

Prof. R. Rojas

Mustererkennung, WS17/18

Übungsblatt 8

Boyan Hristov, Nedeltscho Petrov

10. Dezember 2017

Link zum Git Repository: <https://github.com/BoyanH/FU-MachineLearning-17-18/tree/master/Solutions/Homework8>

Principal Component Analysis

Analyse anhand Klassifizierung von Ziffern mit Hilfe von PCA

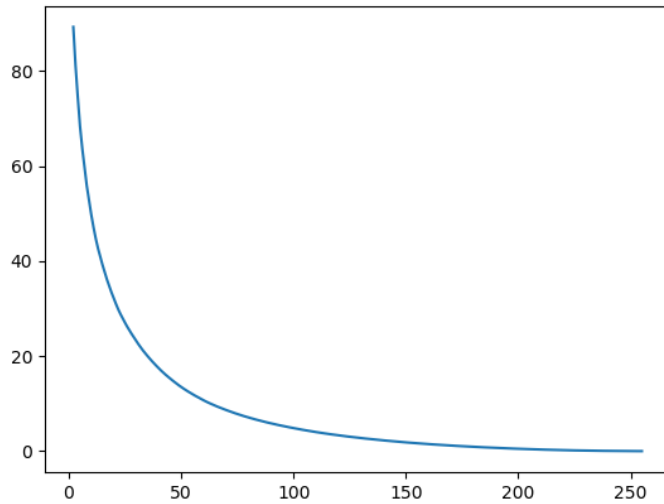
Wir haben unsere PCA Implementierung für die Klassifizierung von Digits angewandt. Wir haben die folgende Ergebnisse bekommen:

```
1 Score LDA without PCA: 0.9168458781362007
2 Score LDA with PCA, k=230: 0.9186379928315412
3 Score LDA with PCA, k=30: 0.9010752688172043
```

Wir wollten mal schauen, was passiert, wenn man einen nicht binären Klassifikator anwendet zusammen mit PCA. Die Ergebnisse waren ziemlich gut. Natürlich bekommt man nicht viel bessere Ergebnisse mit PCA, aber wenn man auch mit etwas schlechtere zufrieden ist, kann man das notwendige Speicherplatz deutlich reduzieren (in unserem Fall um Faktor 8.5).

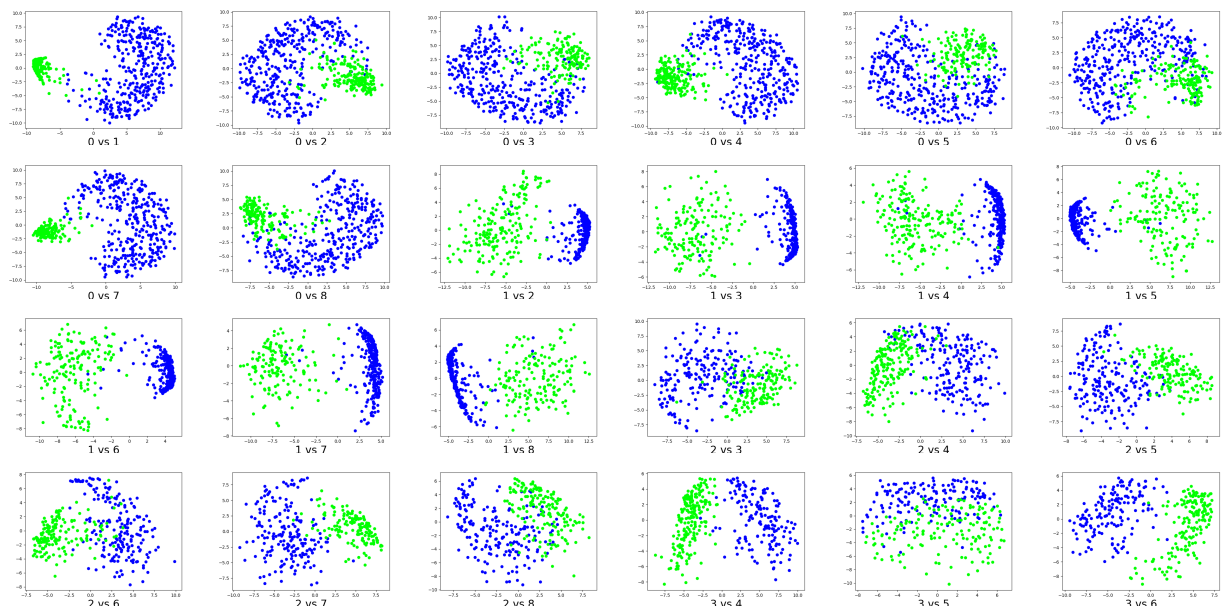
Dabei haben wir LDA von scikit-learn genommen. Wie man sieht, es gibt einige korrelierende Komponenten. Wenn man diese entfernt, kriegt man etwas bessere Ergebnisse. Wichtiger ist, dass man auch ziemlich gute Ergebnisse mit nur 30 Komponenten bekommen kann.

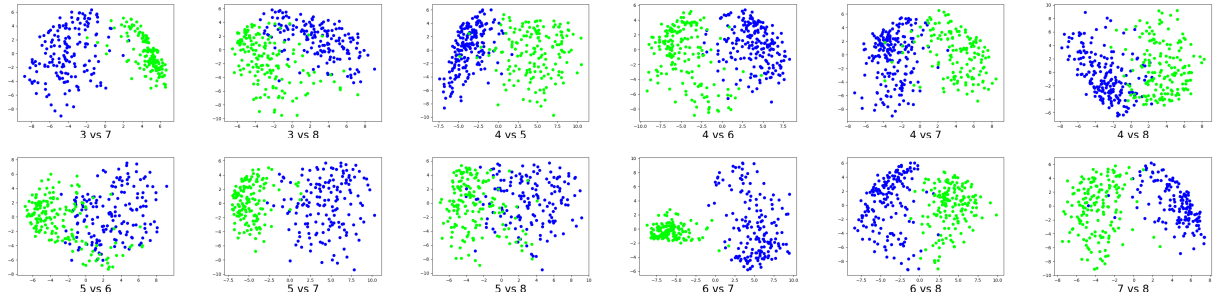
Wie viel Varianz man verliert, wenn man den Datensatz in k Dimensionen reduziert (k auf der x-Achse)



Auf dem Plot ist kein echten Elfbogenbu sehen oder erst bei $k=80$, wobei das initiale Datensatz nur 256 Dimensionen hat.

2D Visualisierung von den einzelnen Cluster bei dem Zifferklassifizierungsproblem



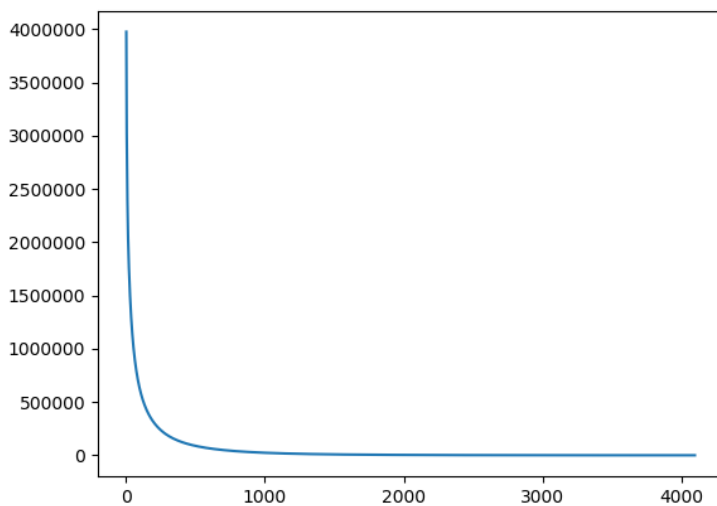


Wie man sieht, sind einige Ziffer linear von einander separierbar (z.B. 1 und 3) und andere gar schwer von einander unterscheidbar (in 2D Raum, 2 und 8). Man sieht, dass es auch einige gibt, die in den meisten Fällen gut unterscheidbar sind, aber haben auch einige Exemplare, die ganz ähnlich sind (z. B. 3 und 5, wenn man das obere Teil von je Ziffer nicht gut geschrieben hat ist das untere ja gleich).

Hier wurde aus dem ganzen Datensatz nur das jeweilige Paar genommen und die Dimensionen davon wurden reduziert. Wenn man aber PCA auf dem ganzen Datensatz anwendet, kriegt man nicht so gute Ergebnisse. Deswegen wurden wir vermuten, das man mit deutlich bessere Ergebnisse bekommen würde, wenn man ein binäres Klassifikator benutzt und damit zusammen PCA anwendet.

Eigenfaces

Wir haben PCA auf dem Datensatz von menschlichen Gesichten angewandt. Die dabei berechnete Hauptkomponenten (die Eigenvektoren der Kovarianzmatrix) nennt man „Eigenfaces“. Wie man auf dem nächsten Plot sieht, ist das „Elfbogen“, also das Punkt, ab dem man nicht viel Varianz mehr gewinnt bei weiteren Hauptkomponenten, erst bei Circa 320 Dimensionen. Die verlorene Varianz hört aber schon bei circa 30 echt drastisch zu sinken, deswegen haben wir 30 Hauptkomponenten genommen.



Hauptkomponenten des Gesichtsdatensatzes (a.k.a Eigenfaces)



Erklärung der PCA Implementierung

Fit Methode

Diese Methode berechnet unsere Hauptkomponenten und damit auch die Transformationsmatrix, mit der man ein Datensatz von n auf k Dimensionen reduziert. Hier wurde PCA genau so wie in der Vorlesung auf eine hohe Abstraktionsebene implementiert, interessanter sind also eigentlich die Hilfsfunktionen.

```
1 def fit(self, X):
2     sorted_k_eig_vectors = PCA.get_sorted_eig_vec(X, self.k)
3     self.principal_components = sorted_k_eig_vectors
4     self.transformation_matrix = sorted_k_eig_vectors.T
```

Hilfsfunktionen

Wir berechnen die Kovarianzmatrix mit `np.cov`. Hier ist wichtig entweder das transponierte Datensatz zu benutzen, oder `rowvar` auf `False` zu setzen, damit man die Verhältnisse der Features im Datensatz und nicht der Samples.

Nachher können wir mit Hilfe von `np.linalg.eigh` die Eigenvektoren und ihrer entsprechende Eigenwerte berechnen. Hier ist ganz wichtig, dass die Eigenvektoren die Spalten der zurückgegebener Matrix sind. Das haben wir initial nicht verstanden und ziemlich viel Zeit deswegen beim Debugging verbracht.

Danach kann man leicht mit `argsort()` die Indizes bekommen, die die Eigenwerte sortieren und damit die Eigenvektoren sortieren und zurückgeben. Da die Eigenvektoren die Spalten sind, muss man hier aufpassen.

Die sortierte Eigenwerte brauchen wir für die Berechnung der verlorenen Kovarianz bei k Dimensionen.

```

1 @staticmethod
2     def get_eig_values_sort_args(eig_values, k=None):
3         if k is None:
4             k = len(eig_values)
5
6         # reverse sorted, first k components
7         return eig_values.argsort()[::-1][:k]
8
9     @staticmethod
10    def get_eig_vec_and_val(X):
11        covariance_matrix = np.cov(X, rowvar=False)
12        # the covariance matrix is always symmetric, we don't need to bother any further
13        # therefore, we can also use eigh instead of eig
14        return np.linalg.eigh(covariance_matrix)
15
16    @staticmethod
17    def get_sorted_eig_vec(X, k=None):
18        eig_values, eig_vectors = PCA.get_eig_vec_and_val(X)
19        sort_args = PCA.get_eig_values_sort_args(eig_values, k)
20
21        return eig_vectors[:, sort_args].T
22
23    @staticmethod
24    def get_sorted_eig_values(X, k=None):
25        eig_values, eig_vectors = PCA.get_eig_vec_and_val(X)
26        # assert(np.all(eig_values >= 0)) # yep, all good
27        sort_args = PCA.get_eig_values_sort_args(eig_values, k)
28
29        return eig_values[sort_args]

```

Plotting der verlorenen Varianz bei k Dimensionen

Die Differenz von der Summe aller Eigenwerte mit der Summe von den benutzten (diese, der benutzten Hauptkomponenten) ist eine gute Möglichkeit zu plotten wie viel von der Varianz verloren wurde bei k Dimension. Ein Idealwert wäre also 0, wir wollen aber sehen, ab wie viele Dimensionen der Wert deutlich langsamer senkt.

```

1 @staticmethod
2     def plot_variance_for_k(X, save_plot_name=None):
3         sorted_eig_values = PCA.get_sorted_eig_values(X)
4         total_variance = sorted_eig_values.sum()
5         ks = np.arange(2, len(X[0]), 1)
6
7         variance_diffs = np.vectorize(lambda k: abs(total_variance - sorted_eig_values[:k].
8 sum()))(ks)
9         plt.plot(ks, variance_diffs)
10
11        if save_plot_name is not None:
12            plt.savefig(save_plot_name + '.png')
13        else:
14            plt.show()

```

Transform bzw. Fit-Transform Methoden

Wie unten erklärt, wenn man je x aus dem Datensatz mit je Hauptkomponente multipliziert bekommt man x in den reduzierten Raum.

Bsp: $x_i * (e_1, e_2, \dots, e_k) = x_i$ repräsentiert durch die k Hauptkomponenten

Wir speichern bei der Fit-Methode die Transformationsmatritze und nutzen diese bei Transform. So können wir dan weitere Daten erstmal transformieren, bevor wir Vorhersagen machen.

```

1  def transform(self, X):
2      '''
3      (x11 x12 ... x1n )      (e11 e12 ... e1n)T      (x1 projected in subspace)
4      (... .. ... .. )      X      (... .. ... .. )      =      (... ..)
5      (xm1 xm2 ... xmn )      (ek1 ek2 ... ekn)      (xm projected in subspace)
6
7      self.transformation_matrix is the already transposed matrix where each row is an
      eigenvector
8
9      :param X: data set, each row represents a sample
10     :return: X in the space defined by the first k eigenvectors of the fit data set
11     '''
12
13     return X.dot(self.transformation_matrix)
14
15     def fit_transform(self, X):
16         self.fit(X)
17         return self.transform(X)

```

Vollständiges Code (auch für Plotting, Eigenfaces, Helpers usw.)

PCA.py

```

1  from Classifier import Classifier
2  import numpy as np
3  from matplotlib import pyplot as plt
4
5
6  class PCA(Classifier):
7      def __init__(self, k):
8          self.transformation_matrix = None
9          self.k = k
10         self.X_mean = None
11         self.principal_components = None
12
13     @staticmethod
14     def get_eig_values_sort_args(eig_values, k=None):
15         if k is None:
16             k = len(eig_values)
17
18         # reverse sorted, first k components
19         return eig_values.argsort()[::-1][:k]
20
21     @staticmethod
22     def get_eig_vec_and_val(X):
23         covariance_matrix = np.cov(X, rowvar=False)
24         # the covariance matrix is always symmetric, we don't need to bother any further
25         # therefore, we can also use eigh instead of eig
26         return np.linalg.eigh(covariance_matrix)
27
28     @staticmethod
29     def get_sorted_eig_vec(X, k=None):
30         eig_values, eig_vectors = PCA.get_eig_vec_and_val(X)
31         sort_args = PCA.get_eig_values_sort_args(eig_values, k)
32
33         return eig_vectors[:,sort_args].T
34
35     @staticmethod
36     def get_sorted_eig_values(X, k=None):

```

```

37     eig_values, eig_vectors = PCA.get_eig_vec_and_val(X)
38     # assert(np.all(eig_values >= 0)) # yep, all good
39     sort_args = PCA.get_eig_values_sort_args(eig_values, k)
40
41     return eig_values[sort_args]
42
43     def fit(self, X):
44         """
45         Finds the largest k eigenvalues and their corresponding eigenvectors from the
46         covariance matrix of the dataset and saves the transformation matrix
47         which can be used to project a data set from it's original space to
48         the space defined by those eigenvectors and their corresponding eigenvalues
49         """
50
51         sorted_k_eig_vectors = PCA.get_sorted_eig_vec(X, self.k)
52         self.principal_components = sorted_k_eig_vectors
53         self.transformation_matrix = sorted_k_eig_vectors.T
54
55     def transform(self, X):
56         """
57         (x11 x12 ... x1n )      (e11 e12 ... e1n)T      (x1 projected in subspace)
58         (... ... ... ... )      X      (... ... ... ...) = (... ...)
59         (xm1 xm2 ... xmn )      (ek1 ek2 ... ekn)      (xm projected in subspace)
60
61         self.transformation_matrix is the already transposed matrix where each row is an
        eigenvector
62
63         :param X: data set, each row represents a sample
64         :return: X in the space defined by the first k eigenvectors of the fit data set
65         """
66
67         return X.dot(self.transformation_matrix)
68
69     def fit_transform(self, X):
70         self.fit(X)
71         return self.transform(X)
72
73     @staticmethod
74     def plot_variance_for_k(X, save_plot_name=None):
75         sorted_eig_values = PCA.get_sorted_eig_values(X)
76         total_variance = sorted_eig_values.sum()
77         ks = np.arange(2, len(X[0]), 1)
78
79         variance_diffs = np.vectorize(lambda k: abs(total_variance - sorted_eig_values[:k].
80 sum()))(ks)
81
82         plt.plot(ks, variance_diffs)
83
84         if save_plot_name is not None:
85             plt.savefig(save_plot_name + '.png')
86         else:
87             plt.show()

```

Helpers.py

```

1 import numpy as np
2 import math
3 import os
4 from matplotlib import pyplot as plt
5
6
7 def plot_multiple_images(title, images, data_set_dimensions, cols=5, fig_name=None):
8     rows = math.ceil(len(images) / cols)
9     plt.figure(figsize=(8,8))
10    plt.suptitle(title, size=20)
11
12    for i, component in enumerate(images):
13        plt.subplot(rows, cols, i + 1)

```

```

14     component += abs(component.min())
15     component *= (1.0 / component.max())
16     image = np.reshape(component, (data_set_dimensions, data_set_dimensions))
17     plt.imshow(image, cmap=plt.cm.binary)
18     plt.xticks(())
19     plt.yticks(())
21
22     if fig_name is not None:
23         file_name = os.path.join(os.path.dirname(__file__), './{0}'.format(fig_name))
24         plt.savefig(file_name)
25     else:
26         plt.show()

```

Parser.py

```

1 import csv
2 import numpy as np
3 import pandas as pd
4 import os
5 from sklearn.model_selection import train_test_split
6
7
8 def parse_data(file_name, sep=' '):
9     file_name = os.path.join(os.path.dirname(__file__), './Dataset/{0}'.format(file_name))
10    return pd.read_csv(file_name, header=None, na_filter=True, sep=sep).as_matrix()
11
12
13 def get_points_and_labels_from_data(data, label_idx=0):
14     points = (np.array(data[:, 1:], dtype=np.float64) if label_idx == 0 else
15              np.array(data[:, :label_idx], dtype=np.float64))
16     labels = np.array(data[:, label_idx])
17
18     return points, labels
19
20
21 def get_data_set(file_name, sep=' ', label_idx=0, seed=4):
22     data = parse_data(file_name, sep)
23     X, y = get_points_and_labels_from_data(data, label_idx)
24     # for determined results we use a seed for random_state, so that data is always split
25     X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, test_size
26                                                         =0.3,
27                                                         random_state=seed)
28
29     return X_train, X_test, y_train, y_test
30
31
32 def extract_classes_from_data_set(X, y, classes):
33     labels_are_numbers = False
34
35     try:
36         a = int(y[0])
37         labels_are_numbers = True
38     except:
39         pass
40
41     is_from_classes = np.vectorize(lambda y: (int(y) if labels_are_numbers else y) in
42                                     classes)
43     filter_arr = is_from_classes(y)
44     return X[filter_arr], y[filter_arr]

```

digits_classification.py

```

1 import numpy as np
2 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
3 from PCA import PCA
4 from Parser import get_data_set

```



```

6 X_train, X_test, y_train, y_test = get_data_set('digits.data')

8 pca_230 = PCA(230)
9 train_transformed_230 = pca_230.fit_transform(X_train)
10 test_transformed_230 = pca_230.transform(X_test)

12 pca_30 = PCA(30)
13 train_transformed_30 = pca_30.fit_transform(X_train)
14 test_transformed_30 = pca_30.transform(X_test)

16 PCA.plot_variance_for_k(X_train)

18 prediction_non_transformed = LDA().fit(X_train, y_train).predict(X_test)
19 score_non_transformed = np.mean(prediction_non_transformed == y_test)
20 print('Score LDA without PCA: {}'.format(score_non_transformed))

22 prediction_transformed_230 = LDA().fit(train_transformed_230, y_train).predict(
    test_transformed_230)
23 score_transformed_230 = np.mean(prediction_transformed_230 == y_test)
24 print('Score LDA with PCA, k={}: {}'.format(230, score_transformed_230))

26 prediction_transformed_30 = LDA().fit(train_transformed_30, y_train).predict(
    test_transformed_30)
27 score_transformed_30 = np.mean(prediction_transformed_30 == y_test)
28 print('Score LDA with PCA, k={}: {}'.format(30, score_transformed_30))

```

demo_digits_PCA.py

```

1 from matplotlib import pyplot as plt
2 import numpy as np
3 from sklearn.decomposition import PCA as PCA_sklearn
4 from PCA import PCA
5 from Parser import parse_data, get_points_and_labels_from_data,
    extract_classes_from_data_set

7 data = parse_data('digits_test.data')
8 X, y = get_points_and_labels_from_data(data)

10 for a in range(10):
11     for b in range(a+1, 10, 1):
12         X_extracted, y_extracted = extract_classes_from_data_set(X, y, [a,b])
13         X_transformed = PCA(2).fit_transform(X_extracted)
14         X_a_idx = y_extracted.astype(int) == a
15         X_a = X_transformed[X_a_idx]
16         X_b = X_transformed[np.invert(X_a_idx)]

18         plt.scatter(X_a[:,0], X_a[:,1], c='#0000FF')
19         plt.scatter(X_b[:,0], X_b[:,1], c='#00FF00')
20         plt.xlabel('{} vs {}'.format(a, b), fontsize=24)
21         plt.savefig('figs/{}vs{}.png'.format(a,b))
22         plt.clf()

```

eigenfaces.py

```

1 from PCA import PCA
2 from Helpers import plot_multiple_images
3 import numpy as np
4 import math
5 import glob
6 import os
7 import cv2

10 def read_pgm(file_name):
11     return cv2.imread(file_name, -1)

13 data_folder = os.path.abspath(os.path.join(os.path.dirname(__file__), './Dataset'))

```

```

14 files = glob.glob(data_folder + '/lfwcrop_grey/faces/*.pgm')
15 data_set = []

17 for file in files:
18     image = read_pgm(file)
19     data_set.append(image.flatten())

21 data_set = np.array(data_set)

23 # drops quickly until about k=320
24 # but really quickly until about 30, plus we can't really submit 320 eigenfaces for this
   homework...
25 # PCA.plot_variance_for_k(data_set, save_plot_name='eigenfaces_variance_for_k')

28 pca_eigenfaces = PCA(30)
29 pca_eigenfaces.fit(data_set)
30 data_set_dimensions = int(math.sqrt(len(data_set[0])))
31 principal_components = pca_eigenfaces.principal_components

33 plot_multiple_images('Eigenfaces', principal_components, data_set_dimensions, fig_name='
   eigenfaces.png')

```