

COURSEWORK 4

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

Mathematics for Machine Learning

Author:

Arnold Cheung (CID: 01184493)

Date: December 2, 2019

Part I:

The file PCA.m, wPCA.m and LDA.m have been completed and used to produce the plots of the recognition error versus the number of components kept for each of the methods. It can be seen that LDA is the best performing method out of the three, achieving the lowest error rate (≈ 0.25) with the smallest number of components kept (≈ 50). This is due to the methods are tested for its classification ability, and LDA is an algorithm designed to maximise the separation between the classes within the total dataset, while the other two methods focuses on maximising the variance in the latent space. LDA can also be understood as a supervised learning technique, while the other two are unsupervised, which is more well trained in classification with labels provided as training data.

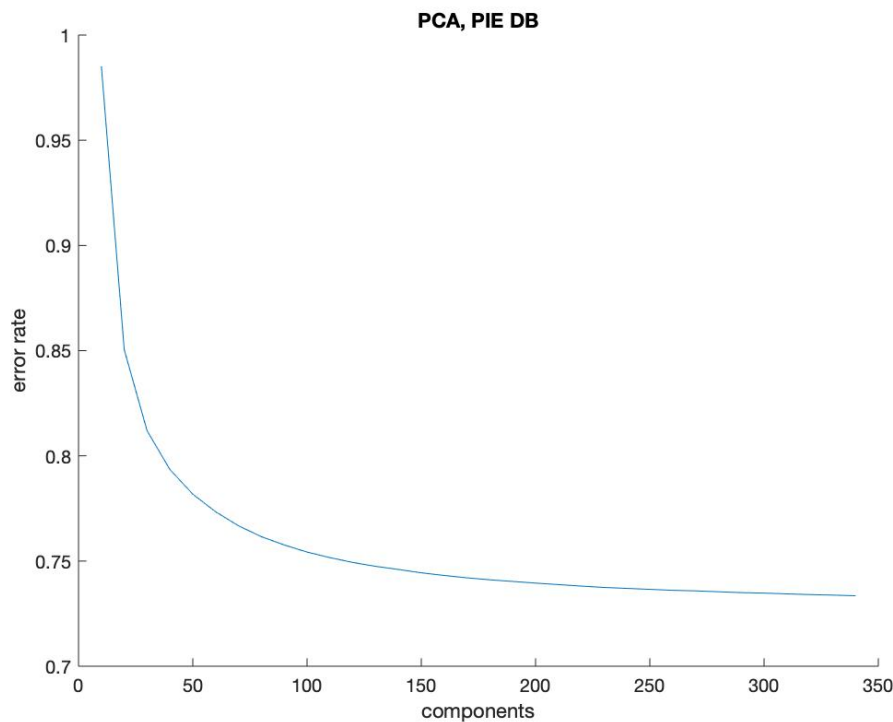


Figure 1: Recognition error versus the number of components kept for PCA

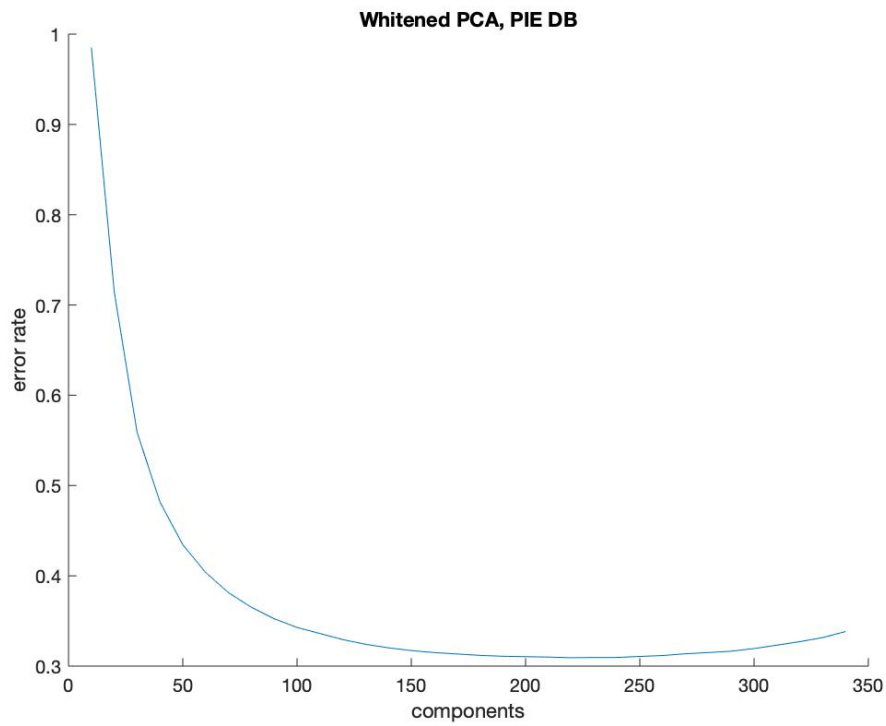


Figure 2: recognition error versus the number of components kept for wPCA

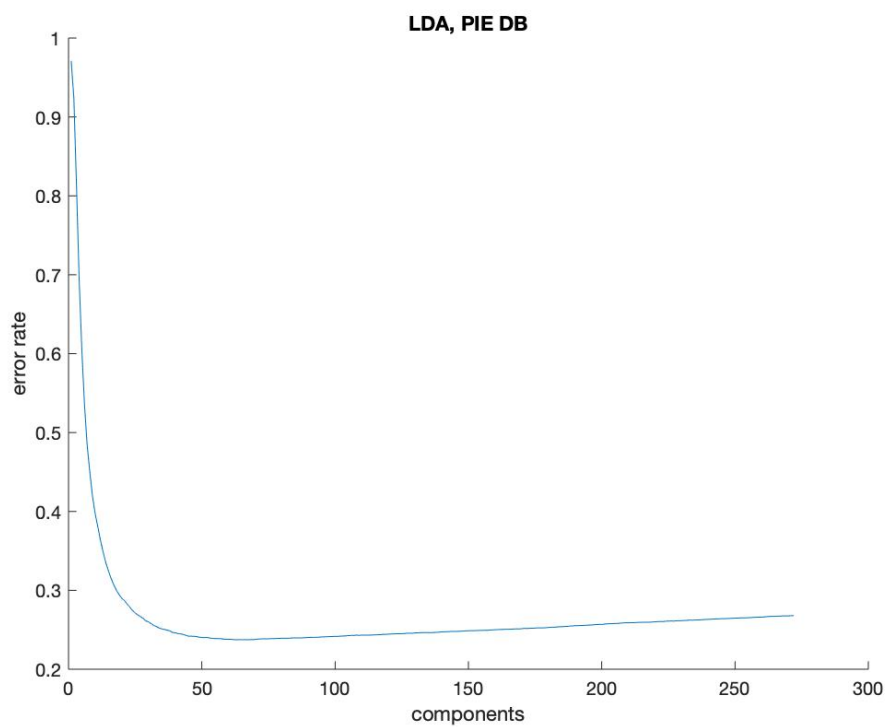


Figure 3: recognition error versus the number of components kept for LDA

Part II:

ia)

For the given optimisation problem, the Lagrangian is formulated as follows:

$$\mathcal{L}(w, b, \xi_i, \alpha, \beta) = \frac{1}{2} w^T S_t w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [y_i (w^T x_i + b) - 1 + \xi_i] - \sum_{i=1}^N \beta_i \xi_i \quad (1)$$

To minimise with respect to w, b and ξ , set the derivative with respect to each variable and set to 0, the optimal w, b and ξ is achieved when the following conditions are satisfied:

$$\begin{aligned} \nabla_w \mathcal{L} &= S_t w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \\ \Rightarrow w &= S_t^{-1} \sum_{i=1}^N \alpha_i y_i x_i \end{aligned} \quad (2)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial b} &= - \sum_{i=1}^N \alpha_i y_i = 0 \\ \Rightarrow \sum_{i=1}^N \alpha_i y_i &= 0 \end{aligned} \quad (3)$$

$$\begin{aligned} \nabla_{\xi} \mathcal{L} &= C - \alpha_i - \beta_i = 0 \\ \Rightarrow \beta_i &= C - \alpha_i \end{aligned} \quad (4)$$

To formulate the dual optimisation problem (maximise with respect to α), apply the above conditions to the original Lagrangian:

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T S_t^{-1} x_j + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i \xi_i - \sum_{i=1}^N \beta_i \xi_i \quad (5)$$

Collecting the terms:

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T S_t^{-1} x_j + \sum_{i=1}^N \alpha_i + \sum_{i=1}^N (C - \alpha_i - \beta_i) \xi_i \quad (6)$$

Since $\beta_i = C - \alpha_i$ (4), the last term becomes 0. Since β_i is also a Lagrangian multiplier and must be ≥ 0 , $C \geq \alpha_i$ can be concluded, and constraints on the dual problem is therefore formulated:

$$\begin{aligned} \max_{\alpha} & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T S_t^{-1} x_j + \sum_{i=1}^N \alpha_i \\ \text{s.t.} & \quad 0 \leq \alpha_i \leq C \\ & \quad \alpha^T y = 0 \end{aligned} \quad (7)$$

The problem can be simplified by changing the problem into a minimisation problem and defining $K_y = [y_i y_j x_i^T S_t^{-1} x_j]$:

$$\begin{aligned} \min_a \quad & \frac{1}{2} \alpha^T K_y \alpha + \mathbf{1}^T \alpha \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \\ & \alpha^T \mathbf{y} = 0 \end{aligned} \quad (8)$$

After computing α , the optimal w^* and b^* can be computed as follows:

$$w^* = S_t^{-1} \sum_{i=1}^N \alpha_i y_i x_i \quad (9)$$

α_i is only greater than 0 for points x_i that lies on the margin, i.e. $w^T x_i + b = y_i$, y_i is known and can only take values from -1, 1, using the optimal w^* , b^* can be recovered as follows:

$$\begin{aligned} b^* &= y_i - w^{*T} x_i \\ &\text{for } i \text{ where } \alpha_i > 0 \end{aligned} \quad (10)$$

The final b^* is the mean of all b retrieved to take into account numerical fluctuations.

ib)

The above optimisation problem has been solved numerically on using Matlab and the classification accuracy has been calculated to be 100% using the given train and test data. The code implementation and label prediction for each sample can be found in the file SVM.m.

ii)

S_t being singular can be dealt with by substituting $w = UQ$ and choose U such that $U^T S_t U = I$. Which transforms the problem from:

$$\min \frac{1}{2} w^T S_t w + C \sum_{i=1}^N \xi_i \quad (11)$$

into:

$$\min \frac{1}{2} Q^T Q + C \sum_{i=1}^N \xi_i \quad (12)$$

Which allowed the problem to be solved with respect to Q and avoiding the need to inverse S_t . The desired w can be retrieved in the end by computing $w = UQ$.

U can be chosen as follows:

$S_t = \frac{1}{N}XX^T$ is a symmetrical matrix and can therefore be diagonalised:

$$S_t = P^T \Lambda P \quad (13)$$

$$P^T = P^{-1} \quad (14)$$

$$U^T S_t U = U^T P \Lambda P^T U = I \quad (15)$$

Where P is the matrix whose columns are the eigenvectors of S_t and Λ is the diagonal matrix containing the eigenvalues of S_t . The required condition is satisfied when U is chosen as:

$$U = P \Lambda^{-\frac{1}{2}} \quad (16)$$

$$U^T P \Lambda P^T U = I$$

$$(P \Lambda^{-\frac{1}{2}})^T P \Lambda P^T P \Lambda^{-\frac{1}{2}} = I$$

$$\Lambda^{-\frac{1}{2}} P^T P \Lambda P^T P \Lambda^{-\frac{1}{2}} = I$$

$$\Lambda^{-\frac{1}{2}} \Lambda \Lambda^{-\frac{1}{2}} = I$$

$$I = I$$

Therefore, by choosing $U = P \Lambda^{-\frac{1}{2}}$ and substituting $w = UQ$, the problem caused by the singularity of S_t can be avoided.