



**Universidad Politécnica de Madrid**

ESCUELA TÉCNICA SUPERIOR DE INGENIEROS INDUSTRIALES

MÁSTER EN AUTOMÁTICA Y ROBÓTICA

**APPLIED ARTIFICIAL INTELLIGENCE**

*Assignment 2.2: Principal Component Analysis*

Josep María BARBERÁ CIVERA (17048)

February 26, 2024

# Principal Component Analysis

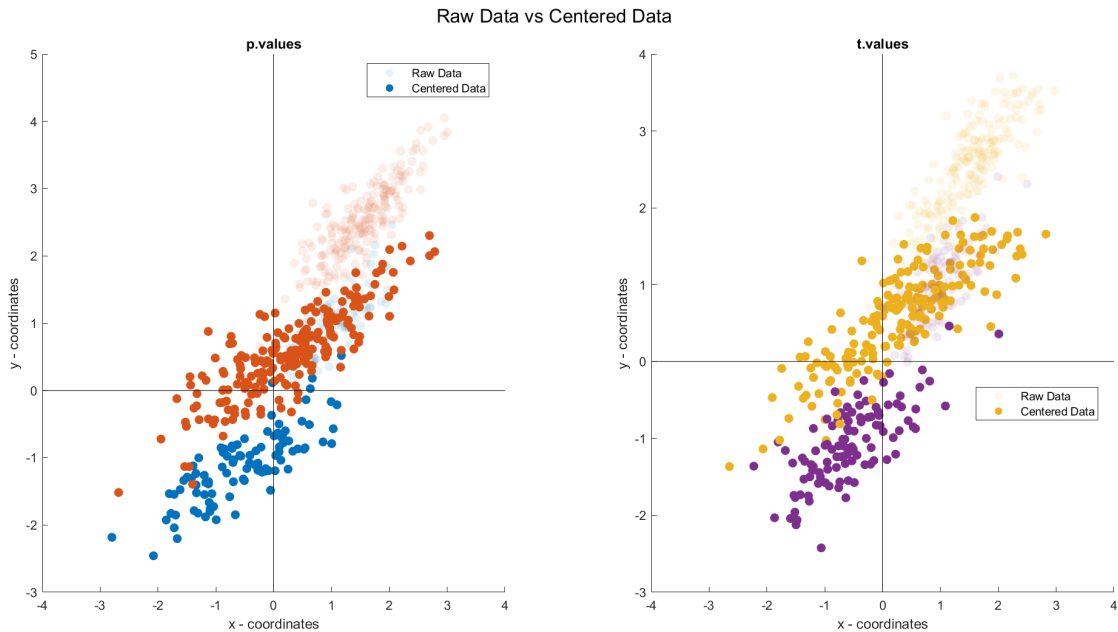
## 1.1 Methodology

Here we describe the steps followed for the PCA analysis.

1. First load the synthetically created “data\_D2\_C2” data set and then proceed to normalize the “p.value” and “t.value” datasets.
2. Next, the 1D and 2D coordinates of the normalized data projected onto the principal component subspace (PCA) are calculated.
3. After denormalizing the projected data, they should be plotted next to the original data set (raw data).
4. Finally, the accuracy of the reconstruction should be evaluated by calculating the expected mean squared error (MSE) of the normalized data and the actual MSE for both the normalized and original data, which should provide information on the effectiveness of the PCA analysis.

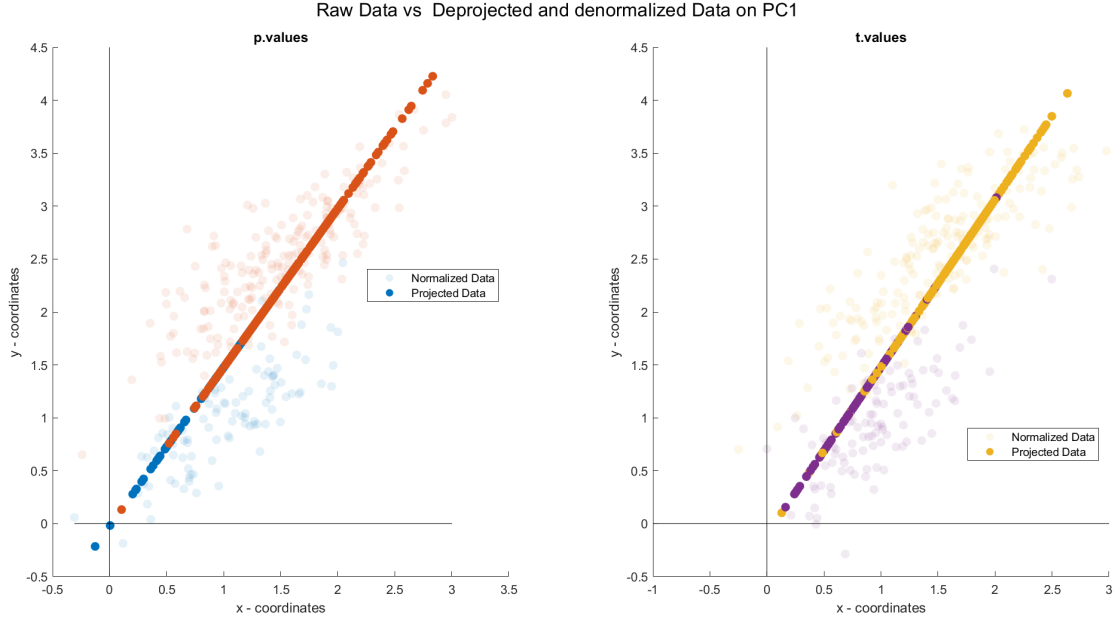
## 1.2 Results

First, the raw data were plotted against the normalized data. It should be noted at this point that normalization involves dividing by the standard deviation and subtracting the mean or simply subtracting the mean without the standard deviation division. As it is usually recommended to use the standard deviation, we have proceeded in this way. The following figure 1 shows the original data and the normalized data.



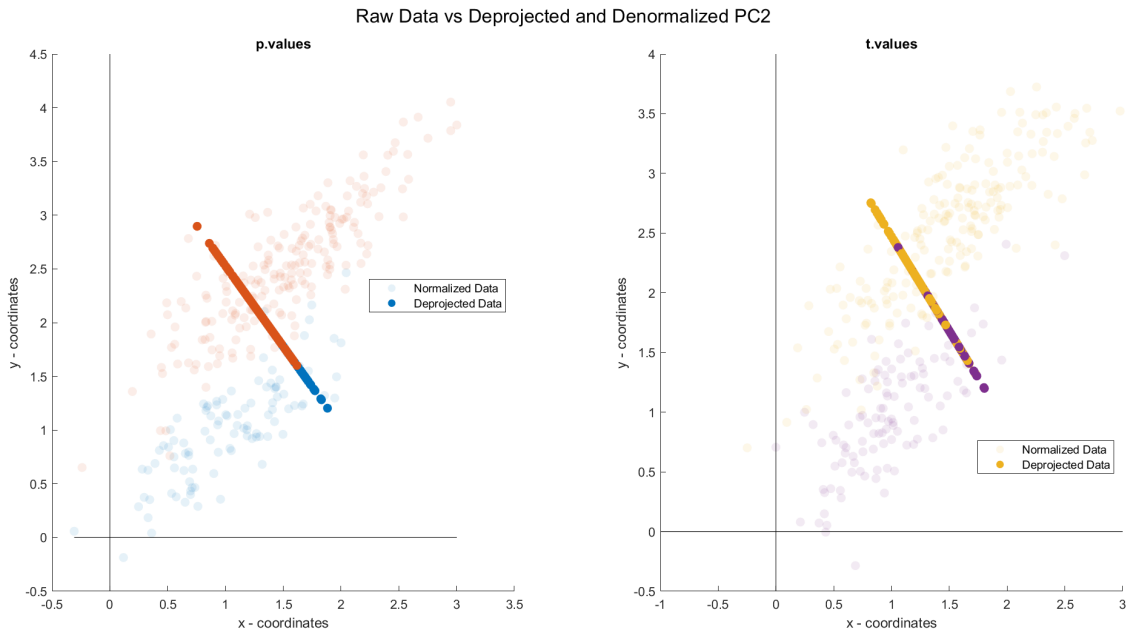
**Figure 1:** Normalized Data: row data in translucent color vs. normalized data. P-values and T-values datasets are plotted the colors refer to the associated labels.

After calculating the principal components, the normalized data are projected onto the component with the most information (the most significant) and, after deprojecting these points, they are plotted together with the original data (see Figure 2).



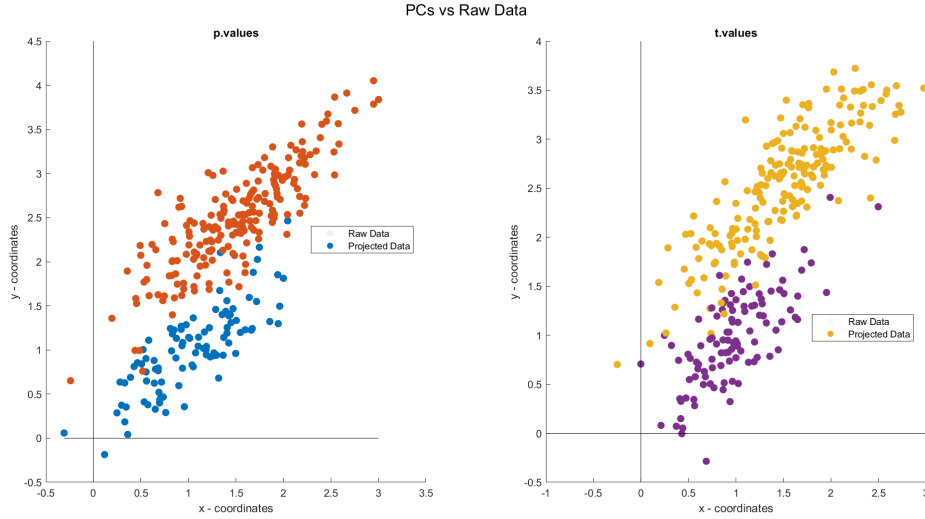
**Figure 2:** *Principal Component: the raw data is plotted again and now the most significant component has been used. Normalizing the data, projecting onto the eigenvector with the highest eigenvalue, deprojecting and denormalizing.*

A similar procedure was followed for the second component. The figure 3 shows the result.



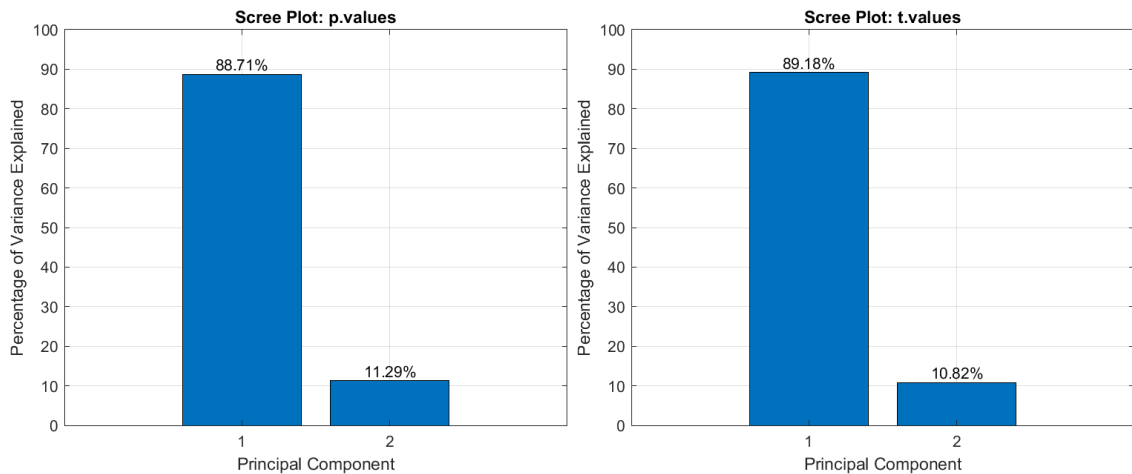
**Figure 3:** *The Second Component: as previous figure 2 the projected and deprojected data is plotted. Now the second eigenvector is used.*

Finally, it is projected on both components. The graph shows the result of this projection against the original data before deprojecting. In the figure 4, the deprojected data can be observed, which coincide with the original data (what we call **Raw Data**), since there are only two main components and there is no loss of information.



**Figure 4:** *PC 1 and PC 2: both components are analyzed, for this data it is meaningless because we only have two dimensions and therefore the reconstruction give us the same points than the original data (raw data).*

In addition, it has been plotted in a bar chart (Figure 5 with the value of the relative eigenvalues and we can see how the principal component is the most significant and allows us to represent  $\sim 90\%$  of the information with a single axis.



**Figure 5:** *Scree Plots for P-values and T-values data set. It is rapidly conclude how only one component is enough in order to meaningfully represent the information.*

### 1.3 Discussions and Results

We have been able to learn from several reference readings what PCA analysis consists of and its relevance in the case of image compression [1, 2].

In addition, it has been proven with synthetic data how the principal components of a point cloud can be found and how to project the information on these.

The calculation of errors in this analysis has been done in three way. First we have computed the **expected MSE for PC1** as

$$E[MSE] = \sum_{j=m+1}^D \lambda_j = \lambda_2$$

where  $D$  is the dimension of our data (in this case is 2), and  $m$  is the number of components we use (in this case 1). So for our case, the expected error corresponds directly with the second eigenvalue. In the case of **p-values** the  $E[MSE] = 0.22587$  and in the case of **t-values** it is  $E[MSE] = 0.21643$ .

The other two errors computed are the **actual MSE of the normalized data** which for the **p-values** corresponds to  $MSE_{norm} = 0.22512$  and for the **t-values** results  $MSE_{norm} = 0.21571$ . On the other hand, the **actual MSE of the not normalized data** is  $MSE = 0.12886$  for the **p-values** and  $MSE = 0.13140$  for the **t-values**.

## References

- [1] J. Shlens. *A Tutorial on Principal Component Analysis*. 2014. arXiv: [1404.1100 \[cs.LG\]](#).
- [2] L. I. Smith. *A tutorial on principal components analysis*. Tech. rep. Cornell University, USA, Feb. 2002. URL: [http://www.cs.otago.ac.nz/cosc453/student\\_tutorials/principal\\_components.pdf](http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf).

## 1.4 Relevant Code

```
1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 %                               Master in Robotics
3 %                               Applied Artificial Intelligence
4 %
5 % Assinment 2.2: Principal Component Analysis
6 % Student: Josep Barbera Civera
7 % ID: 17048
8 % Date: 20/02/2024
9 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
10
11 % 0. Load data_D2_C2 and rnrmalize the data (p.value and
    t.value)
12 % 1. Compute the 1D coordinates of the normalized test
    data
13 % projected on the PCA sub-space
14 % 2. Compute the 2D coordinates of the normalized test
    data
15 % projected on the PCA sub-space
16 % 3. De-normalize the projected data and plot it in the
    same
17 % graphic as original data
18 % 4. Compute the reconstruction, the expected MSE of
    normalized
19 % data, the actual MSE of the normalized data and the
    actual
20 % MSE of the not normalized data.
21
22 load data_D2_C2.mat
23
24 %% Accesing Data
25 pvalues = p.value;
26 plabels = p.class;
27 tvalues = t.value;
28 tlabels = t.class;
29
30 %% Normalizing data
31 disp("----- Normalizing Data
    -----");
32 disp("Normalizing P-values");
33 [Mp, Np] = size(pvalues);
34 mn_p = mean(pvalues')';
35 std_p = std(pvalues')';
36 for i = 1:Np
37     pn(:,i) = (pvalues(:,i) - mn_p)./std_p;
38 end
39 disp("Normalizing T-values");
```

```

40 [Mt, Nt] = size(tvalues);
41 mn_t = mean(tvalues');
42 std_t = std(tvalues');
43 for i = 1:Nt
44     tn(:,i) = (tvalues(:,i) - mn_t)./std_t;
45 end
46
47 %% Plotting Normalized Data vs Raw Data
48 two_plots('Raw Data vs Centered Data', 'p.values', 't.
    values', ...
49         'x - coordinates', 'y - coordinates', 'x -
        coordinates', ...
50         'y - coordinates', 'Raw Data', 'Centered Data
        ', ...
51         'Raw Data', 'Centered Data', pvalues, tvalues
        , pn, tn, ...
52         plabels, tlabels, 'centered_data_vs_raw_data.
        png');
53
54 %% Computing the Covariance Matrix
55 disp("----- Computing the Covariance Matrix
    -----");
56 cov_pn = 1 / (Np-1) * pn * pn';
57 cov_tn = 1 / (Nt-1) * tn * tn';
58
59 %% Computing the Eigenvalues and Eigenvectors
60 disp("----- Computing the Eigenvalues and Eigenvectors
    -----");
61 [PC_p, Vp] = eig(cov_pn);
62 [PC_t, Vt] = eig(cov_tn);
63
64 %% Sort the variances in decreasing order
65 disp("----- Sorting variances in decreasing order
    -----");
66 % Extract diagonal of matrix as vector
67 Vp = diag(Vp);
68 Vt = diag(Vt);
69 % Sort PC_p and convert Vp to a column vector with the
    eigenvalues
70 [~, p_rindices] = sort(-1*Vp);
71 Vp = Vp(p_rindices);
72 PC_p= PC_p(:, p_rindices);
73 % Sort PC_t and convert Vt to a column vector with the
    eigenvalues
74 [~, t_rindices] = sort(-1*Vt);
75 Vt = Vt(t_rindices);
76 PC_t= PC_t(:,t_rindices);
77

```

```

78 %% Projection 1D-PCA1
79 disp("----- Computing PC1
    -----");
80 % We project on the most significant PC
81 PC = PC_p(:,1);
82 signals_p = PC' * pn;
83 original_pvalues_1D_1 = (PC * signals_p);
84 for i = 1:Np
85     original_pvalues_1D_1_desnorm(:,i) =
        original_pvalues_1D_1(:,i).*std_p + mn_p;
86 end
87
88 PC = PC_t(:,1);
89 signals_t = PC' * tn;
90 original_tvalues_1D_1 = (PC * signals_t);
91 for i = 1:Nt
92     original_tvalues_1D_1_desnorm(:,i) =
        original_tvalues_1D_1(:,i).*std_t + mn_t;
93 end
94
95 %% Plotting Normalized Data vs PC1
96 % two_plots('Raw Data vs Deprojected and denormalized
    Data on PC1', ...
97 %         'p.values', 't.values', 'x - coordinates',
    ...
98 %         'y - coordinates', 'x - coordinates', 'y -
    coordinates', ...
99 %         'Normalized Data', 'Projected Data', '
    Normalized Data', ...
100 %         'Projected Data', pvalues, tvalues,
    original_pvalues_1D_1_desnorm, ...
101 %         original_tvalues_1D_1_desnorm, plabels,
    tlabels, ...
102 %         'Deprojected_and_denormalized_PC1_vs_Raw_Data
    .png');
103
104 %% Projection 1D-PCA2
105 disp("----- Computing PC2
    -----");
106 % We project on the **second** most significant PC
107 PC = PC_p(:,2);
108 signals_p = PC' * pn;
109 original_pvalues_1D_2 = (PC * signals_p);
110 for i = 1:Np
111     original_pvalues_1D_2_desnorm(:,i) =
        original_pvalues_1D_2(:,i).*std_p + mn_p;
112 end
113

```



```

114 PC = PC_t(:,2);
115 signals_t = PC' * tn;
116 original_tvalues_1D_2 = (PC * signals_t);
117 for i = 1:Nt
118     original_tvalues_1D_2_desnorm(:,i) =
119         original_tvalues_1D_2(:,i).*std_t + mn_t;
120 end
121 %% Plotting Normalized Data vs PC2
122 two_plots('Raw Data vs Deprojected and Denormalized PC2',
123     'p.values', ...
124     't.values', 'x - coordinates', 'y -
125         coordinates', ...
126     'x - coordinates', 'y - coordinates', '
127         Normalized Data', ...
128     'Deprojected Data', 'Normalized Data', '
129         Deprojected Data', ...
130     pvalues, tvalues,
131     original_pvalues_1D_2_desnorm, ...
132     original_tvalues_1D_2_desnorm, plabels,
133     tlabels, ...
134     'Deprojected_and_Denormalized_PC2_vs_Raw_Data
135         .png');
136
137 %% Plotting Scree Plot
138 disp("----- Plotting Scree Plots
139     -----");
140 % s = scree_plot(Vp, ' p.values');
141 % saveas(s, "scree_plot_PCA_p.png");
142 % s = scree_plot(Vt, ' t.values');
143 % saveas(s, "scree_plot_PCA_t.png");
144
145 %% Computing Errors
146 disp("----- Computing Error
147     -----");
148 % Expected Error
149 mse_expected = [Vp(2) Vt(2)];
150 sprintf('Expected MSE for PC1: p= %0.5f, t= %0.5f',
151     mse_expected(1), mse_expected(2))
152
153 % Actual Error with the normalized data
154 mse_actual_normalized_p = (pn - original_pvalues_1D_1)
155     .^2;
156 for i=1:Np
157     mse_sum_p_norm(i) = mse_actual_normalized_p(1,i) +
158         mse_actual_normalized_p(2,i);
159 end
160 mse_p_norm = mean(mse_sum_p_norm);

```

```

149
150 mse_actual_normalized_t = (tn - original_tvalues_1D_1)
    .^2;
151 for i=1:Nt
152     mse_sum_t_norm(i) = mse_actual_normalized_t(1,i) +
        mse_actual_normalized_t(2,i);
153 end
154 mse_t_norm = mean(mse_sum_t_norm);
155 sprintf('Actual MSE for Normalized PC1: p= %0.5f, t= %0.5
    f', mse_p_norm, mse_t_norm)
156
157 % Actual Error with the de-normalized data
158 mse_actual_de_normalized_p = (pvalues -
    original_pvalues_1D_1_desnorm).^2;
159 for i=1:Np
160     mse_sum_p_de_norm(i) = mse_actual_de_normalized_p(1,i
        ) + mse_actual_de_normalized_p(2,i);
161 end
162 mse_p_de_norm = mean(mse_sum_p_de_norm);
163
164 mse_actual_de_normalized_t = (tvalues -
    original_tvalues_1D_1_desnorm).^2;
165 for i=1:Np
166     mse_sum_t_de_norm(i) = mse_actual_de_normalized_t(1,i
        ) + mse_actual_de_normalized_t(2,i);
167 end
168 mse_t_de_norm = mean(mse_sum_t_de_norm);
169 sprintf('Actual MSE for Normalized PC1: p= %0.5f, t= %0.5
    f', mse_p_de_norm, mse_t_de_norm)
170
171
172 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
173 %% Plotting Functions
174 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
175 function two_plots(generic_title, title_subplot_1,
    title_subplot_2, ...
176                     x_1_label, y_1_label, x_2_label,
                        y_2_label, ...
177                     legend_1_light, legend_1,
                        legend_2_light, ...
178                     legend_2, light_data_1, light_data_2,
                        data_1, ...
179                     data_2, labels_1, labels_2,
                        saved_name)
180     disp("----- Plotting
        -----");
181     disp(generic_title);
182

```

```

183 figure;
184 % Generic Title
185 sgtitle(generic_title);
186 % First Subplot
187 subplot(1, 2, 1);
188 for i=1:length(labels_1)
189     if labels_1(i) == 1
190         scatter(light_data_1(1,i), light_data_1(2,i),
191                 50, 'o', ...
192                 'MarkerEdgeColor', 'none', '
193                 MarkerFaceColor', [0 0.4470
194                 0.7410], ...
195                 'MarkerFaceAlpha', 0.1); hold on;
196         scatter(data_1(1,i), data_1(2,i), 50, 'o', ...
197                 'MarkerEdgeColor', 'none', '
198                 MarkerFaceColor', [0 0.4470
199                 0.7410], ...
200                 'MarkerFaceAlpha', 1); hold on;
201     else
202         scatter(light_data_1(1,i), light_data_1(2,i),
203                 50, 'o', ...
204                 'MarkerEdgeColor', 'none', '
205                 MarkerFaceColor', [0.8500 0.3250
206                 0.0980], ...
207                 'MarkerFaceAlpha', 0.1); hold on;
208         scatter(data_1(1,i), data_1(2,i), 50, 'o', ...
209                 'MarkerEdgeColor', 'none', '
210                 MarkerFaceColor', [0.8500 0.3250
211                 0.0980], ...
212                 'MarkerFaceAlpha', 1); hold on;
213     end
214 end
215 % Plot vertical lines for x=0 and y=0
216 plot([0 0], ylim, 'k-');
217 plot(xlim, [0 0], 'k-');
218 % Subplot title
219 title(title_subplot_1);
220 % Axis labels
221 xlabel(x_1_label);
222 ylabel(y_1_label);
223 legend({legend_1_light, legend_1}, 'Location', 'best'
224        );
225 % Second Subplot
226 subplot(1, 2, 2);
227 for i=1:length(labels_2)
228     if labels_2(i) == 1

```

```

219         scatter(light_data_2(1,i), light_data_2(2,i),
220                 50, 'o', ...
221                 'MarkerEdgeColor', 'none', '
                MarkerFaceColor', [0.4940 0.1840
222                 0.5560], ...
223                 'MarkerFaceAlpha', 0.1); hold on;
224         scatter(data_2(1,i), data_2(2,i), 50, 'o', '
225                 MarkerEdgeColor', 'none', ...
226                 'MarkerFaceColor', [0.4940 0.1840
227                 0.5560], 'MarkerFaceAlpha', 1);
228         hold on;
229     else
230         scatter(light_data_2(1,i), light_data_2(2,i),
231                 50, 'o', ...
232                 'MarkerEdgeColor', 'none', '
233                 MarkerFaceColor', [0.9290 0.6940
234                 0.1250], ...
235                 'MarkerFaceAlpha', 0.1); hold on;
236         scatter(data_2(1,i), data_2(2,i), 50, 'o', '
237                 MarkerEdgeColor', 'none', ...
238                 'MarkerFaceColor', [0.9290 0.6940
239                 0.1250], 'MarkerFaceAlpha', 1);
240         hold on;
241     end
242 end
243 % Plot vertical lines for x=0 and y=0
244 plot([0 0], ylim, 'k-');
245 plot(xlim, [0 0], 'k-');
246 % Subplot title
247 title(title_subplot_2);
248 % Axis labels
249 xlabel(x_2_label);
250 ylabel(y_2_label);
251 legend({legend_2_light, legend_2}, 'Location', 'best'
252        );
253
254 pos = get(gcf, 'Position');
255 set(gcf, 'Position', pos+[-900 -300 900 300])
256 saveas(gca, saved_name);
257 end
258
259 function s = scree_plot(eigvalues, mytitle)
260     s = figure;
261     explained_variance = eigvalues / sum(eigvalues) *
262         100;
263     bar(explained_variance);
264     ylim([0, 100]);
265     title(strcat('Scree Plot: ', mytitle));

```

```

255 xlabel('Principal Component');
256 ylabel('Percentage of Variance Explained');
257 grid on;
258 % Add value on top of each bar
259 for i = 1:length(explained_variance)
260     text(i, explained_variance(i), sprintf('%.2f%%',
261         explained_variance(i)), ...
262         'HorizontalAlignment', 'center', '
263         VerticalAlignment', 'bottom');
262 end
263 tightfig;
264 end

```