

CS 5785: APPLIED MACHINE LEARNING

HOMEWORK 2

September 27, 2017

Sarah Le Cam - sdl83

Yunie Mao - ym224

Cornell Tech

Contents

Eigenface for Face Recognition	3
Summary	3
Data	3
Procedure & Insights	4
Question 1 (a)	4
Question 1 (b)	4
Question 1 (c)	4
Question 1 (d)	5
Question 1 (e)	6
Question 1 (f)	7
Question 1 (g)	7
Question 1 (h)	7
What's Cooking?	9
Summary	9
Data	9
Procedure & Insights	10
Question 1 (a)	10
Question 1 (b)	10
Question 1 (c)	10
Question 1 (d)	10
Question 1 (e)	10
Question 1 (f)	11
Question 1 (g)	11
Written Exercises	12
Question 1	12
Question 2	12
2. a.	12

2. b.	12
2. c.	12
2. d.	12
2. e.	12
Question 3	12
2. a.	12
2. b.	12
2. c.	13
2. d.	13
2. e.	13
Sources & External libraries	14

EIGENFACE FOR FACE RECOGNITION

Summary

We were given a set of black and white pictures of faces with training and testing data files containing the links to those images and corresponding labels identifying the individuals. Using this data, we calculated the mean image and subtracted it from each of the images in our training set. We then performed a Singular-Value Decomposition to find the set of Eigenfaces. Using our Eigenfaces, we computed the Eigenfeatures and the ranked r -dimensional feature vectors for both the training images and test images. Finally, we fitted the Eigenfeatures and labels of our training set to a logistic regression using *scikit-learn*'s logistic regression model. We used this model to find predicted labels for our test data using the Eigenfeatures of the test images and calculated the mean accuracy on the given test data and labels. To visualize the fit of our model, we plotted our classification's accuracy for the first 200 dimensions in the face space of our test data.

Data

We were provided with a set of 640 pictures total of 10 distinct subjects and two files - a testing and a training text file - matching each image to its respective label. The training set contained the image links and label pairings for 540 images and the testing set contained 100. Each image is 50 x 50 pixels black and white photograph. In order to use this data for model fitting, we converted the images into grayscale and stored the pixel data in a matrix.

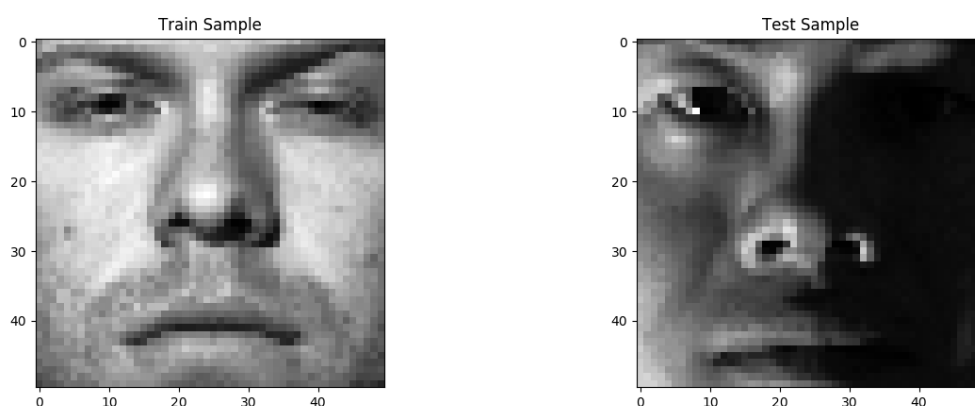
Procedure & Insights

Question 1 (a)

We downloaded and unzipped the faces data file. We then used Anaconda Navigator's Spyder IDE to create a Python project and included our images folder (*faces/images*) and our training and testing data files (*faces/train.txt* & *faces/test.txt*). We generated a Python file (*eigenFaces.py*) and imported the relevant external libraries (NumPy, SciPy, Matplotlib and sklearn).

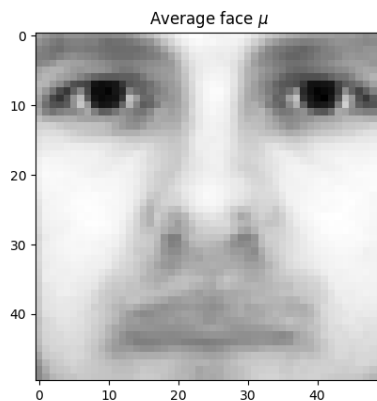
Question 1 (b)

We retrieved each image link from the training and testing datasets using the *split()* function. We then computed the images' grayscale pixel information using Matplotlib's *imread()* and stored the pixel configurations of the training and testing images in two matrices of size, respectively, 540 x 2500 and 100 x 2500. We then displayed a sample image (the 10th in each dataset) using the pixel information stored in each of these matrices using *imshow()*. We saved these the sample image for the training set as *training_image.png* and that for the testing set as *test_sample.png*. For both the training and testing datasets, we also extracted the labels from the text files into 2 flat arrays of size 540 (training labels) and 100 (testing labels) using the *split()* function.

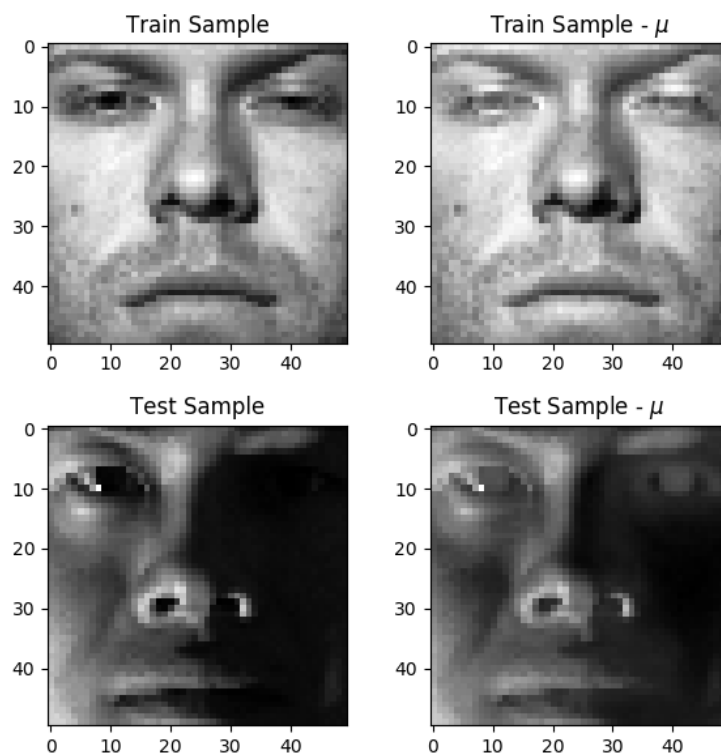


Question 1 (c)

Using the NumPy's *mean()* function along the vertical axis of the training dataset pixel matrix, we found the average face μ and displayed it using the Matplotlib *imshow()* function. We saved the image as *average_image.png*.

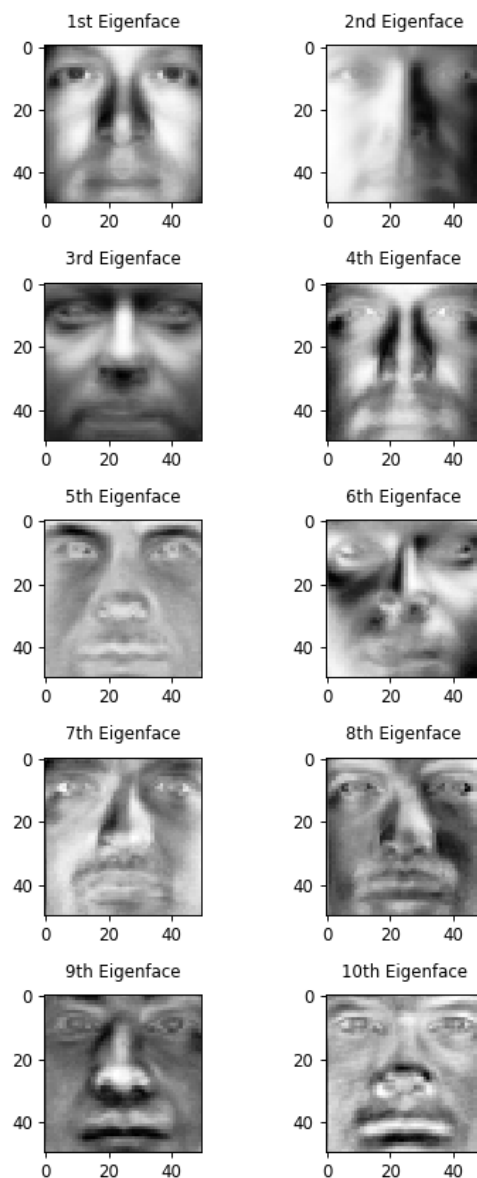
**Question 1 (d)**

We then subtracted our average face μ 's pixel values from those of every image in the training and testing matrices to form new adjusted matrices. These new matrices indicate distance from the mean, allowing us to centralize our data. Again, we displayed a sample (the 10th in each matrix) from each of the new adjusted datasets (see *original_and_adjusted_images.png*).



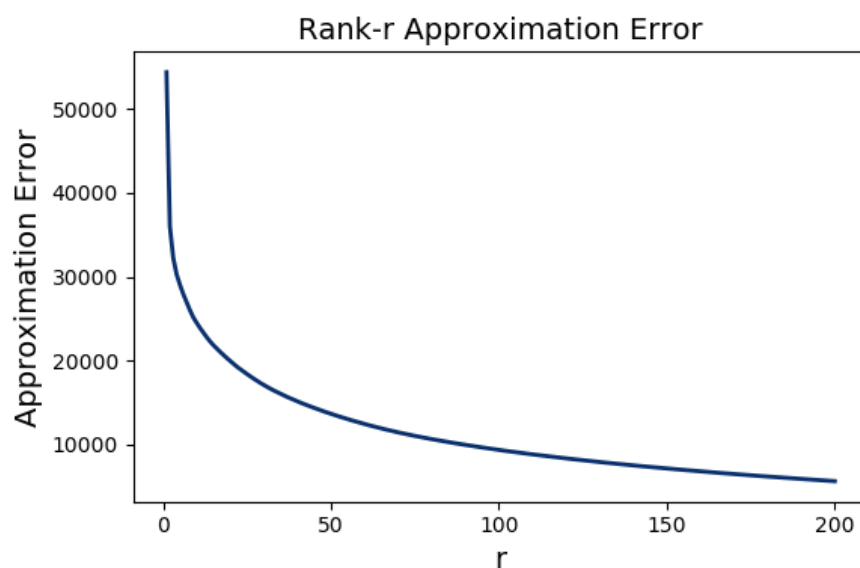
Question 1 (e)

We performed a Singular Value Decomposition (SVD): $X = U\Sigma V^T$ where X is the matrix representation of the adjusted training set. Using NumPy's `linalg.svd()` function, we computed U (the left-singular vector matrix), Σ (the covariance matrix), and V^T (the transpose of the left-singular vectors of X). We then displayed the top ten Eigenfaces from V^T as images in grayscale using `imshow()` (see `first_ten_eigenfaces.png`).



Question 1 (f)

We generated a helper function to compute the rank- r approximation of our adjusted training data by taking the first r columns of U , the first r elements of Σ and the first r rows of V^T . We then computed the low-rank approximation error of our adjusted training data to the rank- r approximation for $r = 1, 2, \dots, 200$ and plotted the results as a function of the value of r (see *low_rank_approximation_err.png*). As the plot shows, as r increases the approximation error decreases exponentially. A value of only 200 for r corresponds to a relatively low approximation error.

**Question 1 (g)**

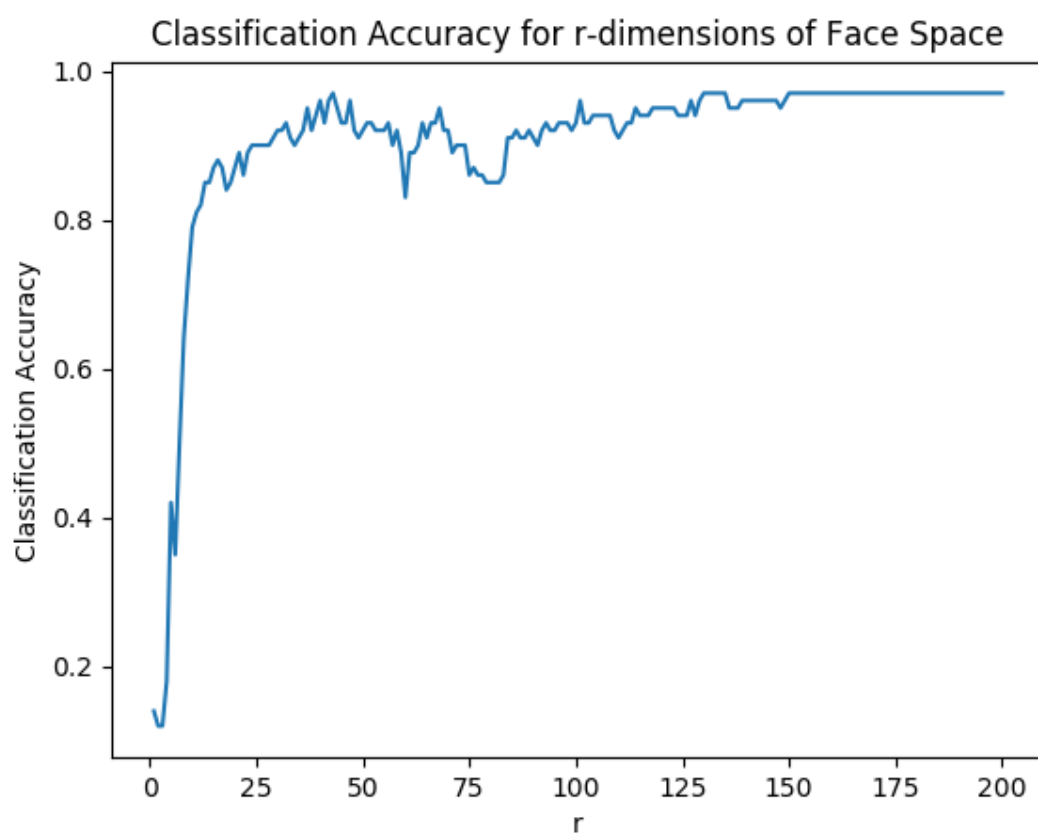
Since the first r Eigenfaces span the r -dimensional subspace of the original image space (*face space*), we can represent a 2500-dimensional face image as an r -dimensional feature vector and thereby reduce the dimensions prior to classification. To compute the r -dimensional feature matrices for the training and test images, we multiplied these images by the transpose of the first r rows of V^T .

Question 1 (h)

Using the function we generated in 1(g), we extracted the Eigenfeatures for our training and test data for $r = 10$. We then fitted our training Eigenfeatures into a logistic regression model provided by scikit-learn and used the model to generate predicted labels for the test data. We

then computed the classification accuracy rate on the test data given the Eigenfeatures and labels. We achieved a classification accuracy of 79%.

To show the classification rate on the test data as a function of r , we generated the Eigenfeatures, trained a logistic regression model, and computed the accuracy rate on our test set for the first 200 values of r . The following plot (*face_recognition_classification_accuracy.png*) shows the classification accuracy for varying r -dimensions. We can see that for values of $r > 30$, we achieve a classification accuracy of 90%.



WHAT'S COOKING?

Summary

We were given descriptive information of many different dishes for training and testing purposes.

[Our goal was to find the best possible method for classifying recipes by cuisine when given their respective ingredients].

We first transformed the data into a usable numeric matrix, then

[TODO]

Data

The Kaggle competition provided us with training and testing json files containing information describing recipes. Both datasets contained recipe identifiers and ingredients lists. The training dataset included an additional cuisine field. The training file included 39,774 dishes with 20 categories and 6,714 unique ingredients. The testing file included 9,944 dishes with 4,484 unique ingredients.

Procedure & Insights

Question 2 (a)

We joined the Kaggle "What's Cooking?" competition and downloaded the training (*train.json*) and testing (*test.json*) data files. We then used Anaconda Navigator's Spyder IDE to create a Python project and included these files. We generated a Python file (*cooking.py*) and imported the relevant external libraries (NumPy, Pandas, IterTools, Matplotlib and sklearn).

Question 2 (b)

We used Pandas' DataFrame to import the json data and find the number of distinct dishes and cuisines. The training file includes 39,774 dishes spanning 20 categories. We then used an encoder to extract the information from the lists contained in each data object. There 6,714 unique ingredients included in total in the training set.

Question 2 (c)

To set up our training set for classification, we generated an $n \times d$ matrix, where n is the number of dish samples and d is the number of unique ingredients for both the training and testing datasets. We represented each dish as a binary ingredient vector x , where $x_i = 1$ if the dish contains ingredient i and $x_i = 0$ otherwise. To do this, we used scikit-learn's CountVectorizer function to generate a map of the ingredient to its frequency for each dish. We then fitted and transformed the ingredients and the labels into numerical vectors based on the non-numerical values.

[TODO]

Question 2 (d)

Using scikit-learn's Naive Bayes Classifier, we performed a 3 fold cross-validation on the training data with the Bernoulli distribution prior and the Gaussian distribution prior assumptions. We achieved an average accuracy rate of 68.2% using the Bernoulli Naive Bayes Classifier and 36.9% using the Gaussian Naive Bayes Classifier. [TODO]

Question 2 (e)

The Bernoulli Naive Bayes Classifier performed much better than the Gaussian Naive Bayes Classifier. This makes sense because we represented each dish as a binary ingredient vector x in our training data. The Bernoulli Naive Bayes best fits our assumptions because it describes

whether or not an ingredient was found in a dish while the Gaussian Naive Bayes Classifier generally describes continuous data that are normally distributed.

[TODO]

Question 2 (f)

We performed a 3 fold cross-validation on the training data using scikit-learn's Logistic Regression model. We achieved an average classification accuracy of 77.5%.

[TODO]

Question 2 (g)

Based on the results from 1(d) and 1(f), the Logistic Regression model had the best accuracy performance over our training data. Using this model, we classified the cuisines based on the dish ingredients in our test data. We generated a csv file that contained the list of dish ids and their corresponding predicted cuisine labels. After submitting our results to Kaggle, we received an accuracy rate of 78.177 %.

Your most recent submission				
Name	Submitted	Wait time	Execution time	Score
test_cooking_predictions.csv	2 minutes ago	0 seconds	0 seconds	0.78177
Complete				
Jump to your position on the leaderboard ▼				

[TODO]

WRITTEN EXERCISES

Question 1

[TODO]

Question 2

2. a.

[TODO]

2. b.

[TODO]

2. c.

[TODO]

2. d.

[TODO]

2. e.

[TODO]

Question 3

3. a.

[TODO]

3. b.

[TODO]

3. c.

[TODO]

3. d.

[TODO]

3. e.

[TODO]

SOURCES & EXTERNAL LIBRARIES

Stéfan van der Walt, S. Chris Colbert and Gaël Varoquaux. *The NumPy Array: A Structure for Efficient Numerical Computation*, Computing in Science & Engineering, 13, 22-30 (2011), DOI:10.1109/MCSE.2011.37

John D. Hunter. *Matplotlib: A 2D Graphics Environment*, Computing in Science & Engineering, 9, 90-95 (2007), DOI:10.1109/MCSE.2007.55

Jones E, Oliphant E, Peterson P, et al. *SciPy: Open Source Scientific Tools for Python*, 2001-, <http://www.scipy.org/>

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, Édouard Duchesnay. *Scikit-learn: Machine Learning in Python*, Journal of Machine Learning Research, 12, 2825-2830 (2011)

Wes McKinney. *Data Structures for Statistical Computing in Python*, Proceedings of the 9th Python in Science Conference, 51-56 (2010)

What's Cooking? | Kaggle, www.kaggle.com/c/whats-cooking.