

COURSEWORK

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

CO496 - Mathematics for Inference and Machine Learning

Author:

Daren Sin (CID: ds2912)

Date: November 30, 2016

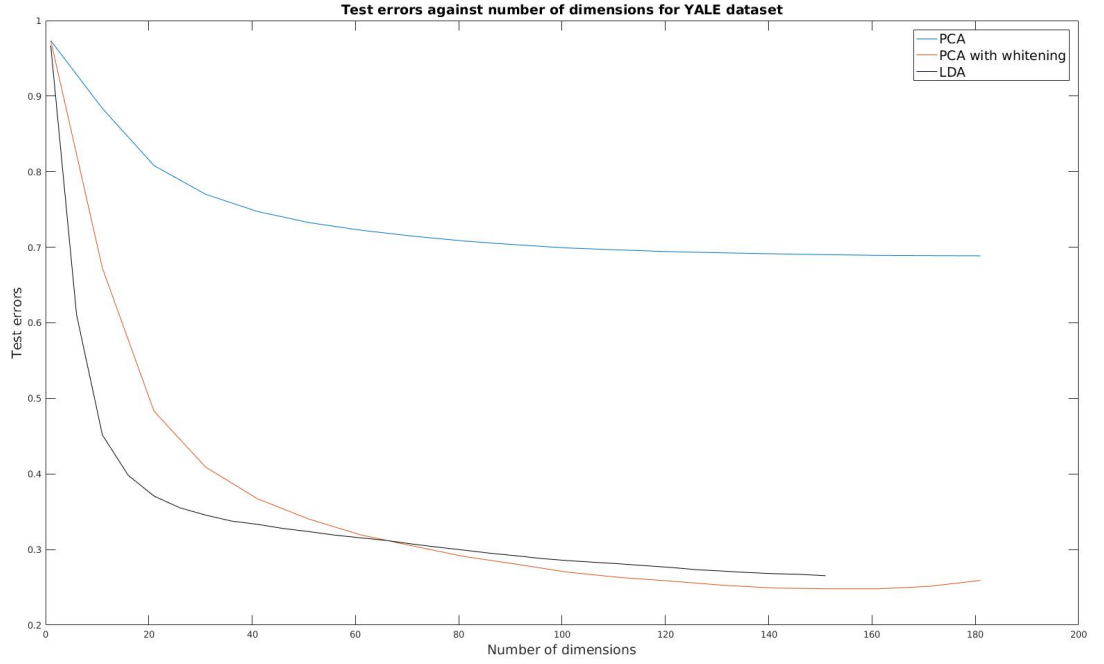


Figure 1: Graph of test errors against number of dimensions for the YALE dataset.

Part I

The plots for the YALE and PIE datasets can be found in Figure 1 and 2 respectively. First, for both datasets, LDA and PCA with whitening performed significantly better, with much lower test errors compared to the default PCA, across all dimensions.

Part II

1.) The Lagrangian is as follows:

$$L(R, \mathbf{a}, \xi_i, \lambda, r) = R^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \lambda_i \left[(\mathbf{x}_i - \mathbf{a})^\top (\mathbf{x}_i - \mathbf{a}) - R^2 - \xi_i \right] - \sum_{i=1}^n r_i \xi_i \quad (1)$$

To find the dual, we need to minimise the Lagrangian with respect to R , \mathbf{a} , and ξ_i . We thus differentiate the Lagrangian with respect to these variables:

$$\frac{\partial L}{\partial R} = 2R - 2R \sum_{i=1}^n \lambda_i \stackrel{!}{=} 0 \Rightarrow 1 - \sum_{i=1}^n \lambda_i = 0 \Rightarrow \sum_{i=1}^n \lambda_i = 1 \quad (2)$$

$$\frac{\partial L}{\partial \mathbf{a}} = \sum_{i=1}^n \lambda_i [-2(\mathbf{x}_i - \mathbf{a})] \stackrel{!}{=} 0 \Rightarrow \sum_{i=1}^n \lambda_i (\mathbf{x}_i - \mathbf{a}) = 0 \quad (3)$$

$$\frac{\partial L}{\partial \xi_i} = C - \lambda_i - r_i \stackrel{!}{=} 0 \Rightarrow C = \lambda_i + r_i \quad (4)$$

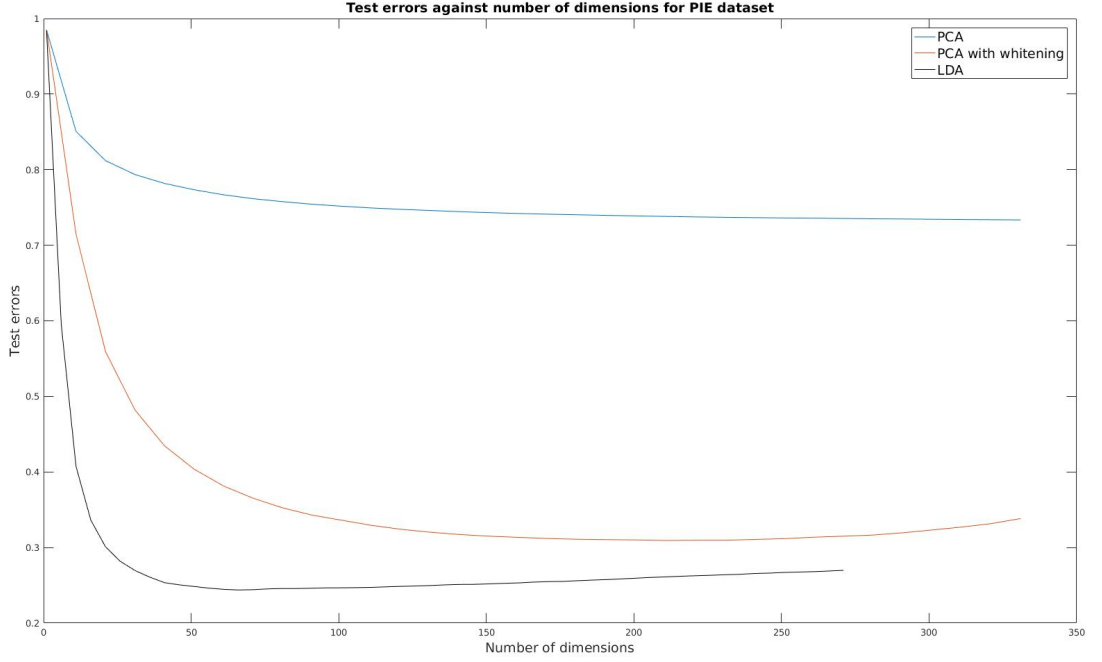


Figure 2: Graph of test errors against number of dimensions for the PIE dataset.

Furthermore, from equation 3, we also have the result:

$$\sum_{i=1}^n (\lambda_i \mathbf{x}_i - \mathbf{a} \lambda_i) = 0 \Rightarrow \sum_{i=1}^n (\lambda_i \mathbf{x}_i) - \mathbf{a} = 0 \Rightarrow \mathbf{a} = \sum_{i=1}^n \lambda_i \mathbf{x}_i \quad (5)$$

using the result that $\sum_{i=1}^n \lambda_i = 1$ from above. So, using the above equations, we can express the Lagrangian as such:

$$L(\lambda) = R^2 + \sum_{i=1}^n (\lambda_i + r_i) \xi_i + \lambda_i \sum_{i=1}^n (\mathbf{x}_i - \mathbf{a})^\top (\mathbf{x}_i - \mathbf{a}) - R^2 \sum_{i=1}^n \lambda_i - \sum_{i=1}^n \lambda_i \xi_i - \sum_{i=1}^n r_i \xi_i \quad (6)$$

$$= R^2 + \sum_{i=1}^n (\lambda_i + r_i) \xi_i + \sum_{i=1}^n [\lambda_i (\mathbf{x}_i - \mathbf{a})^\top (\mathbf{x}_i - \mathbf{a})] - R^2 - \sum_{i=1}^n (\lambda_i + r_i) \xi_i \quad (7)$$

$$= \sum_{i=1}^n [\lambda_i (\mathbf{x}_i^\top - \mathbf{a}^\top) (\mathbf{x}_i - \mathbf{a})] \quad (8)$$

$$= \sum_{i=1}^n \lambda_i [\mathbf{x}_i^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{a} - \mathbf{a}^\top \mathbf{x}_i + \mathbf{a}^\top \mathbf{a}] \quad (9)$$

$$= \sum_{i=1}^n [\lambda_i \mathbf{x}_i^\top \mathbf{x}_i - \lambda_i \mathbf{x}_i^\top \mathbf{a} - \lambda_i \mathbf{a}^\top \mathbf{x}_i + \lambda_i \mathbf{a}^\top \mathbf{a}] \quad (10)$$

$$= \sum_{i=1}^n [\lambda_i \mathbf{x}_i^\top \mathbf{x}_i - \lambda_i \mathbf{x}_i^\top \mathbf{a} - \lambda_i \mathbf{a}^\top (\mathbf{x}_i - \mathbf{a})] \quad (11)$$

$$= \sum_{i=1}^n [\lambda_i \mathbf{x}_i^\top \mathbf{x}_i - \lambda_i \mathbf{x}_i^\top \mathbf{a}] - \mathbf{a}^\top \underbrace{\sum_{i=1}^n \lambda_i (\mathbf{x}_i - \mathbf{a})}_{=0, \text{ from equation 3}} \quad (12)$$

We also note that, from equation 5, we have

$$\mathbf{a} = \sum_{j=1}^n \lambda_j \mathbf{x}_j$$

So,

$$L(\lambda) = \sum_{i=1}^n [\lambda_i \mathbf{x}_i^\top \mathbf{x}_i - \lambda_i \mathbf{x}_i^\top \mathbf{a}] \quad (13)$$

$$= \sum_{i=1}^n (\lambda_i \mathbf{x}_i^\top \mathbf{x}_i) - \sum_{i=1}^n \lambda_i \mathbf{x}_i^\top \left(\sum_{j=1}^n \lambda_j \mathbf{x}_j \right) \quad (14)$$

$$= \sum_{i=1}^n (\lambda_i \mathbf{x}_i^\top \mathbf{x}_i) - \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mathbf{x}_i^\top \mathbf{x}_j \lambda_j \quad (15)$$

$$= \text{diag}(\mathbf{X}^\top \mathbf{X})^\top \boldsymbol{\lambda} - \boldsymbol{\lambda}^\top (\mathbf{X}^\top \mathbf{X}) \boldsymbol{\lambda} \quad (16)$$

$$= \text{diag}(\mathbf{K}_x)^\top \boldsymbol{\lambda} - \boldsymbol{\lambda}^\top (\mathbf{K}_x) \boldsymbol{\lambda} \quad (17)$$

where $\mathbf{K}_x = [\mathbf{x}_i^\top \mathbf{x}_j]$. Hence, we can write the dual as:

$$\begin{aligned} & \max_{\boldsymbol{\lambda}} \quad \text{diag}(\mathbf{K}_x)^\top \boldsymbol{\lambda} - \boldsymbol{\lambda}^\top (\mathbf{K}_x) \boldsymbol{\lambda} \\ & \text{subject to} \quad \lambda_i \geq 0, \quad 0 \leq \lambda_i \leq C, \quad \text{for } i = 1, \dots, n \\ & \quad \quad \quad \sum_{i=1}^n \lambda_i = 1 \quad \Rightarrow \quad \mathbf{1}^\top \boldsymbol{\lambda} = 1 \end{aligned}$$

The second constraint, $0 \leq \lambda_i \leq C$, for $i = 1, \dots, n$ can be derived from equation 4, which implies that $\lambda_i = C - r_i$. Since $r_i \geq 0$ and $\lambda_i \geq 0$, r_i can only take values from 0 to C , inclusive. We then have the constraint, $0 \leq \lambda_i \leq C$.

Furthermore, we can write the above optimisation problem as its equivalent minimisation problem:

$$\begin{aligned} & -\min_{\boldsymbol{\lambda}} \quad -\text{diag}(\mathbf{K}_x)^\top \boldsymbol{\lambda} + \boldsymbol{\lambda}^\top (\mathbf{K}_x) \boldsymbol{\lambda} \\ & \text{subject to} \quad 0 \leq \lambda_i \leq C, \quad \text{for } i = 1, \dots, n \\ & \quad \quad \quad \mathbf{1}^\top \boldsymbol{\lambda} = 1 \end{aligned}$$

-
- 2.) When using arbitrary positive definite kernels, we can write the above optimisation problem as the following:

$$\begin{aligned} & -\min_{\lambda} \quad -\text{diag}(\mathbf{K}_x)^\top \lambda + \lambda^\top (\mathbf{K}_x) \lambda \\ & \text{subject to} \quad 0 \leq \lambda_i \leq C, \quad \text{for } i = 1, \dots, n \\ & \quad \quad \quad \mathbf{1}^\top \lambda = 1 \end{aligned}$$

where $\mathbf{K}_x = [k(\mathbf{x}_i, \mathbf{x}_j)]$ and $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ is the kernel.

- 3.) After identifying the parameters to the MATLAB function `quadprog` (see MATLAB code), we get an optimal λ^* per class, output by the function. We can then substitute the λ^* back into the objective function for the dual, to get, say, $-d^*$:

$$-d^* = -\text{diag}(\mathbf{K}_x)^\top \lambda^* + (\lambda^*)^\top (\mathbf{K}_x) \lambda^*$$

Furthermore, the optimal solution to the primal, p^* , is equal to the optimal solution to the dual, $-d^*$, where the negative sign arises from converting our maximisation problem to a minimisation one. We also note that, at the optimal solution, the slack variables, ξ_i , go to 0. Hence, keeping \mathbf{a} constant, p^* represents the optimal solution to the minimisation problem:

$$R^2 + C \sum_{i=1}^n \xi_i$$

As $\xi_i = 0$ for all $i = 1, \dots, n$, the optimal p^* is simply the optimal R^2 . We can thus find the radius of the optimal enclosing hypersphere by taking the square root of $-d^*$.

We can also find the vector $\mathbf{a}_k \in \mathbb{R}^2$, which represents the center of the optimal enclosing hypersphere of class k by:

$$\mathbf{a}_k = \sum_{j=1}^n \lambda_j^* \mathbf{x}_{j,k}$$

where $\mathbf{x}_{j,k} \in \mathbb{R}^2$ is the j th data point of class k . The plot of the optimal hyperspheres is in Figure 3.

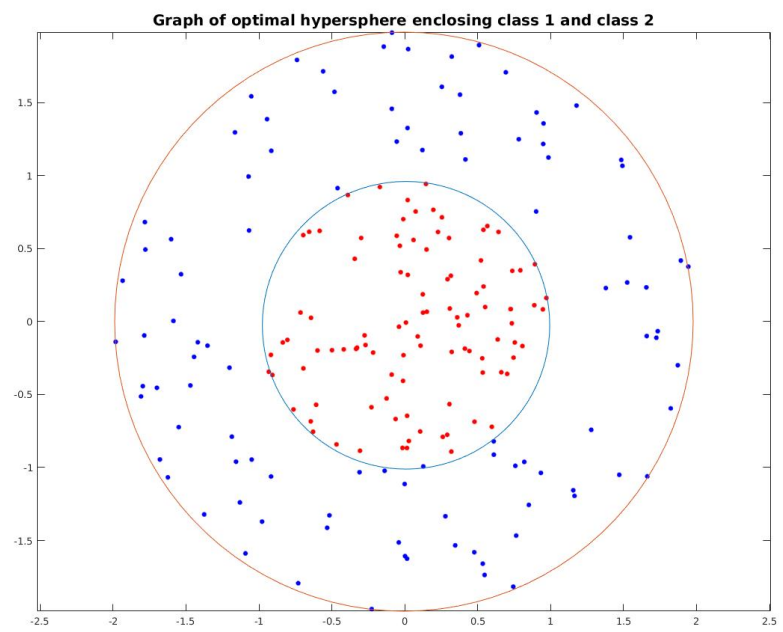


Figure 3: Plot of data points and the optimal enclosing hyperspheres for each class.