

**Math 6350    fall 2019 MSDS    Robert Azencott**

## **Homework 2**

**Due date    Thursday sept 26 at Midnight**

### **Data Set Information:**

**Go to the site of**

" university of california irvine repository of machine learning datasets "

which contains more than 400 data sets for machine learning

The data set to be used for HW2 consists of digitized images of typed characters

the font choices for these characters is quite large (153 fonts).

**Download the   fonts.zip   file** from the following link

<https://archive.ics.uci.edu/ml/machine-learning-databases/00417/>

fonts.zip contains 153 files in format .csv (= comma delimited files), **one file for each font**

for instance AGENCY.csv, ..., TIMES.csv, etc

Each such XYZ.csv file has a "header row" (= top row) which lists the names of each column, ("attribute names")

EXTRACT from font.zip the three font files

**COURIER.csv, CALIBRI.csv, TIMES.csv**

### **Description of a typical   table such as   TIMES.csv   :**

the font type "TIMES" is the same throughout   the   file of TIMES.csv

after the header row , each row of the TIMES.csv file describes one single "example"

there are 12530 such rows = 12530 examples of characters typed in the font TIMES

the number of examples of characters will be **different for each font type**

each "example" corresponds to a digitized image of some specific character (letter or digit essentially) typed in the font TIMES

each image has size 20 x20 pixels

each one of these 400 pixels has a "gray level" which is an integer value between 0 and 255

in each row, the image associated to this row has 400 features (or numerical descriptors) whose values are listed in the 400 columns named r0c0, r0c1, r0c2, ... , r19c18, r19c19

these 400 values represent the gray levels of the 400 pixels of this specific image

"rLcK"= gray level image intensity for pixel in position{Row L, Column K}

the **m\_label** column is an arbitrary integer from 33 to 65535 ; in each row this m-label integer is the standardized ID of the character present in the image associated to this row;

the **strength** column lists values either equal to 0.4 or to 0.7; in each row , strength =0.4 for NORMAL character ; strength= 0.7 for BOLD character;

the **italic** column lists values either equal to 0 or to 1; in each row , italic= 1 for ITALIC character ; italic= 0 for NORMAL character ;

### **preliminary treatment of the data set**

in each font file, the header row begins with 12 names

font, fontVariant, m\_label, strength, italic, orientation, m\_top, m\_left, originalH, originalW, h, w

DISCARD the 9 columns listed below :

fontVariant, m\_label, orientation, m\_top, m\_left, originalH, originalW, h, w

KEEP the 3 columns {font, strength, italic}

KEEP the 400 columns named r0c0, r0c1, r0c2, ... , r19c18, r19c19

any row containing missing numerical data will be discarded

define then three CLASSES of images of "normal" characters as follows

CL1 = all rows of COURIER.csv file for which {row # >1 and strength = 0.4 and italic=0}

CL2 = all rows of CALIBRI.csv file for which {row # >1 and strength = 0.4 and italic=0}

CL3 = all rows of TIME.csv file for which {row # >1 and strength = 0.4 and italic=0}

Display their respective sizes  $n_1, n_2, n_3$

The full data set (denoted DATA ) for the next questions will be the union of the three classes CL1 , CL2, CL3 and hence has size  $N = n_1 + n_2 + n_3$

Example #  $i$  in DATA corresponds to a specific row " $i$ " in the matrix DATA , and will be described by this vector of 400 **features** , namely the 400 numbers listed in row " $i$ "

Recall that these 400 feature values observed for example # $i$  are also viewed as the 400 values taken for this particular example by the random variables

$X_1 = r_{0c0}, X_2 = r_{0c1}, X_3 = r_{0c2}, \dots, X_{399} = r_{19c18}, X_{400} = r_{19c19}$

Each such random variable  $X_j$  is observed  $N$  times, and its  $N$  observed values are listed in the column " $j$ " of DATA.

We will use the basic machine learning tools covered in class so far to attempt a rough automatic classification of DATA into 3 classes CL1 CL2 CL3

## PART 0

Compute the means  $m_1 = \text{mean}(X_1) \dots \text{mean}(X_{400}) = m_{400}$  and the standard deviations

$s_1 = \text{std}(X_1) \dots s_{400} = \text{std}(X_{400})$

Standardize the features matrix DATA by centering and rescaling each random variable  $X_j$  into a new random variable  $Y_j = (X_j - m_j) / s_j$ ; the matrix DATA becomes a standardized data matrix SDATA, with coefficients given by

$\text{SDATA}(i,j) = (\text{DATA}(i,j) - m_j) / s_j$

SDATA is also called a rescaled data matrix for short. The example #  $i$  associated to row " $i$ " of DATA will be from now on described by the vector " $E_i$ " of standardized features defined by row " $i$ " of SDATA

## PART 1

1.1) Compute the correlation matrix COR of the 400 random variables  $Y_1, \dots, Y_{400}$

1.2) For the matrix COR , compute its 400 eigenvalues  $\lambda_1 > \lambda_2 > \dots > \lambda_{400} > 0$  , and its

400 eigenvectors  $v_1, v_2, \dots, v_{400}$

1.3) Plot the decreasing curve  $\lambda_j$  versus  $j$  for  $j=1, 2, \dots, 400$

1.4) for  $j=1, 2, \dots, 400$  compute the successive percentages  $R_j$  given by

$R_j = (\lambda_1 + \lambda_2 + \dots + \lambda_j) / 400$

1.5) Plot the increasing curve  $R_j$  versus  $j$  for  $j=1, 2, \dots, 400$  and compute the smallest

integer " $r$ " such that  $R_r > 90\%$

1.6) Explain the relationship between these computations and the PCA analysis of the set DATA

1.7) Implement the PCA analysis of the rescaled data matrix SDATA either in R or in Python, or in Matlab; explain clearly the inputs and the outputs of the pre-existing standard PCA functions you use

1.8) The standardized example #  $i$  is described by row " $i$ " of SDATA . After matrix transposition, this row becomes a column vector  $E_i$  in  $\mathbb{R}^{400}$  . Compute the first three "scores" of example " $i$ " by

$$\text{scor}_1(i) = \langle E_i, v_1 \rangle, \quad \text{scor}_2(i) = \langle E_i, v_2 \rangle, \quad \text{scor}_3(i) = \langle E_i, v_3 \rangle,$$

These 3 numbers are the coordinates of a 3-dimensional vector  $U_i$  in  $\mathbb{R}^3$ .

The 2 numbers  $\text{scor}_1(i)$ ,  $\text{scor}_2(i)$  define a 2-dimensional vector  $W_i$  in  $\mathbb{R}^3$ . Explain the geometric relationship between  $E_i$ ,  $U_i$ ,  $W_i$ .

1.9) Display graphically the 2 dimensional scatterplot of all the  $W_i$ ,  $i = 1, 2, \dots, N$ , using 3 colors (1 for each class); for instance red for CL1, blue for CL2 ; green for CL3. Interpret this display visually in terms of separability of the 3 classes

1.10) Display graphically the 3 dimensional scatterplot of all the  $U_i$ ,  $i = 1, 2, \dots, N$ , with the same 3 colors; Interpret visually. To facilitate the visual interpretation generate a similar display but with only the two classes CL1 and CL2. Repeat this operation and visual interpretation for CL1 and CL3, and then for CL2 and CL3

## PART 2

2.1) Fix  $k = 15$ . Use the standardized data matrix SDATA to apply the  $k$  nearest neighbor (kNN) algorithm for the automatic classification of arbitrary examples into one of the three classes CL1 CL2 CL3. Compute the percentage  $\text{per}(15)$  of correct classifications on the whole data set of  $N$  examples

2.2) Repeat the preceding operation for  $k = 5, 10, 15, 20, 30, 40, 50, 100, 200$ ; and compute the percentages  $\text{per}(k)$  of correct classifications on the whole data set of  $N$  examples

plot the curve  $\text{per}(k)$  versus  $k$  to try to identify a best range  $[A < k < B]$  of values for the integer  $k$

2.3) Repeat the preceding exploration for a few more values of  $k$  within the range  $[A, B]$ . Conclude by selecting a "best" value  $k^*$  for the integer  $k$

2.4) Compute and interpret the 3x3 confusion matrix for kNN classification using the "best"  $k = k^*$