

```

1   Course: Exploratory_Data_Analysis
2   Lesson: Dimension_Reduction
3
4   - Class: text
5   Output: "Dimension_Reduction. (Slides for this and other Data Science courses may be
6   found at github https://github.com/DataScienceSpecialization/courses/. If you care to
7   use them, they must be downloaded as a zip file and viewed locally. This lesson
8   corresponds to 04_ExploratoryAnalysis/dimensionReduction.)"
9
10
11  - Class: text
12  Output: In this lesson we'll discuss principal component analysis (PCA) and singular
13  value decomposition (SVD), two important and related techniques of dimension
14  reduction. This last entails processes which finding subsets of variables in datasets
15  that contain their essences. PCA and SVD are used in both the exploratory phase and
16  the more formal modelling stage of analysis. We'll focus on the exploratory phase and
17  briefly touch on some of the underlying theory.
18
19  - Class: figure
20  Output: We'll begin with a motivating example - random data.
21  Figure: showRanMat.R
22  FigureType: new
23
24  - Class: cmd_question
25  Output: This is dataMatrix, a matrix of 400 random normal numbers (mean 0 and
26  standard deviation 1). We're displaying it with the R command image. Run the R
27  command head with dataMatrix as its argument to see what dataMatrix looks like.
28  CorrectAnswer: head(dataMatrix)
29  AnswerTests: omnitest(correctExpr='head(dataMatrix)')
30  Hint: Type head(dataMatrix) at the command prompt.
31
32  - Class: cmd_question
33  Output: So we see that dataMatrix has 10 columns (and hence 40 rows) of random
34  numbers. The image here looks pretty random. Let's see how the data clusters. Run the
35  R command heatmap with dataMatrix as its only argument.
36  CorrectAnswer: heatmap(dataMatrix)
37  AnswerTests: omnitest(correctExpr='heatmap(dataMatrix)')
38  Hint: Type heatmap(dataMatrix) at the command prompt.
39
40  - Class: text
41  Output: We can see that even with the clustering that heatmap provides, permuting the
42  rows (observations) and columns (variables) independently, the data still looks
43  random.
44
45  - Class: cmd_question
46  Output: Let's add a pattern to the data. We've put some R code in the file addPatt.R
47  for you. Run the command myedit with the single argument "addPatt.R" (make sure to
48  use the quotation marks) to see the code. You might have to click your cursor in the
49  console after you do this to keep from accidentally changing the file.
50  CorrectAnswer: myedit("addPatt.R")
51  AnswerTests: omnitest(correctExpr='myedit("addPatt.R")')
52  Hint: Type myedit("addPatt.R") at the command prompt.
53
54  - Class: mult_question
55  Output: Look at the code. Will every row of the matrix have a pattern added to it?
56  AnswerChoices: Yes; No
57  CorrectAnswer: No
58  AnswerTests: omnitest(correctVal='No')
59  Hint: What does the coinflip do?
60
61  - Class: mult_question
62  Output: So whether or not a row gets modified by a pattern is determined by a coin
63  flip. Will the added pattern affect every column in the affected row?
64  AnswerChoices: Yes; No
65  CorrectAnswer: No
66  AnswerTests: omnitest(correctVal='No')
67  Hint: The expression rep(c(0,3),each=5) creates the 10-long vector
68  (0,0,0,0,0,3,3,3,3,3) which is added to the rows chosen by the coin flip.
69
70

```

```

51 - Class: text
52 Output: So in rows affected by the coin flip, the 5 left columns will still have a
    mean of 0 but the right 5 columns will have a mean closer to 3.
53
54 - Class: cmd_question
55 Output: Now to execute this code, run the R command source with 2 arguments. The
    first is the filename (in quotes), "addPatt.R", and the second is the argument local
    set equal to TRUE.
56 CorrectAnswer: source("addPatt.R", local=TRUE)
57 AnswerTests: omnitest(correctExpr='source("addPatt.R", local=TRUE)')
58 Hint: Type source("addPatt.R", local=TRUE) at the command prompt.
59
60 - Class: figure
61 Output: Here's the image of the altered dataMatrix after the pattern has been added.
    The pattern is clearly visible in the columns of the matrix. The right half is
    yellower or hotter, indicating higher values in the matrix.
62 Figure: showRanMat.R
63 FigureType: new
64
65 - Class: cmd_question
66 Output: Now run the R command heatmap again with dataMatrix as its only argument.
    This will perform a hierarchical cluster analysis on the matrix.
67 CorrectAnswer: heatmap(dataMatrix)
68 AnswerTests: omnitest(correctExpr='heatmap(dataMatrix)')
69 Hint: Type heatmap(dataMatrix) at the command prompt.
70
71 - Class: text
72 Output: Again we see the pattern in the columns of the matrix. As shown in the
    dendrogram at the top of the display, these split into 2 clusters, the lower numbered
    columns (1 through 5) and the higher numbered ones (6 through 10). Recall from the
    code in addPatt.R that for rows selected by the coinflip the last 5 columns had 3
    added to them. The rows still look random.
73
74 - Class: figure
75 Output: Now consider this picture. On the left is an image similar to the heatmap of
    dataMatrix you just plotted. It is an image plot of the output of hclust(), a
    hierarchical clustering function applied to dataMatrix. Yellow indicates "hotter" or
    higher values than red. This is consistent with the pattern we applied to the data
    (increasing the values for some of the rightmost columns).
76 Figure: showPatt.R
77 FigureType: new
78
79 - Class: text
80 Output: The middle display shows the mean of each of the 40 rows (along the x-axis).
    The rows are shown in the same order as the rows of the heat matrix on the left. The
    rightmost display shows the mean of each of the 10 columns. Here the column numbers
    are along the x-axis and their means along the y.
81
82 - Class: text
83 Output: We see immediately the connection between the yellow (hotter) portion of the
    cluster image and the higher row means, both in the upper right portion of the
    displays. Similarly, the higher valued column means are in the right half of that
    display and lower column means are in the left half.
84
85 - Class: text
86 Output: Now we'll talk a little theory. Suppose you have 1000's of multivariate
    variables  $X_1, \dots, X_n$ . By multivariate we mean that each  $X_i$  contains many
    components, i.e.,  $X_i = (X_{i1}, \dots, X_{im})$ . However, these variables
    (observations) and their components might be correlated to one another.
87
88 - Class: mult_question
89 Output: Which of the following would be an example of variables correlated to one
    another?
90 AnswerChoices: Heights and weights of members of a family; Today's weather and a
    butterfly's wing position; The depth of the Atlantic Ocean and what you eat for
    breakfast
91 CorrectAnswer: Heights and weights of members of a family
92 AnswerTests: omnitest(correctVal='Heights and weights of members of a family')
93 Hint: Which choice is the only one that makes sense?

```

```

94
95 - Class: text
96 Output: As data scientists, we'd like to find a smaller set of multivariate variables
    that are uncorrelated AND explain as much variance (or variability) of the data as
    possible. This is a statistical approach.
97
98 - Class: text
99 Output: In other words, we'd like to find the best matrix created with fewer
    variables (that is, a lower rank matrix) that explains the original data. This is
    related to data compression.
100
101 - Class: text
102 Output: Two related solutions to these problems are PCA which stands for Principal
    Component Analysis and SVD, Singular Value Decomposition. This latter simply means
    that we express a matrix X of observations (rows) and variables (columns) as the
    product of 3 other matrices, i.e.,  $X = UDV^t$ . This last term ( $V^t$ ) represents the
    transpose of the matrix V.
103
104 - Class: text
105 Output: Here U and V each have orthogonal (uncorrelated) columns. U's columns are the
    left singular vectors of X and V's columns are the right singular vectors of X. D is
    a diagonal matrix, by which we mean that all of its entries not on the diagonal are
    0. The diagonal entries of D are the singular values of X.
106
107 - Class: cmd_question
108 Output: To illustrate this idea we created a simple example matrix called mat.
    Look at it now.
109 CorrectAnswer: mat
110 AnswerTests: omnitest(correctExpr='mat')
111 Hint: Type mat at the command prompt.
112
113 - Class: cmd_question
114 Output: So mat is a 2 by 3 matrix. Lucky for us R provides a function to perform
    singular value decomposition. It's called, unsurprisingly, svd. Call it now with a
    single argument, mat.
115 CorrectAnswer: svd(mat)
116 AnswerTests: omnitest(correctExpr='svd(mat)')
117 Hint: Type svd(mat) at the command prompt.
118
119 - Class: cmd_question
120 Output: We see that the function returns 3 components, d which holds 2 diagonal
    elements, u, a 2 by 2 matrix, and v, a 3 by 2 matrix. We stored the diagonal entries
    in a diagonal matrix for you, diag, and we also stored u and v in the variables matu
    and matv respectively. Multiply matu by diag by  $t(matv)$  to see what you get.
    (This last expression represents the transpose of matv in R). Recall that in R
    matrix multiplication requires you to use the operator %%.
121 CorrectAnswer: matu %% diag %% t(matv)
122 AnswerTests: omnitest(correctExpr='matu %% diag %% t(matv)')
123 Hint: Type matu %% diag %% t(matv) at the command prompt.
124
125 - Class: text
126 Output: So we did in fact get mat back. That's a relief! Note that this type of
    decomposition is NOT unique.
127
128 - Class: text
129 Output: Now we'll talk a little about PCA, Principal Component Analysis, "a simple,
    non-parametric method for extracting relevant information from confusing data sets."
    We're quoting here from a very nice concise paper on this subject which can be found
    at http://arxiv.org/pdf/1404.1100.pdf. The paper by Jonathon Shlens of Google
    Research is called, A Tutorial on Principal Component Analysis.
130
131 - Class: text
132 Output: Basically, PCA is a method to reduce a high-dimensional data set to its
    essential elements (not lose information) and explain the variability in the data. We
    won't go into the mathematical details here, (R has a function to perform PCA), but
    you should know that SVD and PCA are closely related.
133
134 - Class: cmd_question
135 Output: We'll demonstrate this now. First we have to scale mat, our simple example

```

data matrix. This means that we subtract the column mean from every element and divide the result by the column standard deviation. Of course R has a command, `scale`, that does this for you. Run `svd` on `scale` of `mat`.

CorrectAnswer: `svd(scale(mat))`

AnswerTests: `omnittest(correctExpr='svd(scale(mat))')`

Hint: Type `svd(scale(mat))` at the command prompt.

- **Class:** `cmd_question`

Output: Now run the R program `prcomp` on `scale(mat)`. This will give you the principal components of `mat`. See if they look familiar.

CorrectAnswer: `prcomp(scale(mat))`

AnswerTests: `omnittest(correctExpr='prcomp(scale(mat))')`

Hint: Type `prcomp(scale(mat))` at the command prompt.

- **Class:** `text`

Output: Notice that the principal components of the scaled matrix, shown in the Rotation component of the `prcomp` output, ARE the columns of `V`, the right singular values. Thus, PCA of a scaled matrix yields the `V` matrix (right singular vectors) of the same scaled matrix.

- **Class:** `text`

Output: Now that we covered the theory let's return to our bigger matrix of random data into which we had added a fixed pattern for some rows selected by coinflips. The pattern effectively shifted the means of the rows and columns.

- **Class:** `figure`

Output: Here's a picture showing the relationship between PCA and SVD for that bigger matrix. We've plotted 10 points (5 are squished together in the bottom left corner). The x-coordinates are the elements of the first principal component (output from `prcomp`), and the y-coordinates are the elements of the first column of `V`, the first right singular vector (gotten from running `svd`). We see that the points all lie on the 45 degree line represented by the equation $y=x$. So the first column of `V` IS the first principal component of our bigger data matrix.

Figure: `showRel.R`

FigureType: `new`

- **Class:** `cmd_question`

Output: To prove we're not making this up, we've run `svd` on `dataMatrix` and stored the result in the object `svd1`. This has 3 components, `d`, `u` and `v`. look at the first column of `V` now. It can be viewed by using the `svd1$v[,1]` notation.

CorrectAnswer: `svd1$v[,1]`

AnswerTests: `omnittest(correctExpr='svd1$v[,1]')`

Hint: Type `svd1$v[,1]` at the command prompt.

- **Class:** `text`

Output: See how these values correspond to those plotted? Five of the entries are slightly to the left of the point $(-0.4, -0.4)$, two more are negative (to the left of $(0,0)$), and three are positive (to the right of $(0,0)$).

- **Class:** `figure`

Output: Here we again show the clustered data matrix on the left. Next to it we've plotted the first column of the `U` matrix associated with the scaled data matrix. This is the first LEFT singular vector and it's associated with the ROW means of the clustered data. You can see the clear separation between the top 24 (around -0.2) row means and the bottom 16 (around 0.2). We don't show them but note that the other columns of `U` don't show this pattern so clearly.

Figure: `showUV.R`

FigureType: `new`

- **Class:** `text`

Output: The rightmost display shows the first column of the `V` matrix associated with the scaled and clustered data matrix. This is the first RIGHT singular vector and it's associated with the COLUMN means of the clustered data. You can see the clear separation between the left 5 column means (between -0.1 and 0.1) and the right 5 column means (all below -0.4). As with the left singular vectors, the other columns of `V` don't show this pattern as clearly as this first one does.

- **Class:** `text`

Output: So the singular value decomposition automatically picked up these patterns,

the differences in the row and column means.

```
176
177 - Class: cmd_question
178 Output: Why were the first columns of both the U and V matrices so special? Well as
it happens, the D matrix of the SVD explains this phenomenon. It is an aspect of SVD
called variance explained. Recall that D is the diagonal matrix sandwiched in between
U and V^t in the SVD representation of the data matrix. The diagonal entries of D are
like weights for the U and V columns accounting for the variation in the data.
They're given in decreasing order from highest to lowest. Look at these diagonal
entries now. Recall that they're stored in svd1$d.
179 CorrectAnswer: svd1$d
180 AnswerTests: omnitest(correctExpr='svd1$d')
181 Hint: Type svd1$d at the command prompt.
182
183 - Class: figure
184 Output: Here's a display of these values (on the left). The first one (12.46) is
significantly bigger than the others. Since we don't have any units specified, to the
right we've plotted the proportion of the variance each entry represents. We see that
the first entry accounts for about 40% of the variance in the data. This explains why
the first columns of the U and V matrices respectively showed the distinctive
patterns in the row and column means so clearly.
185 Figure: showVar.R
186 FigureType: new
187
188 - Class: cmd_question
189 Output: Now we'll show you another simple example of how SVD explains variance.
We've created a 40 by 10 matrix, constantMatrix. Use the R command head with
constantMatrix as its argument to see the top rows.
190 CorrectAnswer: head(constantMatrix)
191 AnswerTests: omnitest(correctExpr='head(constantMatrix)')
192 Hint: Type head(constantMatrix) at the command prompt.
193
194 - Class: cmd_question
195 Output: The rest of the rows look just like these. You can see that the left 5
columns are all 0's and the right 5 columns are all 1's. We've run svd with
constantMatrix as its argument for you and stored the result in svd2. Look at the
diagonal component, d, of svd2 now.
196 CorrectAnswer: svd2$d
197 AnswerTests: omnitest(correctExpr='svd2$d')
198 Hint: Type svd2$d at the command prompt.
199
200 - Class: mult_question
201 Output: Which index holds the largest entry of the svd2$d?
202 AnswerChoices: 9; 1; 10; 5
203 CorrectAnswer: 1
204 AnswerTests: omnitest(correctVal='1')
205 Hint: Which choice has a positive exponent? Notice that the entries are given in
scientific notation, where xey means  $x * 10^y$ .
206
207 - Class: figure
208 Output: So the first entry by far dominates the others. Here the picture on the left
shows the heat map of constantMatrix. You can see how the left columns differ from
the right ones. The middle plot shows the values of the singular values of the
matrix, i.e., the diagonal elements which are the entries of svd2$d. Nine of these
are 0 and the first is a little above 14. The third plot shows the proportion of the
total each diagonal element represents.
209 Figure: showSimple.R
210 FigureType: new
211
212 - Class: mult_question
213 Output: According to the plot, what percentage of the total variation does the first
diagonal element account for?
214 AnswerChoices: 100%; 90%; 50%; 0%
215 CorrectAnswer: 100%
216 AnswerTests: omnitest(correctVal='100%')
217 Hint: The first element is at 1.0. Which percent does the number 1.0 represent?
218
219 - Class: text
220 Output: So what does this mean? Basically that the data is one-dimensional. Only 1
```

```

221     piece of information, namely which column an entry is in, determines its value.
222 - Class: figure
223 Output: Now let's return to our random 40 by 10 dataMatrix and consider a slightly
more complicated example in which we add 2 patterns to it. Again we'll choose which
rows to tweak using coinflips. Specifically, for each of the 40 rows we'll flip 2
coins. If the first coinflip is heads, we'll add 5 to each entry in the right 5
columns of that row, and if the second coinflip is heads, we'll add 5 to just the
even columns of that row.
224 Figure: twoPatts.R
225 FigureType: new
226
227 - Class: figure
228 Output: So here's the image of the scaled data matrix on the left. We can see both
patterns, the clear difference between the left 5 and right 5 columns, but also,
slightly less visible, the alternating pattern of the columns. The other plots show
the true patterns that were added into the affected rows. The middle plot shows the
true difference between the left and right columns, while the rightmost plot shows
the true difference between the odd numbered and even-numbered columns.
229 Figure: showTrue.R
230 FigureType: new
231
232 - Class: text
233 Output: The question is, "Can our analysis detect these patterns just from the data?"
Let's see what SVD shows. Since we're interested in patterns on columns we'll look at
the first two right singular vectors (columns of V) to see if they show any evidence
of the patterns.
234
235 - Class: figure
236 Output: Here we see the 2 right singular vectors plotted next to the image of the
data matrix. The middle plot shows the first column of V and the rightmost plot the
second. The middle plot does show that the last 5 columns have higher entries than
the first 5. This picks up, or at least alludes to, the first pattern we added in
which affected the last 5 columns of the matrix. The rightmost plot, showing the
second column of V, looks more random. However, closer inspection shows that the
entries alternate or bounce up and down as you move from left to right. This hints at
the second pattern we added in which affected only even columns of selected rows.
237 Figure: showVs.R
238 FigureType: new
239
240 - Class: cmd_question
241 Output: To see this more closely, look at the first 2 columns of the v component. We
stored the SVD output in the svd object svd2.
242 CorrectAnswer: svd2$V[,1:2]
243 AnswerTests: ANY_of_exprs('svd2$V[,1:2]', 'svd2$V[,c(1,2)]', 'svd2$V[,c(1:2)]')
244 Hint: Type svd2$V[,1:2] at the command prompt.
245
246 - Class: figure
247 Output: Seeing the 2 columns side by side, we see that the values in both columns
alternately increase and decrease. However, we knew to look for this pattern, so
chances are, you might not have noticed this pattern if you hadn't known it was
there. This example is meant to show you that it's hard to see patterns, even
straightforward ones.
248 Figure: clearPlot.R
249
250 - Class: cmd_question
251 Output: Now look at the entries of the diagonal matrix d resulting from the svd.
Recall that we stored this output for you in the svd object svd2.
252 CorrectAnswer: svd2$d
253 AnswerTests: omnitest(correctExpr='svd2$d')
254 Hint: Type svd2$d at the command prompt.
255
256 - Class: figure
257 Output: We see that the first element, 14.55, dominates the others. Here's the plot
of these diagonal elements of d. The left shows the numerical entries and the right
show the percentage of variance each entry explains.
258 Figure: showDvar.R
259 FigureType: new
260

```



```

261 - Class: mult_question
262 Output: According to the plot, how much of the variance does the second element
      account for?
263 AnswerChoices: 18%; 53%; 11%; .1%
264 CorrectAnswer: 18%
265 AnswerTests: omnitest(correctVal='18%')
266 Hint: At what height does the second point fall in the chart on the right?
267
268 - Class: text
269 Output: So the first element which showed the difference between the left and right
      halves of the matrix accounts for roughly 50% of the variation in the matrix, and the
      second element which picked up the alternating pattern accounts for 18% of the
      variance. The remaining elements account for smaller percentages of the variation.
      This indicates that the first pattern is much stronger than the second. Also the two
      patterns confound each other so they're harder to separate and see clearly. This is
      what often happens with real data.
270
271 - Class: text
272 Output: Now you're probably convinced that SVD and PCA are pretty cool and useful as
      tools for analysis, but one problem with them that you should be aware of, is that
      they cannot deal with MISSING data. Neither of them will work if any data in the
      matrix is missing. (You'll get error messages from R in red if you try.) Missing data
      is not unusual, so luckily we have ways to work around this problem. One we'll just
      mention is called imputing the data.
273
274 - Class: text
275 Output: This uses the k nearest neighbors to calculate a values to use in place of
      the missing data. You may want to specify an integer k which indicates how many
      neighbors you want to average to create this replacement value. The bioconductor
      package (http://bioconductor.org) has an impute package which you can use to fill in
      missing data. One specific function in it is impute.knn.
276
277 - Class: text
278 Output: We'll move on now to a final example of the power of singular value
      decomposition and principal component analysis and how they work as a data
      compression technique.
279
280 - Class: figure
281 Output: Consider this low resolution image file showing a face. We'll use SVD and see
      how the first several components contain most of the information in the file so that
      storing a huge matrix might not be necessary.
282 Figure: showFace.R
283 FigureType: new
284
285 - Class: cmd_question
286 Output: The image data is stored in the matrix faceData. Run the R command dim on
      faceData to see how big it is.
287 CorrectAnswer: dim(faceData)
288 AnswerTests: omnitest(correctExpr='dim(faceData)')
289 Hint: Type dim(faceData) at the command prompt.
290
291 - Class: figure
292 Output: So it's not that big of a file but we want to show you how to use what you
      learned in this lesson. We've done the SVD and stored it in the object svd1 for you.
      Here's the plot of the variance explained.
293 Figure: showFaceSVD.R
294 FigureType: new
295
296 - Class: mult_question
297 Output: According to the plot what percentage of the variance is explained by the
      first singular value?
298 AnswerChoices: 23; 15; 40; 100
299 CorrectAnswer: 40
300 AnswerTests: omnitest(correctVal='40')
301 Hint: At what height does the first (leftmost) point fall in the plot?
302
303 - Class: text
304 Output: So 40% of the variation in the data matrix is explained by the first
      component, 22% by the second, and so forth. It looks like most of the variation is

```

contained in the first 10 components. How can we check this out? Can we try to create an approximate image using only a few components?

```
305
306 - Class: text
307 Output: Recall that the data matrix X is the product of 3 matrices, that is  $X=UDV^t$ .
      These are precisely what you get when you run svd on the matrix X.
308
309 - Class: text
310 Output: Suppose we create the product of pieces of these, say the first columns of U
      and V and the first element of D. The first column of U can be interpreted as a 32 by
      1 matrix (recall that faceData was a 32 by 32 matrix), so we can multiply it by the
      first element of D, a 1 by 1 matrix, and get a 32 by 1 matrix result. We can multiply
      that by the transpose of the first column of V, which is the first principal
      component. (We have to use the transpose of V's column to make it a 1 by 32 matrix in
      order to do the matrix multiplication properly.)
311
312 - Class: text
313 Output: Alas, that is how we do it in theory, but in R using only one element of d
      means it's a constant. So we have to do the matrix multiplication with the %%
      operator and the multiplication by the constant (svd1$d[1]) with the regular
      multiplication operator *.
314
315 - Class: cmd_question
316 Output: Try this now and put the result in the variable a1. Recall that svd1$u,
      svd1$d, and svd1$v contain all the information you need. NOTE that because of the
      peculiarities of R's casting, if you do the scalar multiplication with the * operator
      first (before the matrix multiplication with the %% operator) you MUST enclose the 2
      arguments (svd1$u[,1] and svd1$d[1]) in parentheses.
317 CorrectAnswer: a1 <- svd1$u[,1] %% t(svd1$v[,1]) * svd1$d[1]
318 AnswerTests: expr_creates_var("a1"); ANY_of_exprs('a1 <- svd1$u[,1] %% t(svd1$v[,1])
      * svd1$d[1]', 'a1 <- (svd1$d[1] * svd1$u[,1]) %% t(svd1$v[,1])', 'a1 <- (svd1$u[,1] *
      svd1$d[1]) %% t(svd1$v[,1])')
319 Hint: Type a1 <- (svd1$u[,1] * svd1$d[1]) %% t(svd1$v[,1]) OR a1 <- svd1$u[,1] %%
      t(svd1$v[,1]) * svd1$d[1] at the command prompt.
320
321 - Class: cmd_question
322 Output: Now to look at it as an image. We wrote a function for you called myImage
      which takes a single argument, a matrix of data to display using the R function
      image. Run it now with a1 as its argument.
323 CorrectAnswer: myImage(a1)
324 AnswerTests: omnittest(correctExpr='myImage(a1)')
325 Hint: Type myImage(a1) at the command prompt.
326
327 - Class: text
328 Output: It might not look like much but it's a good start. Now we'll try the same
      experiment but this time we'll use 2 elements from each of the 3 SVD terms.
329
330 - Class: cmd_question
331 Output: Create the matrix a2 as the product of the first 2 columns of svd1$u, a
      diagonal matrix using the first 2 elements of svd1$d, and the transpose of the first
      2 columns of svd1$v. Since all of your multiplicands are matrices you have to use
      only the operator %% AND you DON'T need parentheses. Also, you must use the R
      function diag with svd1$d[1:2] as its sole argument to create the proper diagonal
      matrix. Remember, matrix multiplication is NOT commutative so you have to put the
      multiplicands in the correct order. Please use the 1:2 notation and not the c(m:n),
      i.e., the concatenate function, when specifying the columns.
332 CorrectAnswer: a2 <- svd1$u[,1:2] %% diag(svd1$d[1:2]) %% t(svd1$v[,1:2])
333 AnswerTests: expr_creates_var("a2"); omnittest(correctExpr='a2 <- svd1$u[,1:2] %%
      diag(svd1$d[1:2]) %% t(svd1$v[,1:2])')
334 Hint: Type a2 <- svd1$u[,1:2] %% diag(svd1$d[1:2]) %% t(svd1$v[,1:2]) at the
      command prompt.
335
336 - Class: cmd_question
337 Output: Use myImage again to see how a2 displays.
338 CorrectAnswer: myImage(a2)
339 AnswerTests: omnittest(correctExpr='myImage(a2)')
340 Hint: Type myImage(a2) at the command prompt.
341
342 - Class: cmd_question
```



```

343 Output: We're starting to see slightly more detail, and maybe if you squint you see a
    grimacing mouth. Now let's see what image results using 5 components. From our plot
    of the variance explained 5 components covered a sizeable percentage of the
    variation. To save typing, use the up arrow to recall the command which created a2
    and replace the a2 and assignment arrow with the call to myImage, and change the
    three occurrences of 2 to 5.
344 CorrectAnswer: myImage(svd1$u[,1:5] %*% diag(svd1$d[1:5]) %*% t(svd1$v[,1:5]))
345 AnswerTests: omnitest(correctExpr='myImage(svd1$u[,1:5] %*% diag(svd1$d[1:5]) %*%
    t(svd1$v[,1:5]))')
346 Hint: Type myImage(svd1$u[,1:5] %*% diag(svd1$d[1:5]) %*% t(svd1$v[,1:5])) at the
    command prompt.
347
348 - Class: cmd_question
349 Output: Certainly much better. Clearly a face is appearing with eyes, nose, ears, and
    mouth recognizable. Again, use the up arrow to recall the last command (calling
    myImage with a matrix product argument) and change the 5's to 10's. We'll see how
    this image looks.
350 CorrectAnswer: myImage(svd1$u[,1:10] %*% diag(svd1$d[1:10]) %*% t(svd1$v[,1:10]))
351 AnswerTests: omnitest(correctExpr='myImage(svd1$u[,1:10] %*% diag(svd1$d[1:10])
    %*% t(svd1$v[,1:10]))')
352 Hint: Type myImage(svd1$u[,1:10] %*% diag(svd1$d[1:10]) %*% t(svd1$v[,1:10])) at
    the command prompt.
353
354 - Class: text
355 Output: Now that's pretty close to the original which was low resolution to begin
    with, but you can see that 10 components really do capture the essence of the image.
    Singular value decomposition is a good way to approximate data without having to
    store a lot.
356
357 - Class: text
358 Output: We'll close now with a few comments. First, when reducing dimensions you have
    to pay attention to the scales on which different variables are measured and make
    sure that all your data is in consistent units. In other words, scales of your data
    matter. Second, principal components and singular values may mix real patterns, as we
    saw in our simple 2-pattern example, so finding and separating out the real patterns
    require some detective work. Let's do a quick review now.
359
360 - Class: mult_question
361 Output: Which of the following cliches LEAST captures the essence of dimension
    reduction?
362 AnswerChoices: find the needle in the haystack; see the forest through the trees;
    separate the wheat from the chaff; a face that could launch a 1000 ships
363 CorrectAnswer: a face that could launch a 1000 ships
364 AnswerTests: omnitest(correctVal='a face that could launch a 1000 ships')
365 Hint: Which choice fails to deal with discerning differences between the valuable and
    the invaluable.
366
367 - Class: mult_question
368 Output: A matrix X has the singular value decomposition  $UDV^t$ . The principal
    components of X are ?
369 AnswerChoices: the columns of U; the rows of U; the columns of V; the rows of V
370 CorrectAnswer: the columns of V
371 AnswerTests: omnitest(correctVal='the columns of V')
372 Hint: Recall the simple example where we ran prcomp and svd on the same scaled matrix
    and saw that the columns of V matched the rotations of the prcomp output.
373
374 - Class: mult_question
375 Output: A matrix X has the singular value decomposition  $UDV^t$ . The singular values of
    X are found where?
376 AnswerChoices: the columns of U; the columns of D; the columns of V; the diagonal
    elements of D;
377 CorrectAnswer: the diagonal elements of D
378 AnswerTests: omnitest(correctVal='the diagonal elements of D')
379 Hint: Recall that U and V give us vectors and D gave us values.
380
381 - Class: mult_question
382 Output: True or False? PCA and SVD are totally unrelated.
383 AnswerChoices: False; True
384 CorrectAnswer: False

```

```
385     AnswerTests: omnitest(correctVal='False')
386     Hint: Recall the question about principal components and their relationship to V.
387
388 - Class: mult_question
389     Output: True or False? D gives the singular values of a matrix in decreasing order of
weight.
390     AnswerChoices: True; False
391     CorrectAnswer: True
392     AnswerTests: omnitest(correctVal='True')
393     Hint: Recall that the first value accounted for the highest percentage of variation
in the data.
394
395
396 - Class: text
397     Output: Congratulations! We hope you enjoyed making faces and that this lesson didn't
reduce the dimensions of your understanding.
398
399 - Class: mult_question
400     Output: "Would you like to receive credit for completing this course on
401 Coursera.org?"
402     CorrectAnswer: NULL
403     AnswerChoices: Yes;No
404     AnswerTests: coursera_on_demand()
405     Hint: ""
406
```