Error rate definition** In our definition of TPM in (1), we made the assumption that the cost of misclassifying an object from either class is the same, i.e. $c(1|2) = c(2|1)$. We can adjust the cost values if we are more concerned about correctly classifying a certain class of observations.

**Evaluating reproducibility of experiments** Currently, a typical differential expression analysis is conducted in a gene-wise manner, i.e. genes are treated as observations and the treatment conditions as features. In our study, we took the same approach because our goal was to differentiate expression pattern between two groups of genes. However, with the increase in the availability of RNA-Seq data thanks to advances in information technology, we can also study the comparability and reproducibility of RNA-Seq experiments. In this sense, we will be exploring the relationship between treatment conditions or experiments, with genes acting as features/variables. Evaluation of experiment reproducibility is usually accomplished by performing the same experiment using the same setting, which is, unfortunately, not a common practice in RNA-Seq studies. In light of this, one of our long-term goal is the quantification of similarity between RNA-Seq experiments, which not only accounts for differences in experimental designs and parameter settings, but also utilize the information hidden in the expression of genes.

# 6.   ACKNOWLEDGMENTS

# REFERENCES

448  Carlos L. Ballaré. Jasmonate-induced defenses: a tale of intelligence, collaborators and

449      rascals. *Trends in Plant Science*, 16(5), 2010.

450  Tanya Z. Berardini, Leonore Reiser, Donghui Li, Yarik Mezheritsky, Robert Muller, Emily

451      Strait, and Eva Huala. The Arabidopsis Information Resource: Making and mining the

452      "gold standard" annotated reference plant genome. *genesis*, 2015. doi: 10.1002/dvg.

453      22877.

454  Winston Chang, Joe Cheng, JJ Allaire, Yihui Xie, and Jonathan McPherson. *shiny: Web

455      Application Framework for R*, 2017. URL `https://CRAN.R-project.org/package=`

456      `shiny`. R package version 1.0.3.

457  Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin. Maximum likelihood from

458      incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series

459      B (Methodological)*, 39(1):1–38, 1977.

460  Yong Ding, Michael Fromm, and Zoya Avramova. Multiple exposures to drought train

461      transcriptional responses in *Arabidopsis. Nature Communications*, 2012. doi: 10.1038/

462      ncomms1732.

463  Yong Ding, Ning Liu, Laetitia Virlouvet, Jean-Jack Riethoven, Michael Fromm, and Zoya

464      Avramova. Four distinct types of dehydration stress memory genes in *Arabidopsis

465      thaliana. BMC Plant Biology*, 13(229):1–38, 2013.

466  Niklas Elmqvist, Pierre Dragicevic, and Jean-Daniel Fekete. Rolling the Dice: Multidi-

467      mensional Visual Exploration using Scatterplot Matrix Navigation. *IEEE Transactions

468      on Visualization and Computer Graphics*, 14(6):1141–1148, 2008.

469  Chris Fraley and Adrian E. Raftery. Model-based clustering, discriminant analysis, and

470      density estimation. *Journal of the American Statistical Association*, 97:611–631, 2002.

471  Ji-Hyun Kim. Estimating classification error rate: Repeated cross-validation, repeated
472     hold-out and bootstrap. *Computational Statistics and Data Analysis*, 53(11):3735–3745,
473     2009.

474  Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model
475     selection. *Ijcai*, 14(2):1137–1145, 1995.

476  Zhefeng Lin, Silin Zhong, and Don Grierson. Recent advances in ethylene research. *Jour-*
477     *nal of Experimental Botany*, 60(12):3311–3336, 2009.

478  Guillaume Moissiard, Sylvain Bischof, Dylan Husmann, William A. Pastor, Christo-
479     pher J. Hale, Linda Yen, Hume Stroud, Ashot Papikian, Ajay A. Vashisht, James A.
480     Wohlschlegel, and Steven E. Jacobsen. Transcriptional gene silencing by *Arabidop-*
481     *sis* microrchidia homologues involves the formation of heteromers. *Proceedings of the*
482     *National Academy of Sciences*, 111(20):7474–7479, 2014.

483  David M. W. Powers and Adham Atyabi. The problem of cross-validation: Averaging
484     and bias, repetition and significance. *Engineering and Technology (S-CET)*, May 2012.

485  Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6, 1978.

486  Hume Stroud, Truman Do, Jiamu Du, Xuehua Zhong, Suhua Feng, Lianna Johnson,
487     Dinshao J. Patel, and Steven E. Jacobsen. Non-CG methylation patterns shape the
488     epigenetic landscape in *Arabidopsis*. *Nature Structural & Molecular Biology*, 21(1):
489     64–72, 2014.

490  Deborah F. Swayne, Duncan Temple Lang, Andreas Buja, and Dianne Cook. GGobi:
491     evolving from XGobi into an extensible framework for interactive data visualization.
492     *Computational Statistics & Data Analysis - Data Visualization*, 43:423–444, 2003.

493  Bin Zhuo, Sarah Emerson, Jeffrey H. Chang, and Yanming Di. Identifying stably expressed
494     genes from multiple RNA-Seq data sets. *PeerJ*, 4(e2791), 2016.

# A. CONSTRUCTION OF COVARIANCE ELLIPSES FOR NORMAL COMPONENTS

$_{495}$ In this section, we introduce how the covariance ellipses are constructed by MclustDA

$_{496}$ when a scatterplot or a scatterplot matrix is graphed.

$_{497}$

$_{498}$ For 2D data, suppose the mean and covariance estimates for component $k$ of class $j$

$_{499}$ are $\hat{\mu}_{jk}$ and $\hat{\Sigma}_{jk}$, respectively. Also suppose that $\hat{\Sigma}_{jk}$ has eigenvalues $\lambda_1 \geqslant \lambda_2$ and their

$_{500}$ corresponding eigenvectors $\boldsymbol{e}_1$ and $\boldsymbol{e}_2$. Then MclustDA computes the major and minor

$_{501}$ axes of the ellipse centered at $\hat{\mu}_{jk}$ the following way:

$$\text{major axis} = \hat{\mu}_{jk} \pm \sqrt{\lambda_1}\boldsymbol{e}_1, \quad \text{minor axis} = \hat{\mu}_{jk} \pm \sqrt{\lambda_2}\boldsymbol{e}_2,$$

$_{502}$ and the resulting ellipse has coverage probability of approximately 0.393.

$_{503}$

$_{504}$ In the case of higher dimensional data, MclustDA constructs the scatterplot and

$_{505}$ graphs the ellipses two dimensions at a time. Suppose $\hat{\mu}_{jk}$ and $\hat{\Sigma}_{jk}$ are defined the same

$_{506}$ way as above, and consider data dimensions $p$ and $q$ for visualization via scatterplot.

$_{507}$ Let $\Sigma^{(p,q)} = [\hat{\Sigma}_{jk}]_{(p,q)}$ be the covariance submatrix corresponding to the two dimensions,

$_{508}$ and $\mu^{(p,q)} = [\hat{\mu}_{jk}]_{(p,q)}$ be the corresponding mean vector. Now, suppose $\Sigma^{(p,q)}$ has eigen-

$_{509}$ value/eigenvector pairs $\{\lambda_1^{(p,q)}, \boldsymbol{e}_1^{(p,q)}\}$ and $\{\lambda_2^{(p,q)}, \boldsymbol{e}_2^{(p,q)}\}$ with $\lambda_1^{(p,q)} \geqslant \lambda_2^{(p,q)}$. Then the

$_{510}$ ellipse plotted by MclustDA has major and minor axes as follows:

$$\text{major axis} = \mu^{(p,q)} \pm \sqrt{\lambda_1^{(p,q)}}\boldsymbol{e}_1^{(p,q)}, \quad \text{minor axis} = \mu^{(p,q)} \pm \sqrt{\lambda_2^{(p,q)}}\boldsymbol{e}_2^{(p,q)},$$

$_{511}$ where the ellipse has the same coverage probability as the case above.

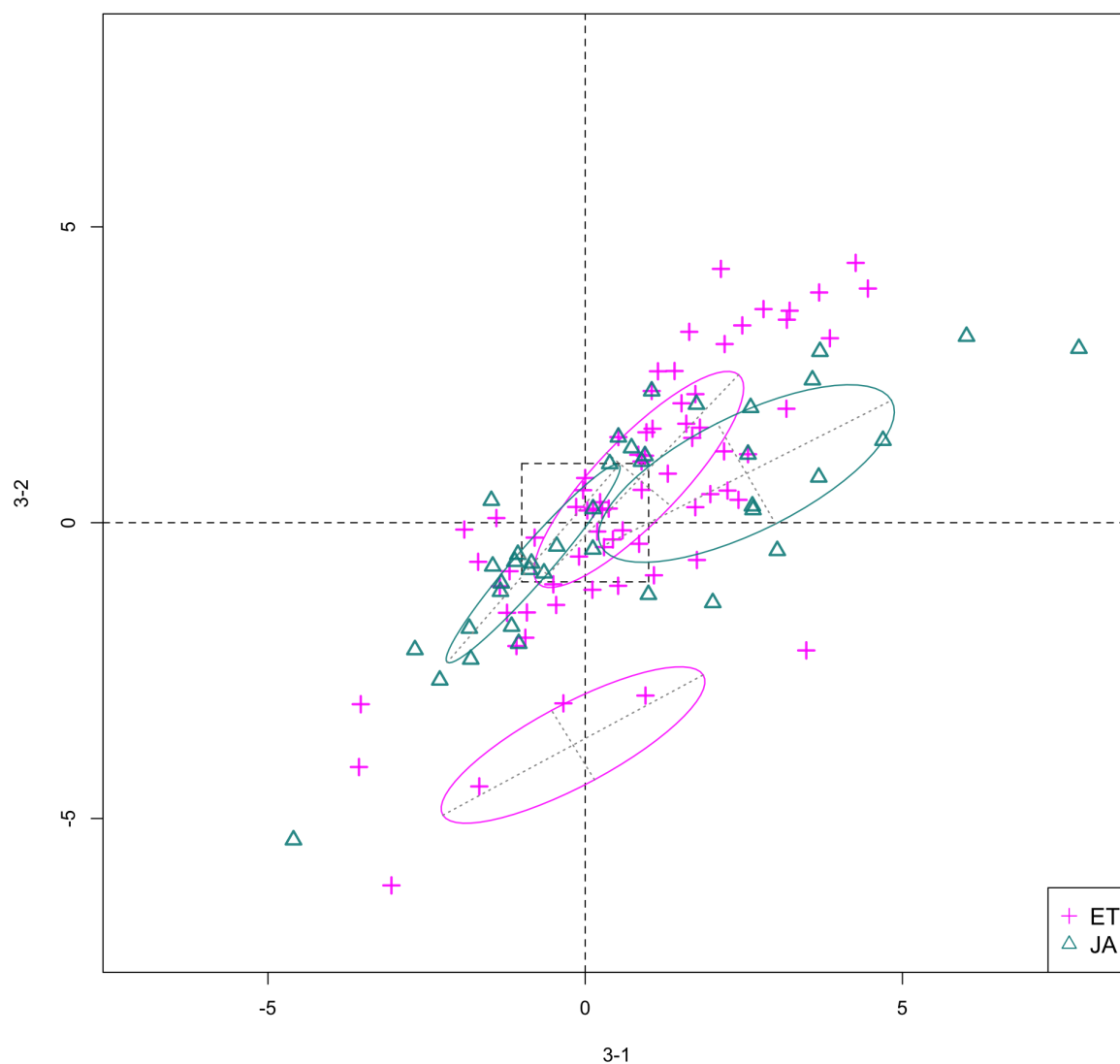# B. SCATTERPLOTS AND SCATTERPLOT MATRICES FOR SELECT TOP RANKED FEATURE COMBINATIONS



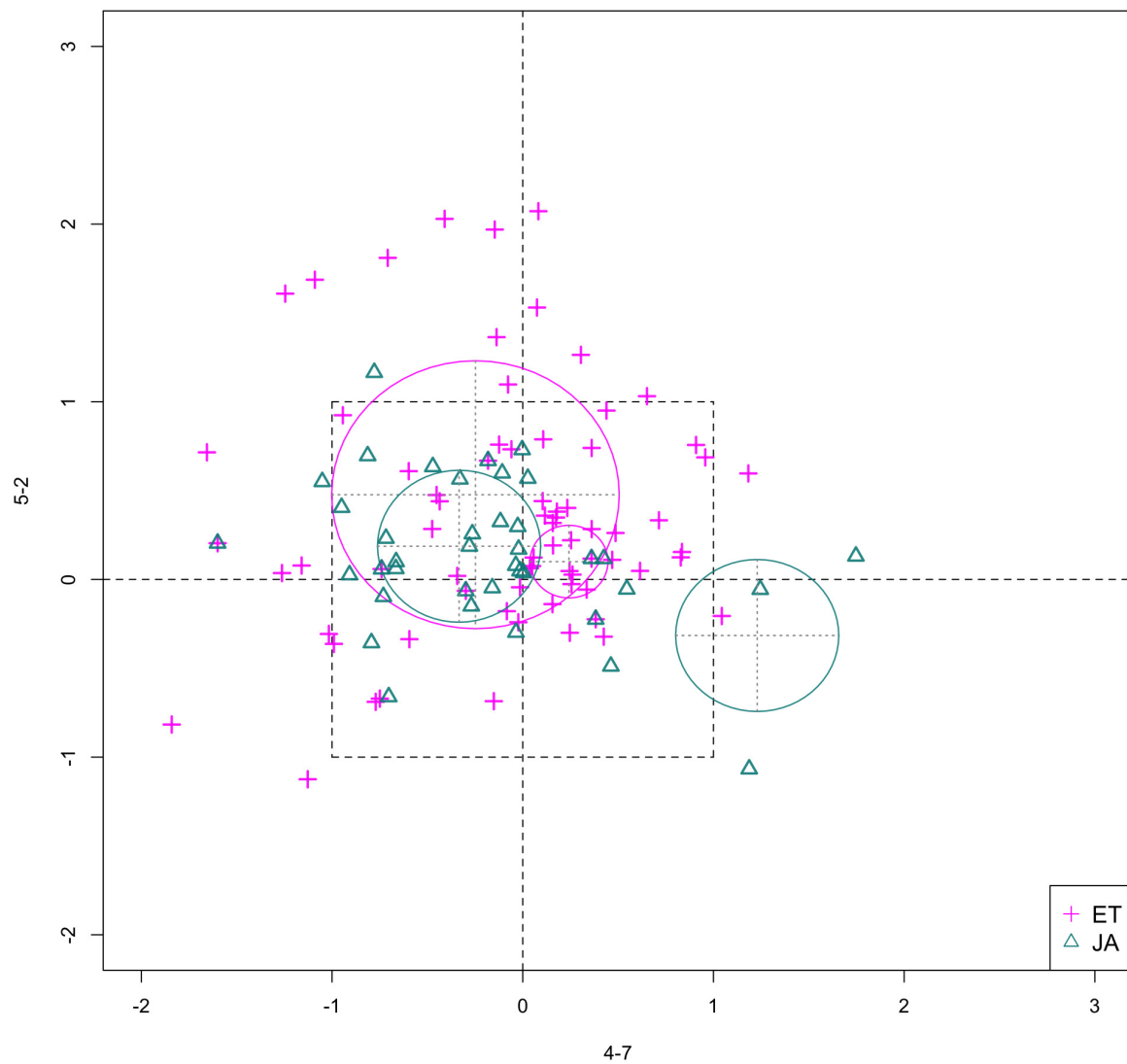**Figure 6:** Scatterplot for 3-1 and 3-2
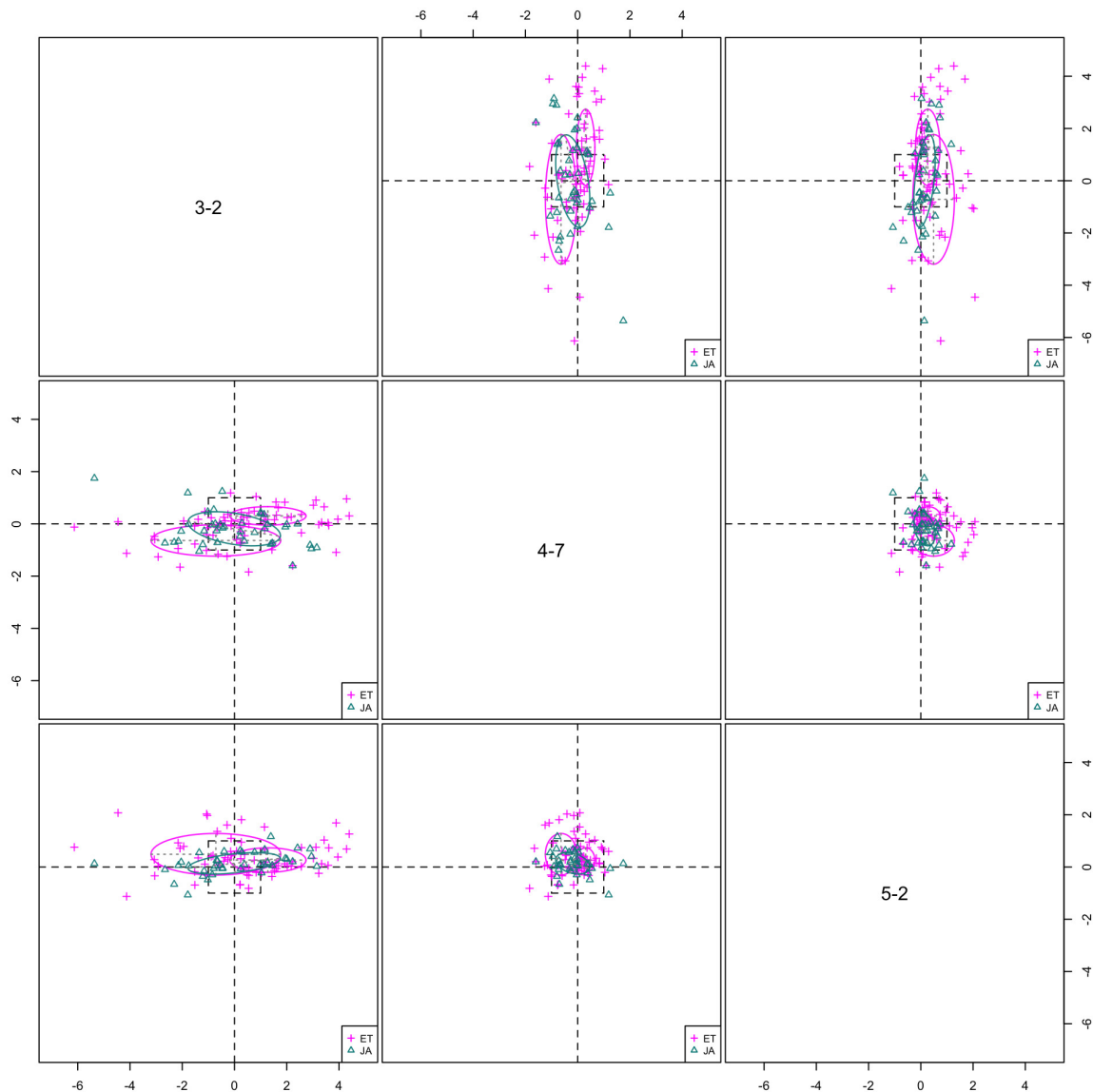
**Figure 7:** Scatterplot for 4-7 and 5-2

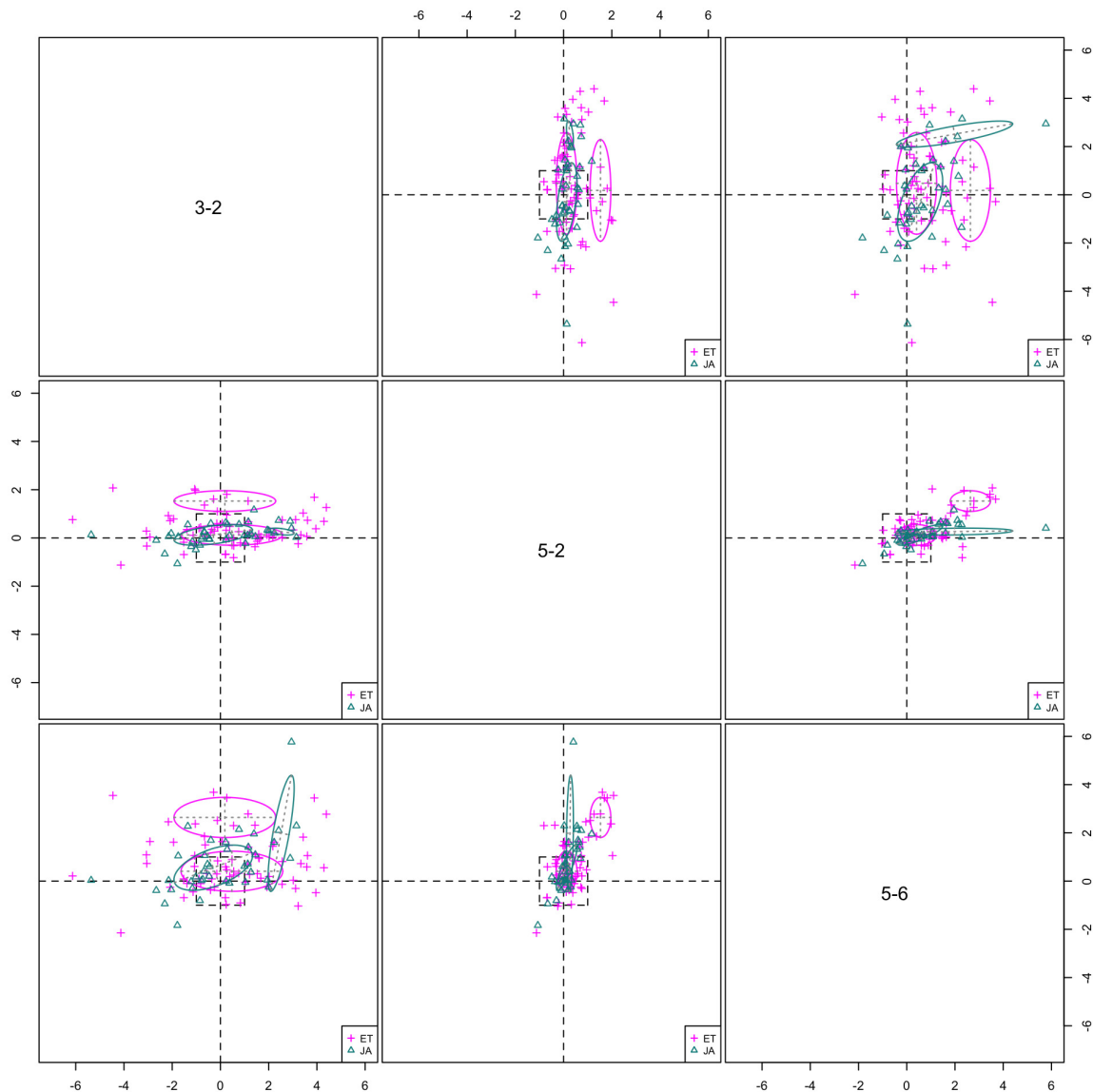**Figure 8:** Scatterplot matrix for [3-2, 4-7, 5-2]

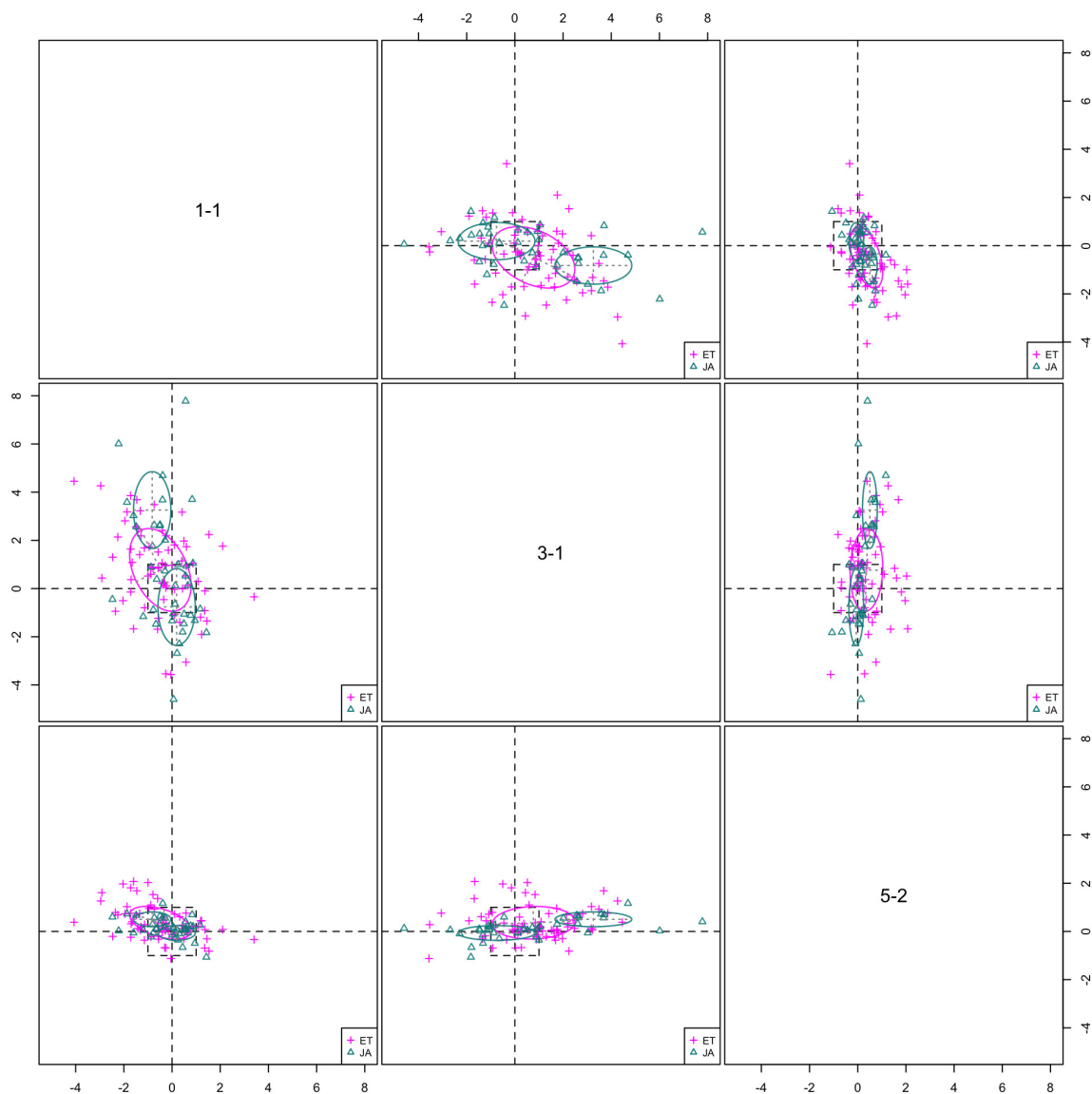**Figure 9:** Scatterplot matrix for [3-2, 5-2, 5-6]

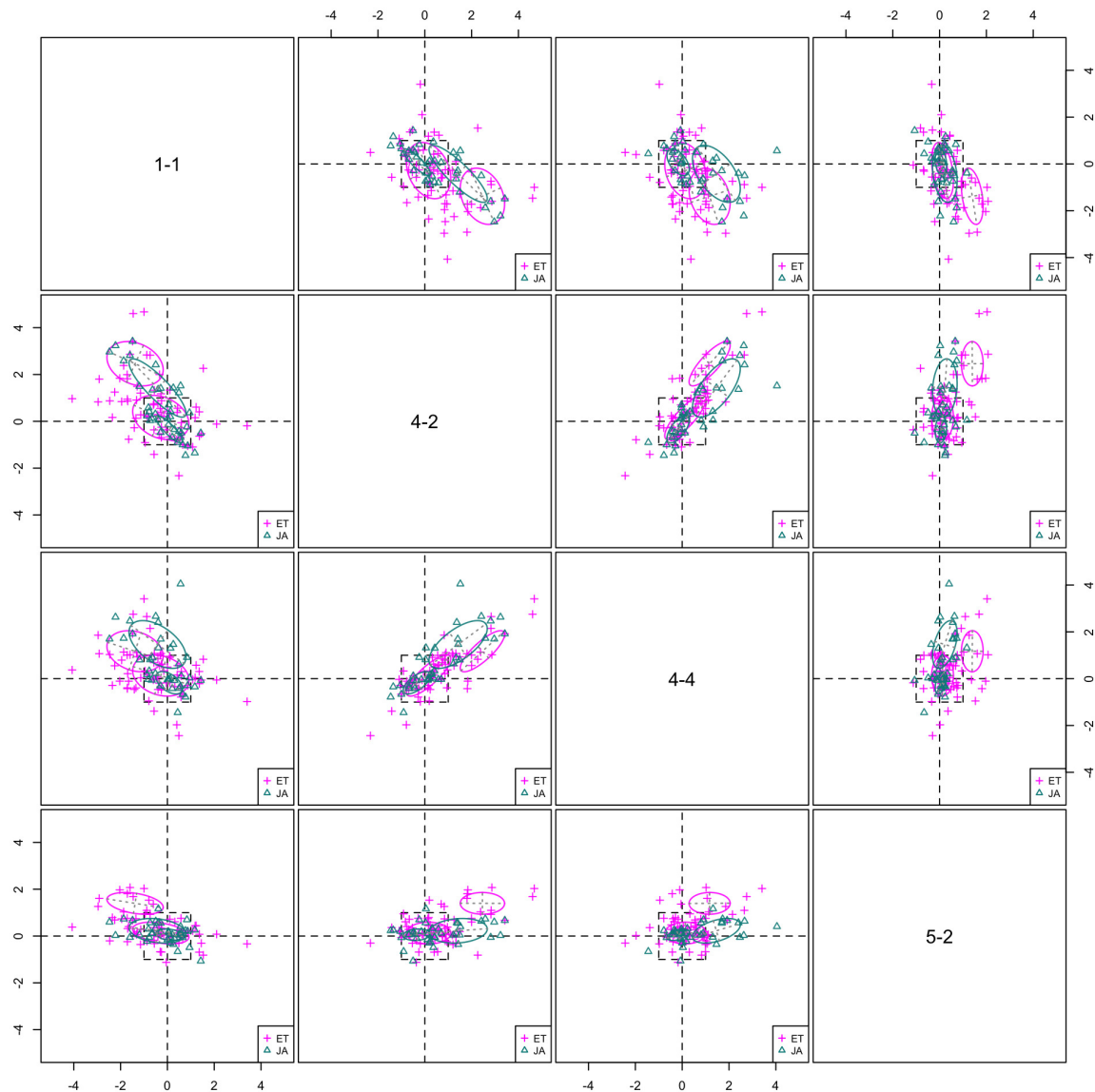**Figure 10:** Scatterplot matrix for [1-1, 3-1, 5-2]

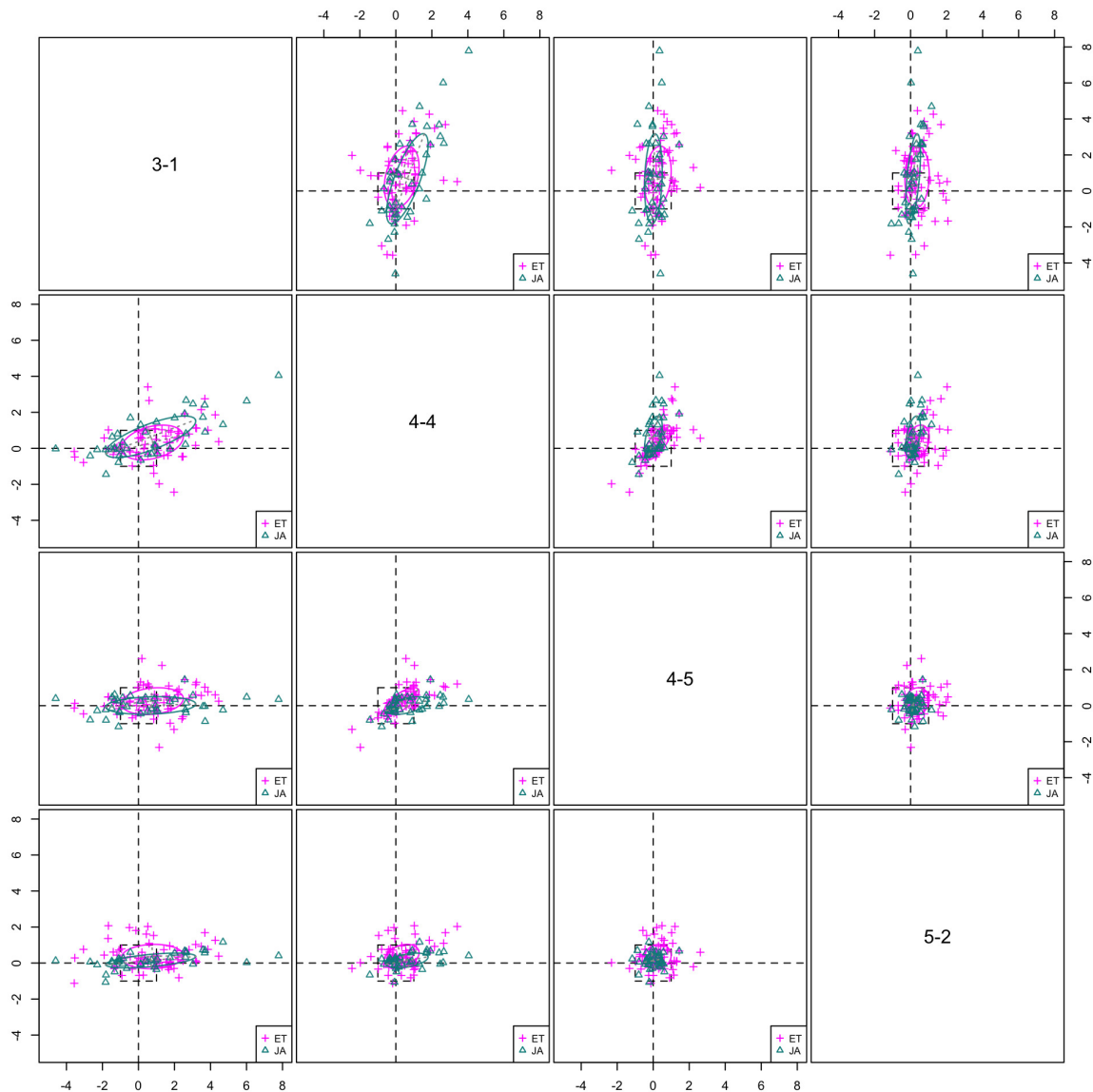**Figure 11:** Scatterplot matrix for [1-1, 4-2, 4-4, 5-2]
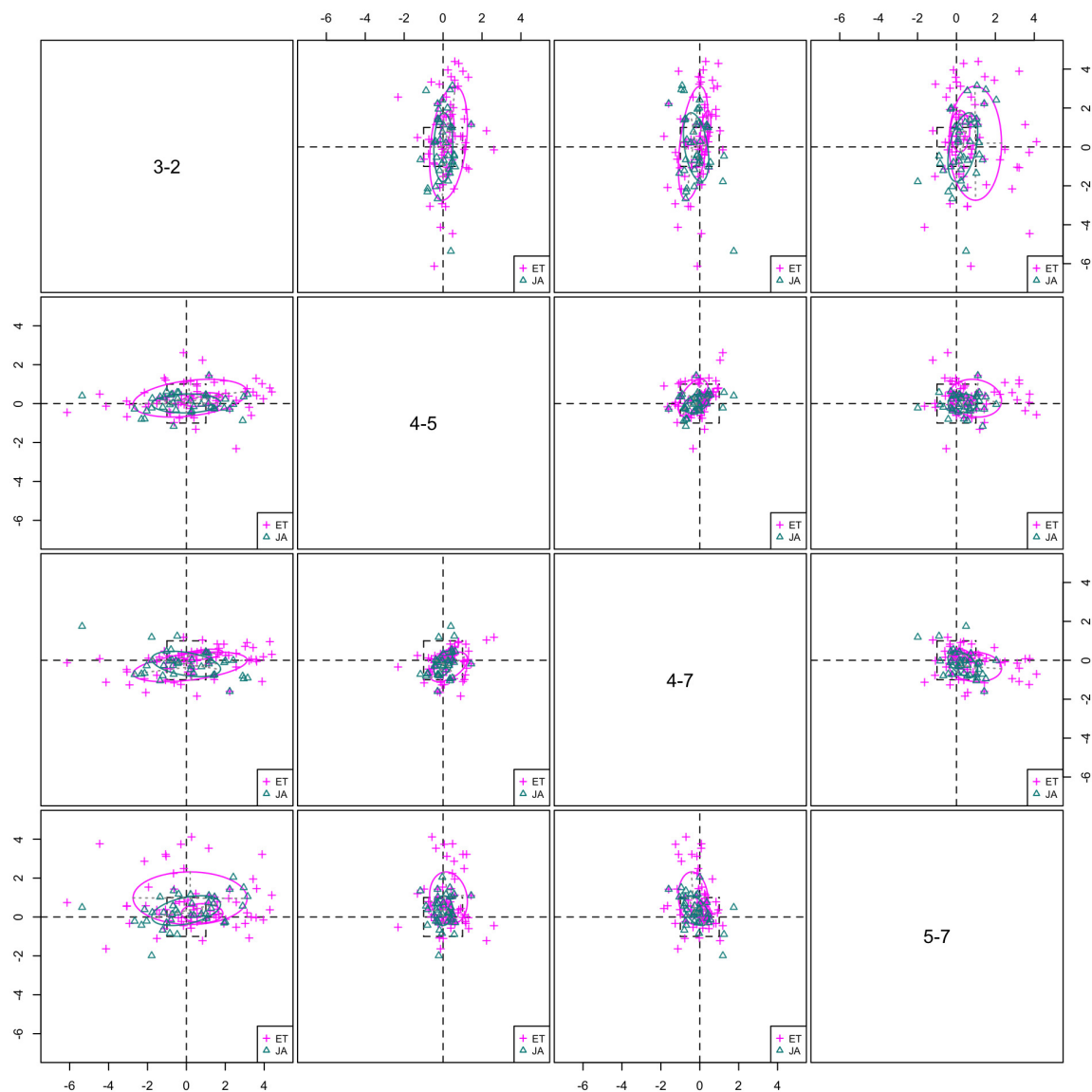
**Figure 12:** Scatterplot matrix for [3-1, 4-4, 4-5, 5-2]

**Figure 13:** Scatterplot matrix for [3-2, 4-5, 4-7, 5-7]