



## Information Retrieval in High Dimensional Data

Lab #7, 7.12.2017

### Linear Discriminant Analysis

Task 1. In this task, we will once again work with the MNIST training set as provided on Moodle. Choose three digit classes, e.g. 1, 2 and 3 and load  $N=1000$  images from each of the classes to the workspace. Store the data in a normalized matrix  $X$  of type size  $(784, 3 \cdot N)$ . Furthermore, generate a color label matrix  $C$  of dimensions  $(3, 3 \cdot N)$ . Each row of  $C$  assigns an RGB color vector to the respective column of  $X$  as an indicator of the digit class. Choose  $[0, 0, 1]$ ,  $[0, 1, 0]$  and  $[1, 0, 0]$  for the three digit classes.

- Compute the principal subspace  $U$  of dimension 2 of  $X$ . Create a  $C$ -colored scatter plot of the scores of  $X$  with respect to this subspace.
- Generate a matrix  $X\_sums$  of size  $(784, 3)$  which consists of the row-wise sums of the centered data samples belonging to each of the three classes, normalized by the square roots of the class sizes. From  $X\_sums$ , compute the principal subspace  $U\_b$  of dimension 2 and generate a  $C$ -colored `scatter` plot of the scores of  $X$  with respect to  $U\_b$ . Compare the plot with the one from a). Which representation would you choose for a k-Nearest-Neighbors classification?
- Write a PYTHON function `sqrtminv` which expects a symmetric positive definite matrix  $A$  as its input and returns the inverse of its square root as its output, without using `scipy.linalg.sqrtm`.
- Divide  $X$  into three matrices of sizes  $(784, N)$ , each containing the samples belonging to one of each of the classes. Center the three matrices to create  $X1c$ ,  $X2c$  and  $X3c$  and compute the within-class-center matrix  $S\_w = (np.dot(X1c, X1c.T) + np.dot(X2c, X2c.T) + np.dot(X3c, X3c.T))$ .
- Calculate the "normalized class sum matrix"  $S\_bw = np.dot(sqrtminv(S\_w), X\_sums)$  and the left singular vectors  $U\_bw$  corresponding to its 2 largest singular values. Create a  $C$ -colored `scatter` plot of the scores of  $X$  with respect to the subspace described by  $sqrtminv(S\_w) * U\_bw$  (Hint: orthogonalize the basis vectors). Compare the plot with the ones from a) and b). Which representation would you choose for a k-Nearest-Neighbor classification? What weaknesses does this implementation of LDA have?

## Helpful Python/Numpy functions

<code>np.linalg.qr(A)</code>	Returns the QR orthogonalization of A
<code>np.sqrt(X)</code>	Square root