

```

1: //////////////////////////////////////
2: How to deal with Principal Component Analysis
3:
4: maXbox Starter 58 - Data Science with PCA
5:
6: "Errors should never pass silently. (Unless explicitly silenced.)"
7:   - The Zen of Python
8:
9: Principal component analysis (PCA) is often the first thing to try out if
   you want to cut down the number of features and do not know what feature
   extraction method to use.
10: PCA is limited as its a linear method, but chances are that it already
   goes far enough for your model to learn well enough.
11: Add to this the strong mathematical properties it offers and the speed at
   which it finds the transformed feature data space and is later able to
   transform between original and transformed features; we can almost
   guarantee that it also will become one of your frequently used machine
   learning tools.
12:
13: This tutor will go straight to an overview to PCA.
14:
15: The script 811_mXpcatest_dmath_datascience.pas (pcatest.pas) (located in
   the demo\console\curfit subdirectory) performs a principal component
   analysis on a set of 4 variables.
16: (Example taken from: P. DAGNELIE, Analyse statistique la plusieurs
   variables, Presses Agronomiques de Gembloux, Belgique, 1982). The program
   prints:
17:
18: ✕ The mean vector and variance-covariance matrix of the original variables
19: ✕ The correlation coefficients between the original variables
20: ✕ The eigenvalues and eigenvectors of the correlation matrix
21: ✕ The correlation coefficients between principal factors and the original
22:   variables
23: ✕ The values of the principal factors for each point
24:
25: Summarizing it, given the original feature space, PCA finds a linear
   projection of itself in a lower dimensional space that has the following
   two properties:
26:
27: • The conserved variance is maximized.
28: • The final reconstruction error (when trying to go back from transformed
29:   features to the original ones) is minimized.
30:
31: As PCA simply transforms the input data, it can be applied both to
   classification and regression problems. In this section, we will use a
   classification task to discuss the method.
32:
33: The script can be found at:
34:
35:   http://www.softwareschule.ch/examples/811_mXpcatest_dmath_datascience.pas
36:   ..\examples\811_mXpcatest_dmath_datascience.pas
37: It may be seen that:
38:
39: • High correlations exist between the original variables, which are
40:   therefore not independent
41:
42: • According to the eigenvalues, the last two principal factors may be
43:   neglected since they represent less than 11 % of the total variance. So,
44:   the original variables depend mainly on the first two factors
45:
46: • The first principal factor is negatively correlated with the second and

```

```

47:  fourth variables, and positively correlated with the third variable
48:
49:  • The second principal factor is positively correlated with the first
50:    variable
51:
52:  • The table of principal factors show that the highest scores are usually
53:    associated with the first two principal factors, in agreement with the
54:    previous results
55:
56:
57: Const
58:   N      = 11; { Number of observations }
59:   Nvar   = 4;  { Number of variables }
60:
61: { Data }
62: var X : array[1..N] of array[1..Nvar] of Float;
63: =
64: (( 87.9, 19.6,   1 ,   1661),
65:  ( 89.9, 15.2,  90.1,   968),
66:  (153 , 19.7,  56.6,  1353),
67:  (132.1, 17 ,   91 ,  1293),
68:  ( 88.8, 18.3,  93.7,  1153),
69:  (220.9, 17.8, 106.9,  1286),
70:  (117.7, 17.8,  65.5,  1104),
71:  (109 , 18.3,  41.8,  1574),
72:  (156.1, 17.8,  57.4,  1222),
73:  (181.5, 16.8, 140.6,   902),
74:  (181.4, 17 ,   74.3,  1150));//}
75:
76: var
77:   XX      : TMatrix; { Data }
78:   M       : TVector; { Mean vector }
79:   V       : TMatrix; { Variance-covariance matrix }
80:   R       : TMatrix; { Correlation matrix }
81:   S       : TVector; { Standard deviations }
82:   Lambda  : TVector; { Eigenvalues of correlation matrix }
83:   C       : TMatrix; { Eigenvectors of correlation matrix }
84:   Rc      : TMatrix; { Correlation factors/variables }
85:   Z       : TMatrix; { Scaled variables }
86:   F       : TMatrix; { Principal factors }
87:   I, J    : Integer; { Loop variables }
88:
89:
90:  PCA involves a lot of linear algebra, which we do not want to go into.
91:  Nevertheless, the basic algorithm can be easily described as follows:
92:
93:  1. Center the data by subtracting the mean from it.
94:  2. Calculate the covariance matrix.
95:  3. Calculate the eigenvectors of the covariance matrix.
96:
97:  If we start with N features, then the algorithm will return a transformed
98:    feature space again with N dimensions (we gained nothing so far). The nice
99:    thing about this algorithm, however, is that the eigenvalues indicate how
100:    much of the variance is described by the corresponding eigenvector.
101:
102:  Lets assume we start with N = 2000 features and we know that our model
103:    does not work well with more than 15 features. Then, we simply pick the 15
104:    eigenvectors with the highest eigenvalues.
105:
106:
107:  >>>
108:
109:  Mean vector:

```

```

103:
104: 138.027272727273
105: 17.7545454545455
106: 74.4454545454545
107: 1242.36363636364
108:
109: Variance-covariance matrix
110:
111:      1793.9074      -4.4733      726.8760      -2217.5917
112:      -4.4733       1.4879      -26.6170       197.4529
113:      726.8760      -26.6170     1224.4170      -6202.4802
114:     -2217.5917     197.4529     -6202.4802     48104.2314
115:
116: Correlation matrix
117:
118:      1.0000      -0.0866       0.4905      -0.2387
119:     -0.0866       1.0000      -0.6236       0.7380
120:       0.4905     -0.6236       1.0000      -0.8082
121:     -0.2387       0.7380     -0.8082       1.0000
122:
123:
124: Eigenvalues of correlation matrix:
125:
126: 2.59317860210044
127: 0.978288683253008
128: 0.285238674676994
129: 0.143294039969554
130:
131: Eigenvectors (columns) of correlation matrix:
132:
133:      0.2908      0.8713      0.3322      0.2142
134:     -0.5062      0.4248     -0.7423      0.1107
135:      0.5773      0.1360     -0.4184     -0.6879
136:     -0.5709      0.2047      0.4043     -0.6846
137:
138: Correlations between factors (columns) and variables (lines):
139:
140:      0.4682      0.8618      0.1774      0.0811
141:     -0.8152      0.4201     -0.3965      0.0419
142:      0.9296      0.1345     -0.2235     -0.2604
143:     -0.9194      0.2024      0.2159     -0.2592
144:
145: Principal factors:
146:
147:     -3.4114     -0.2835      0.1338      0.0511
148:      1.7022     -2.0748      0.4841      0.0733
149:     -1.2870      1.0193     -0.6491      0.2578
150:      0.4137     -0.2731      0.3081     -0.5820
151:     -0.0141     -0.8313     -1.1131     -0.2990
152:      0.9719      1.8876      0.3148     -0.3509
153:      0.0542     -0.5662     -0.3352      0.5091
154:     -1.8275     -0.2247      0.4421     -0.4907
155:     -0.1230      0.3023      0.2804      0.4942
156:      2.6719      0.5015     -0.4967     -0.1048
157:      0.8489      0.5427      0.6309      0.4420
158:
159:
160: By the way the well known dmath.dll is not statically linked:
161: Put the file dmath.dll therefore in your PATH and you can change to newer
    versions.
162: By default, DMath is intended to be used as a shared library, such as a
    Windows DLL. Compilation scripts are supplied in the dll subdirectory to
    generate the shared library file, the interface file and (in the case of
    FPC/Lazarus) the object file.

```

```

163: It is sometimes more efficient to link statically with individual units
    rather than calling the shared library. This may be done by replacing the
    reference to unit dmath by references to the individual units.
164:
165:   procedure VecMean(X           : TMatrix;
166:                     Lb, Ub, Nvar : Integer;
167:                     M           : TVector); external 'VecMean@dmath.dll';
168: { Computes the mean vector M from matrix X }
169:
170:   procedure MatVarCov(X           : TMatrix;
171:                      Lb, Ub, Nvar : Integer;
172:                      M           : TVector;
173:                      V           : TMatrix); external
'MatVarCov@dmath.dll';
174: { Computes the variance-covariance matrix V from matrix X }
175:
176:   procedure MatCorrel(V : TMatrix;
177:                      Nvar : Integer;
178:                      R   : TMatrix); external 'MatCorrel@dmath.dll';
179: { Computes the correlation matrix R from the var-cov matrix V }
180:
181:   procedure VecSD(X           : TMatrix;
182:                  Lb, Ub, Nvar : Integer;
183:                  M, S         : TVector); external 'VecSD@dmath.dll';
184: { Computes the vector of standard deviations S from matrix X }
185:
186:   procedure ScaleVar(X           : TMatrix;
187:                     Lb, Ub, Nvar : Integer;
188:                     M, S         : TVector;
189:                     Z           : TMatrix); external 'ScaleVar@dmath.dll';
190: { Scales a set of variables by subtracting means and dividing by SD's }
191:
192:   procedure PCA(R           : TMatrix;
193:                Nvar        : Integer;
194:                Lambda      : TVector;
195:                C, Rc       : TMatrix); external 'PCA@dmath.dll';
196: { Performs a principal component analysis of the correlation matrix R }
197:
198:   procedure PrinFac(Z           : TMatrix;
199:                    Lb, Ub, Nvar : Integer;
200:                    C, F         : TMatrix); external 'PrinFac@dmath.dll';
201: { Computes principal factors }
202:
203:
204: Of course, its not always this and that simple. Often, we dont know what
    number of dimensions is advisable in upfront. In that case, we leave
    n_components or Nvar parameter unspecified when initializing PCA to let it
    calculate the full transformation. After fitting the data,
    explained_variance_ratio_ contains an array of ratios in decreasing order:
    The first value is the ratio of the basis vector describing the direction
    of the highest variance, the second value is the ratio of the direction of
    the second highest variance, and so on.
205:
206: Being a linear method, PCA has, of course, its limitations when we are
    faced with strange data that has non-linear relationships. We wont go into
    much more details here, but its sufficient to say that there are
    extensions of PCA.
207:
208:
209: Ref:
210:   Building Machine Learning Systems with Python
211:   Second Edition March 2015

```

```
212:
213:     DMath
214:     Math library for Delphi, FreePascal and Lazarus May 14, 2011
215:
216:     http://fann.sourceforge.net
217:
218: Doc:
219:     Neural Networks Made Simple: Steffen Nissen
220:     http://fann.sourceforge.net/fann\_en.pdf
221:     http://www.softwareschule.ch/examples/datascience.txt
222:     https://maxbox4.wordpress.com
223:     https://www.tensorflow.org/
224:
225:
226: https://sourceforge.
    net/projects/maxbox/files/Examples/13\_General/811\_mXpcatest\_dmath\_datascien
    pas/download
227: https://sourceforge.
    net/projects/maxbox/files/Examples/13\_General/809\_FANN\_XorSample\_traindata.
    pas/download
228:
229:
230: Plots displaying the explained variance over the number of components is
    called a Scree plot. A nice example of combining a Screeplot with a grid
    search to find the best setting for the classification problem can be
    found at
231:
232: http://scikit-
    learn.sourceforge.net/stable/auto\_examples/plot\_digits\_pipe.html.
233:
234: Although, PCA tries to use optimization for retained variance,
    multidimensional scaling (MDS) tries to retain the relative distances as
    much as possible when reducing the dimensions. This is useful when we have
    a high-dimensional dataset and want to get a visual impression.
235: >>> Building Machine Learning Systems with Python
```