# Solving Anscombe's Quartet using a Transfer Learning Approach

Kevin Bu[1] and Jose Clemente[1,2]

1. Department of Genetics and Data Science. Graduate School of Biomedical Sciences at the Icahn School of Medicine at Mount Sinai. 1 Gustave L. Levy Pl, New York, NY 10029.

**Emails:** kevin.bu@icahn.mssm.edu, jose.clemente@mssm.edu

ORCID: 0000-0002-9030-8062 (Kevin Bu)

2. To whom correspondence may be addressed. Email: jose.clemente@mssm.edu

**Classification**

PHYSICAL SCIENCES: Statistics

**Keywords**

Anscombe's Quartet, transfer learning, correlation analysis

**Author Contributions**

K.B. and J.C. designed and performed research, J.C. contributed data, K.B. analyzed data; K.B. and J.C. wrote the paper.

**This PDF file includes:**

> Abstract
> Main Text
> Figure Legends 1, 2, 3

**Abstract**

Analysis of high-dimensional datasets often involves usage of summary statistics, one of which is the correlation coefficient. These values are then used to inform downstream analysis, whether in feature selection or in subsequent construction of networks and heatmaps. Condensing pairwise scatterplots into these singular values however, often results in a loss of information. Originally proposed by F. J. Anscombe in his famous 'Anscombe's Quartet,' this phenomenon has been canonically used to demonstrate the importance of plotting and the limitations of summary statistics such as correlation or variance [F.J. Anscombe, (1973) *American Statistician*. 27 (1), 17-21]. While numerous methods exist for the generation of visually distinct datasets that share similar summary statistics, the converse has not been extensively studied. To address this gap, we propose ICLUST (Image CLUSTering), an image classifier tool that can visually distinguish correlations with similar summary statistics in simulations and identify meaningful clusters in real data. Such a tool can potentially benefit those performing exploratory analysis or feature selection in a complementary fashion by identifying relationships between variables that traditional summary metrics cannot provide.

**Significance Statement**

Distilling large-scale, multidimensional datasets via analysis of pairwise relationships often employs a single value to describe the relationship between variables. However, as demonstrated through simulations, such summarization fails to retain the nuances of the data. Characteristics such as the type of relationship (linear versus nonlinear, etc.) and the spread of the data are commonly lost when using correlations. Here we propose a transfer learning framework, borrowing from image clustering and classification software, to visually classify graphs. We apply our method towards separation of scatterplots with similar correlation statistics

28   but visually distinctive patterns in both simulations and real data, demonstrating its broad

29   applicability.

30

31

32   **Main Text**

33

34   **Introduction**

35

36          Exploratory analysis of large multidimensional datasets often relies on summary statistics

37   such as correlation coefficients for the construction of networks and heatmaps. However, the

38   usage of such summary statistics results in the loss of information encoded in the scatterplots of

39   pairwise relationships.  Anscombe's quartet has canonically been used to illustrate the

40   importance of graphing and the limitations of summary statistics such as correlation or variance -

41   Anscombe himself stated, "make both calculations and graphs. Both sorts of output should be

42   studied; each will contribute to understanding" [1]. This is especially critical in biological fields as

43   Pearson and Spearman correlation are the default analytical tools when performing exploratory

44   analysis in the gene expression and microbiome domains respectively [2-4].

45          Several methods have been developed to generate these kinds of datasets, analogous to

46   Anscombe's Quartet. The Datasaurus is one such dataset, generated using either a genetic

47   algorithm or a simulated annealing method [5, 6]. However, there is a lack of tools that can

48   separate these plots once they have been generated. Even the more modern exploratory data

49   analysis tools still collapse pairwise relationships into summary statistics such as the s-Corrplot

50   package or the MIC, which like Spearman, only quantifies strength of relationship without

51   specifying the nature of that association [7, 8]

52          Here we propose ICLUST, a tool that employs transfer learning based on the pre-trained

53   VGG16 convolutional neural network. Although the model had been trained to distinguish images

54   of cats and dogs, by extracting the last layer of the network (4096 features), we can use the pre-

55    trained weights to distinguish images of plotted pairwise correlations in an automated fashion,

56    thus seeking to find 'visual' similarities in a way that would be impossible manually. We apply this

57    tool to the separation of pairwise correlations from simulations and real data with the hypothesis

58    that ICLUST can visually distinguish correlations with similar summary statistics (with

59    performance inversely proportional to noise) and identify clusters in real data, some of which

60    would have been masked by using correlation coefficients alone as a clustering criterion.

61

62

63    **Results**

64

65         We first applied ICLUST to Anscombe's quartet, taking the original data and adding to

66    each point a specified amount of noise according to a bivariate normal distribution. Five plots

67    were created for each class at each level of noise; the resulting set of images was then passed

68    through ICLUST and the PCA plots are shown in **Fig. 1a.** A v-measure score (VMS) for each

69    level of noise was computed to quantitatively assess the quality of clustering in accordance with

70    the true labels. VMS as a function of noise (orange) is shown in **Fig. 1b** with error bars reflecting

71    the standard deviation over one hundred such trials. The baseline for comparison (shown in blue)

72    is the VMS obtained using clustering based on distances of the Pearson correlation summary

73    statistic alone. Consistent with our hypothesis, increasing the level of noise reduces the accuracy

74    of clustering as plot classes begin to overlap upon visual examination (**Fig. 1c**).

75         Given the relative efficacy of ICLUST on distinguishing clusters in canonical simulated

76    data, we tested whether or not ICLUST could identify distinct clusters in real data. We applied

77    ICLUST to data obtained from the WHO on a variety of health statistics for each country by

78    computing pairwise correlations between all variables and arbitrarily choosing a window of

79    Pearson correlation values in which to examine scatterplots. By doing so, we emulate the

80    simulation approach described earlier, generating a dataset with similar values but potentially

81    differing shapes and relationships. Here, we arbitrarily choose a window of correlation magnitude

4

82    and select all correlations with Pearson's *r* with a magnitude between 0.8975 to 0.9025.

83    Hierarchical clustering based on Euclidean distance between correlation strength yields the

84    dendrogram in **Figure 2a**, while clustering using ICLUST 4096-component feature vectors yields

85    the structure in **Figure 2b.** Clustering assignment was determined by the best silhouette score,

86    which corresponded to k = 2 clusters. The average image of the scatterplots in clusters 1 (red)

87    and 2 (teal) are shown for correlation strength-based clustering and ICLUST in **Figure 2c** and

88    **Figure 2d**, respectively. The PCA plot obtained based on Euclidean distance of the image

89    fingerprints is shown in **Figure 2e.** Notably, variables that fall in cluster 1 tend to be normalized

90    rates (e.g. immunization per 1000), while variables that fall in cluster 2 tend to be less uniformly

91    distributed because of the presence of outliers. An example of this is population of a country, as

92    countries such as India and China that are expected to be outliers skew the distribution.

93          We then applied ICLUST to an airline delays dataset, containing various metrics for

94    flights (such as time spent taxiing). In this dataset, we can not only distinguish visual differences

95    in shape (across a variety of correlation strengths, from *r* = 0 to *r* = 1) but also observe

96    correlations that share similar correlation coefficients but distinct visual structure. When

97    performing clustering analysis, the algorithm chooses k = 2 as the best silhouette score both

98    when using correlation strength (**Fig. 3a**) or image fingerprints (**Fig. 3b**). The average image

99    corresponding to these clusters for correlation strength and image fingerprints are shown in

100   **Figure 23** and **Figure 3d** respectively. In **Figure 3c**, cluster 1 corresponds to the teal cluster in

101   **Figure 3a** while and cluster 2 corresponds to the red cluster. In **Figure 3d**, Cluster 1 is the teal

102   portion of the dendrogram in **Figure 3b.** The PCA plot of these clusters based on the neural

103   network fingerprints is shown in **Figure 3e**. , correlations with similar strength can appear

104   drastically different, while correlations with different strength can appear more similar (**Fig. 3f-g**).

105   Thus with both real examples and simulation, we demonstrate how Anscombe's observation is

106   indeed applicable to real world settings and that ICLUST can both separate visually distinct

107   graphs that share summary statistics and cluster similar graphs with different correlation

108   coefficients.

109

110

**Discussion**

112

113     Given the prevalent usage of summary statistics in constructing models, networks, and

114     other meaningful representations of data, we propose a transfer learning based image-clustering

115     approach to the separation of scatterplots. Through simulations of Anscombe's Quartet as well as

116     representative real datasets (WHO, airline), we demonstrate the efficacy of ICLUST in identifying

117     clusters of distinct patterns where summary statistics would otherwise fail to do so. Going

118     forward, ICLUST can aid in exploratory data analysis in a complementary fashion to traditional

119     methods, in a way consistent with Anscombe's axiom of combining both graphs and calculations

120     to arrive at the most accurate representation of data.

121

122

**Methods**

124

Plotting Simulated Data

126     Bivariate independent uniform displacement was added to Anscombe's Quartet in the following

127     manner. Let (xi, yi) be a datapoint from the dataset. We define $\sigma_{max} \in \{0.1, 0.25, 0.5, 0.75, 1\}$;

128     values were arbitrarily chosen to yield a representative range of noises. A new simulated dataset

129     for each $\sigma_{max}$ was generated by computing $[x_i + e_1, y_i + e_2]$ where $e_1, e_2 \sim Unif(0, \sigma_{max})$ for each

130     dataset. Python's Matplotlib and Seaborn libraries with were used to construct plots. Opacity of

131     points was set to alpha = 0.1 such that overlapping points were treated differently when plotted.

132     The default sns.lmplot function was used with palette='set1' and default marker size=36,

133     shape='o'. The origin of each plot was fixed at the center of the coordinate axes (which is

134     hidden). The scales of the plots are allowed to vary per default plotting parameters and the

135     method is thus scale invariant. For each set of parameters, 100 simulations were generated. Note

136 that images shown in Fig. 1. are enlarged and include the axes for better visibility; however the

137 clustering analysis was performed on the raw images.

138

139

140 Evaluating Performance on Simulation Data

141 An unweighted v-measure score (VMS) was used to assess the performance of ICLUST on the

142 labeled simulated data, as defined by:

143

$$VMS = \frac{2(Homogeneity * Completeness)}{Homogeneity + Completeness}$$

144 Where homogeneity is defined as:

$$Homogeneity = 1 - \frac{H(C|K)}{H(C)}$$

145 where

$$H(C|K) = -\sum_{c,k} \frac{n_{ck}}{N} \log\left(\frac{n_{ck}}{n_k}\right)$$

146

$$H(C) = -\sum_{c} \frac{\sum_k n_{ck}}{C} \log\left(\frac{\sum_k n_{ck}}{C}\right)$$

147 And completeness is defined as:

$$Completeness = 1 - \frac{H(K|C)}{H(K)}$$

148 Where

$$H(K|C) = -\sum_{c,k} \frac{n_{ck}}{N} \log\left(\frac{n_{ck}}{n_c}\right)$$

$$H(K) = -\sum_{k} \frac{\sum_c n_{ck}}{C} \log\left(\frac{\sum_c n_{ck}}{C}\right)$$

149 Where N is the total number of points, C is the total number of labels, and $n_c, n_k, n_{ck}$ represent the

150 number of elements with true label C, in cluster K, and in cluster K with label C, respectively.

151

152    This is a generalization of the weighted VMS, given by:

$$V_\beta = \frac{(1 + \beta)\,hc}{\beta h + c}$$

153    Where $\beta$ scales the VMS by a weighting towards homogeneity; here we set $\beta = 1$.

154

155

156    <u>Plotting WHO and Airline Delay Data</u>

157    For real-world datasets, Python's Matplotlib and Seaborn libraries were used to construct

158    scatterplots. Opacity of points was set to alpha = 0.1 such that overlapping points were treated

159    differently when plotted. The default sns.lmplot function was used with palette='set1' and default

160    marker size=36, shape='o'. The upper and lower bounds for the x and y axes are dynamic and

161    vary on a scatterplot by scatterplot basis, thus using the default parameters for determination of

162    scaling and display. All correlations with a Pearson's r between 0.8975 and 0.9025 in magnitude

163    were plotted, yielding n = 51 scatterplots. The window was chosen based on a range likely to

164    contain various shapes as described by Anscombe's Quartet. Two plots were removed from the

165    WHO dataset as outliers (identified via initial PCA), resulting in 49 plots. The outliers were

166    removed to best demonstrate the two distinct clusters; visually, the outliers appeared distinct from

167    the other plots consistent with ICLUST's ability to distinguish visual differences between

168    scatterplots. For the airline data, all scatterplots were plotted in the dataset (n = 80 scatterplots).

169    WHO data and airline delay data were obtained from sources [8] and [9] respectively.

170

171

172    <u>Image Classification, Transfer Learning and Image Clustering</u>

173    Image classification in ICLUST uses the VGG16 model, a convolutional neural network trained on

174    the ImageNet dataset [10]. Briefly, input images (.PDFs, .PNGs, etc.) are scaled via the keras PIL

175    image library which converts them into VGG16 inputs i.e. RGB (3-channel) images, each of

176    dimensions 224 x 224 pixels. With each successive layer of the network, these pixels are

8

177    converted into features using pre-trained functions. Instead of using the original output layer

178    however, in a transfer learning setting, we adopt the penultimate layer (4096 features) as the

179    feature map for our problem and use these as fingerprints for each image. Unsupervised

180    clustering is performed using UPGMA after calculating the Euclidean distance between feature

181    vectors corresponding to each image. Silhouette score is computed for each possible number of

182    clusters, iterating from 2 through max_clust (default=10). Code for image processing and transfer

183    learning were obtained from an open-source GitHub repository (see acknowledgements). The

184    software was adapted from an earlier version and streamlined for use with the addition of new

185    functionality such as concatenation of images, creation of dendrograms, generation of average

186    images, and clustering based off of silhouette score. The average image for a given cluster is

187    obtained by averaging the pixel intensities across the entire image for all members in the cluster.

188    If the true class labels are given, clustering accuracy is assessed using VMS; otherwise,

189    unsupervised clustering is performed in which the program iterates through cluster numbers

190    (default range is 2-20) with the cluster number chosen based on the k that yields the highest

191    silhouette score. All raw data and code used to generate analysis and figures are located at

192    https://github.com/kbpi314/ICLUST.

193

194

195    **Acknowledgments**

196

199

200

201    **References**

202    1.    Anscombe, F.J., *Graphs in Statistical Analysis.* The American Statistician, 1973. **27**(1): p.

203        17-21.

204   2.      Langfelder, P. and S. Horvath, *WGCNA: an R package for weighted correlation network*

205           *analysis.* BMC Bioinformatics, 2008. **9**: p. 559.

206   3.      Song, W.M. and B. Zhang, *Multiscale Embedded Gene Co-expression Network Analysis.*

207           PLoS Comput Biol, 2015. **11**(11): p. e1004574.

208   4.      Caporaso, J.G., et al., *QIIME allows analysis of high-throughput community sequencing*

209           *data.* Nat Methods, 2010. **7**(5): p. 335-6.

210   5.      Matejka, J. and G. Fitzmaurice, *Same Stats, Different Graphs*, in *Proceedings of the*

211           *2017 CHI Conference on Human Factors in Computing Systems - CHI '17.* 2017. p.

212           1290-1294.

213   6.      Chatterjee, S. and A. Firat, *Generating Data with Identical Statistics but Dissimilar*

214           *Graphics.* The American Statistician, 2007. **61**(3): p. 248-254.

215   7.      McKenna, S., et al., *s-CorrPlot: an interactive scatterplot for exploring correlation.* Journal

216           of Computational and Graphical Statistics, 2016. **25**(2): p. 445-463.

217   8.      Reshef, D.N., et al., *Detecting novel associations in large data sets.* Science, 2011.

218           **334**(6062): p. 1518-24.

219   9.      Wickham, H. *ASA Sections on Statistical Computing.* 2018  June 23, 2020]; Available

220           from: http://stat-computing.org/dataexpo/2009/the-data.html.

221   10.     Simonyan, K. and A. Zisserman, *Very Deep Convolutional Networks for Large-Scale*

222           *Image Recognition*  International Conference on Learning Representations, 2015.

223

224

225   **Figure Legends**

226

227   **Fig. 1.** ICLUST can resolve Anscombe's Quartet. (*A*) Principal coordinate analysis (PCA) plots

228   (with PCA computed on the transfer learning features) at varying noise levels where points

229   represent images of scatterplots derived from Anscombe's quartet with the addition of noise. (*B*)

230   V-measure score (VMS) as a function of noise level for the clustering structure in Anscombe's
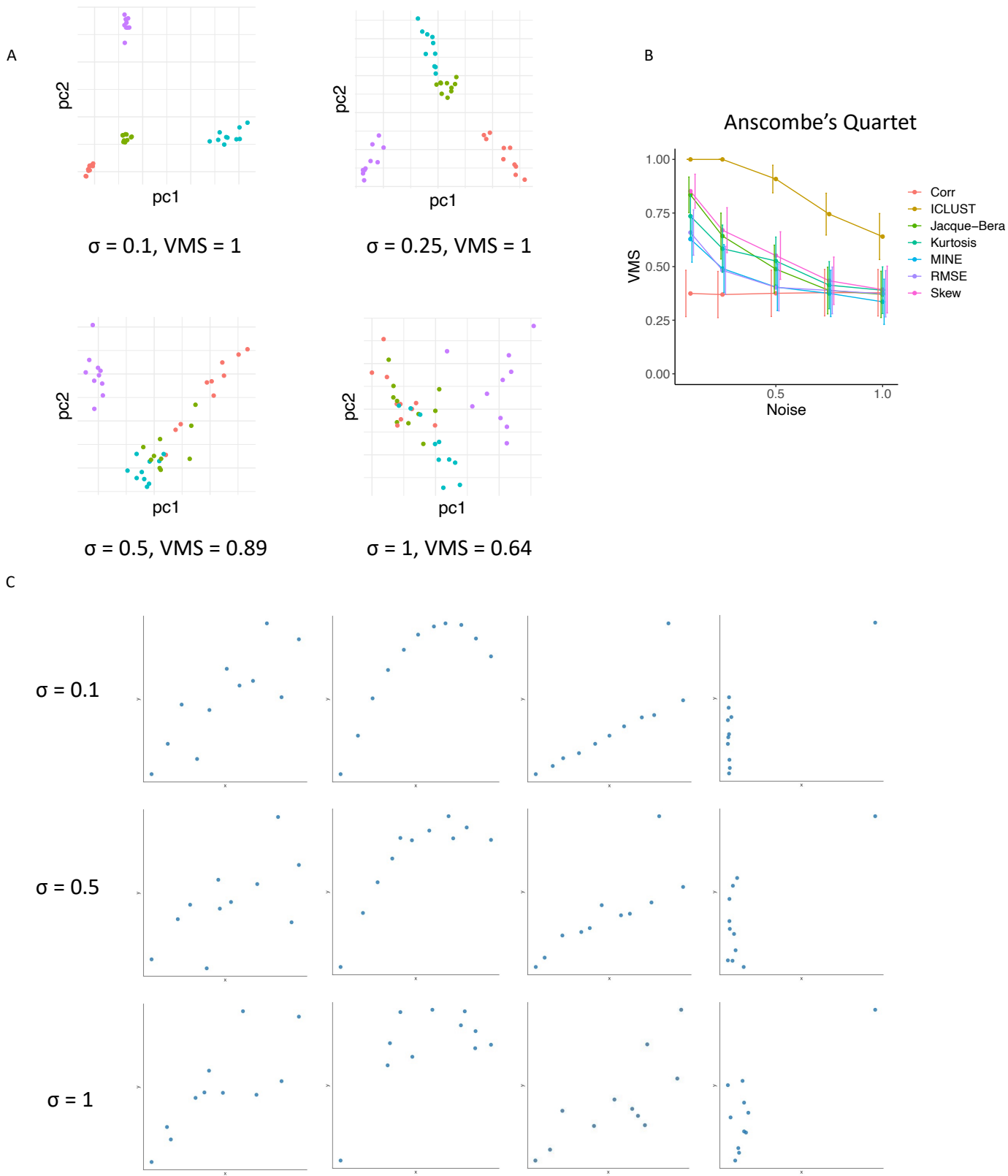
10

231    Quartet (obtained via cutting the hierarchical clustering tree at k = 4). (*C*) Examples of how the

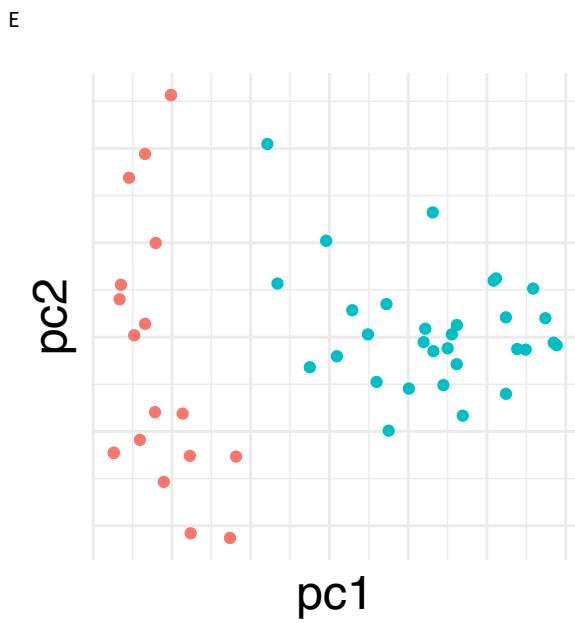232    plots become distorted as noise levels increase.
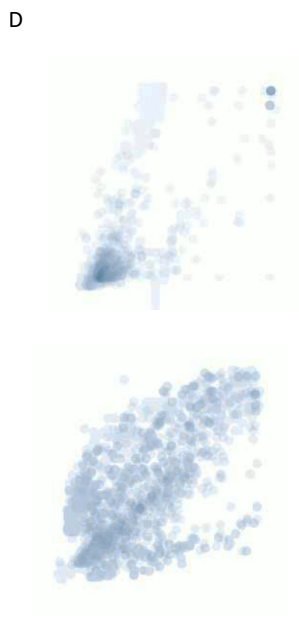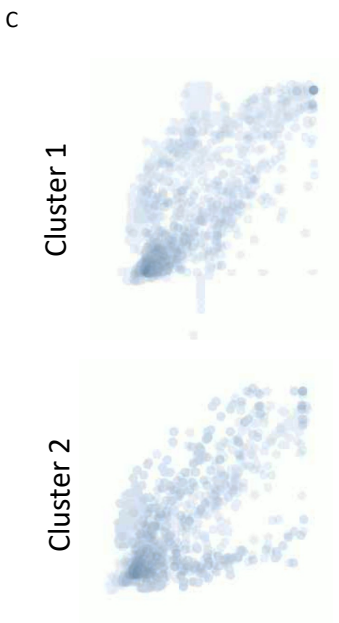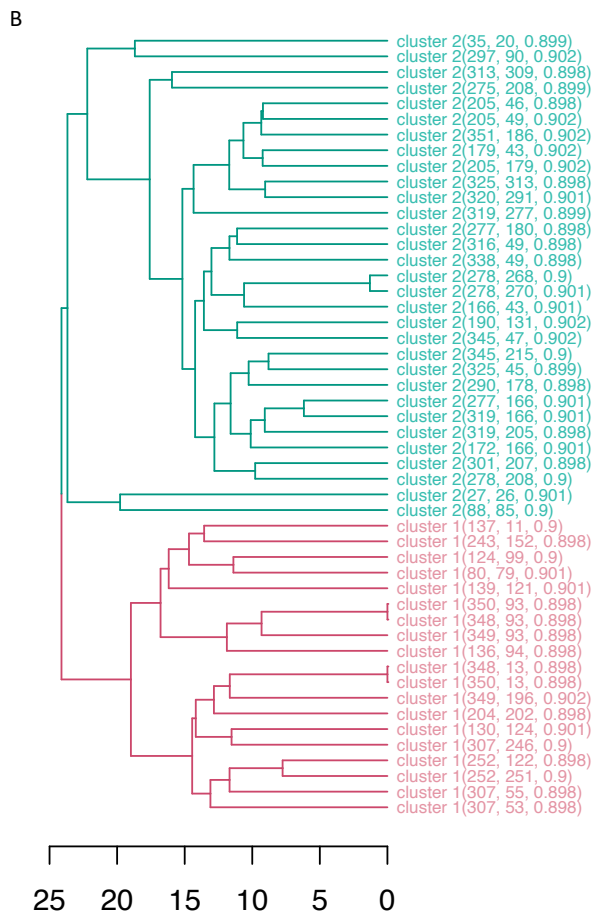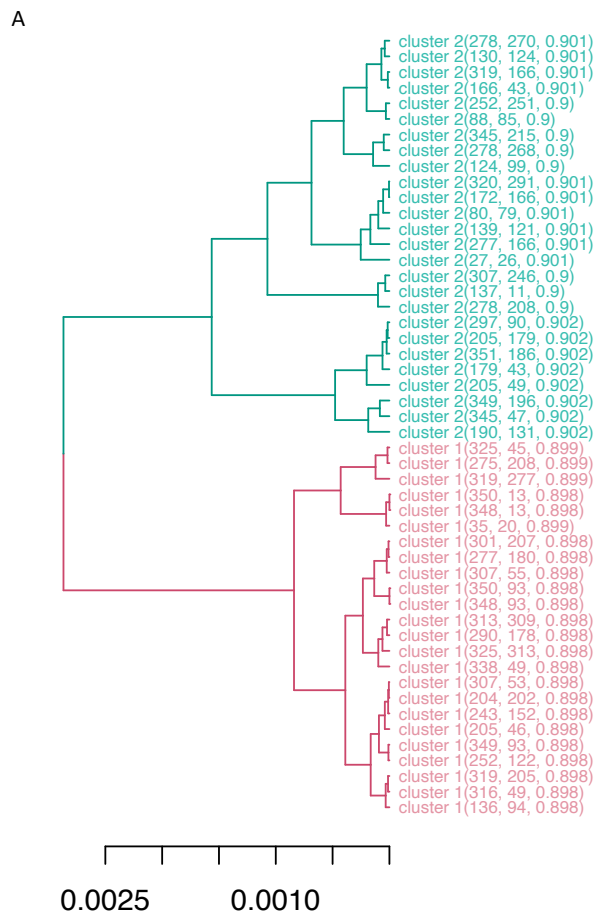
233

234    **Fig. 2.** ICLUST identifies distinct clustering in WHO data. A subset of scatterplots were obtained

235    by selecting all pairwise correlations in the WHO dataset with Pearson correlation between

236    0.8975 and 0.9025, with two outlier plots removed. Clustering assignment was determined by

237    selecting the number of clusters with the highest silhouette score. (*A*) Dendrogram obtained by

238    hierarchical clustering of scatterplots based on correlation strength alone. (*B*) Dendrogram

239    obtained by hierarchical clustering of scatterplots based on Euclidean distance between 4096-

240    component feature vectors of the images as processed by ICLUST. (*C*) Average image in each

241    cluster as determined by correlation strength-based clustering, corresponding to the dendrogram

242    in (A). (*D*) Average image in each cluster according to visual similarity clustering via ICLUST. (E).

243    Principal Coordinate Analysis (PCA) of the scatterplots based on the 4096-component feature

244    vector for each image with colors pertaining to the clustering obtained in (B).

245

246    **Fig. 3.** ICLUST identifies distinct clustering in airline data. All scatterplots in the dataset were

247    plotted and clustered using (*A*) correlation strength alone and (*B*) image 4096-component feature

248    vectors. (*C*) Average image in each cluster as determined by correlation strength-based

249    clustering, corresponding to the dendrogram in (A). (*D*) Average image in each cluster according

250    to visual similarity clustering via ICLUST. (E). Principal Coordinate Analysis (PCA) of the

251    scatterplots based on the 4096-component feature vector for each image with colors pertaining to

252    the clustering obtained in (B). (*F*) Examples of scatterplots with similar correlation size but

253    different visual shape. (*G*) Example of correlations with similar shape but different correlation

254    strength.

A

σ = 0.1, VMS = 1

σ = 0.25, VMS = 1

σ = 0.5, VMS = 0.89

σ = 1, VMS = 0.64

B

Anscombe's Quartet

C

σ = 0.1

σ = 0.5

σ = 1

**Figure 1**

A

cluster 2(278, 270, 0.901)
cluster 2(130, 124, 0.901)
cluster 2(319, 166, 0.901)
cluster 2(166, 43, 0.901)
cluster 2(252, 251, 0.9)
cluster 2(88, 85, 0.9)
cluster 2(345, 215, 0.9)
cluster 2(278, 268, 0.9)
cluster 2(124, 99, 0.9)
cluster 2(320, 291, 0.901)
cluster 2(172, 166, 0.901)
cluster 2(80, 79, 0.901)
cluster 2(139, 121, 0.901)
cluster 2(277, 166, 0.901)
cluster 2(27, 26, 0.901)
cluster 2(307, 246, 0.9)
cluster 2(137, 11, 0.9)
cluster 2(278, 208, 0.9)
cluster 2(297, 90, 0.902)
cluster 2(205, 179, 0.902)
cluster 2(351, 186, 0.902)
cluster 2(179, 43, 0.902)
cluster 2(205, 49, 0.902)
cluster 2(349, 196, 0.902)
cluster 2(345, 47, 0.902)
cluster 2(190, 131, 0.902)
cluster 1(325, 45, 0.899)
cluster 1(275, 208, 0.899)
cluster 1(319, 277, 0.899)
cluster 1(350, 13, 0.898)
cluster 1(348, 13, 0.898)
cluster 1(35, 20, 0.899)
cluster 1(301, 207, 0.898)
cluster 1(277, 180, 0.898)
cluster 1(307, 55, 0.898)
cluster 1(350, 93, 0.898)
cluster 1(348, 93, 0.898)
cluster 1(313, 309, 0.898)
cluster 1(290, 178, 0.898)
cluster 1(325, 313, 0.898)
cluster 1(338, 49, 0.898)
cluster 1(307, 53, 0.898)
cluster 1(204, 202, 0.898)
cluster 1(243, 152, 0.898)
cluster 1(205, 46, 0.898)
cluster 1(349, 93, 0.898)
cluster 1(252, 122, 0.898)
cluster 1(319, 205, 0.898)
cluster 1(316, 49, 0.898)
cluster 1(136, 94, 0.898)

0.0025   0.0010

B

cluster 2(35, 20, 0.899)
cluster 2(297, 90, 0.902)
cluster 2(313, 309, 0.898)
cluster 2(275, 208, 0.899)
cluster 2(205, 46, 0.898)
cluster 2(205, 49, 0.902)
cluster 2(351, 186, 0.902)
cluster 2(179, 43, 0.902)
cluster 2(205, 179, 0.902)
cluster 2(325, 313, 0.898)
cluster 2(320, 291, 0.901)
cluster 2(319, 277, 0.899)
cluster 2(277, 180, 0.898)
cluster 2(316, 49, 0.898)
cluster 2(338, 49, 0.898)
cluster 2(278, 246, 0.9)
cluster 2(278, 270, 0.901)
cluster 2(166, 43, 0.901)
cluster 2(190, 131, 0.902)
cluster 2(345, 47, 0.898)
cluster 2(345, 215, 0.9)
cluster 2(325, 45, 0.899)
cluster 2(290, 178, 0.898)
cluster 2(277, 166, 0.901)
cluster 2(319, 166, 0.901)
cluster 2(319, 205, 0.898)
cluster 2(172, 166, 0.901)
cluster 2(301, 207, 0.898)
cluster 2(278, 208, 0.9)
cluster 2(27, 26, 0.901)
cluster 2(88, 85, 0.9)
cluster 1(137, 11, 0.9)
cluster 1(243, 152, 0.898)
cluster 1(124, 99, 0.9)
cluster 1(80, 79, 0.901)
cluster 1(139, 121, 0.901)
cluster 1(350, 93, 0.898)
cluster 1(348, 93, 0.898)
cluster 1(349, 93, 0.898)
cluster 1(136, 94, 0.898)
cluster 1(348, 13, 0.898)
cluster 1(350, 13, 0.898)
cluster 1(349, 196, 0.902)
cluster 1(204, 202, 0.898)
cluster 1(130, 124, 0.901)
cluster 1(307, 246, 0.9)
cluster 1(252, 122, 0.898)
cluster 1(252, 251, 0.9)
cluster 1(307, 55, 0.898)
cluster 1(307, 53, 0.898)

25  20  15  10  5  0

C

Cluster 1

Cluster 2

D

E

pc2

pc1

**Figure 2**

**Figure 3**