

CW4

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

CO496

Author:

Chengyi Zhang (CID: 01604256)

Date: November 29, 2018

Q1 From Fig. 1 and Fig. 2, the error rate of PCA is higher than whitening PCA.

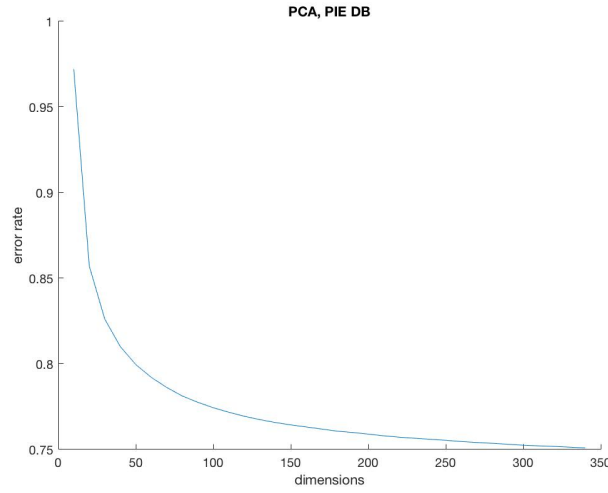


Figure 1: PCA

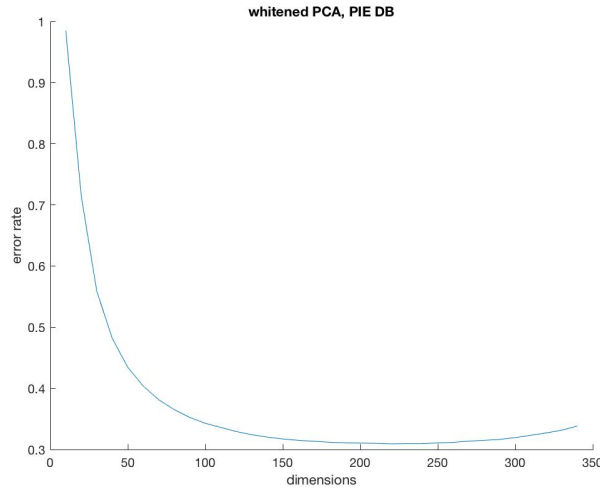


Figure 2: wPCA

When the dimension grows, both of them have similar change. The error rate decreases rapidly at start, and then the change gets slower as dimension is getting larger. For WPCA, the error rate even increases at the end. The rapid reduction at start is because the column eigenvectors in the base are placed corresponding to the eigenvalues sorted in descending order, which means that the eigenvectors in lower dimensions (top principle components) contain more information of the data than the ones in higher dimension. In WPCA, the data is normalised to have unit variance (the co-variance matrix is identity). Thus, the model with WPCA has better performance. Also, in WPCA, the smallest eigenvalues are close to zeros, and thus may cause numerical instability. The increase of error rate at the end of the curve WPCA may be because of the numerical instability by the small eigenvalues.

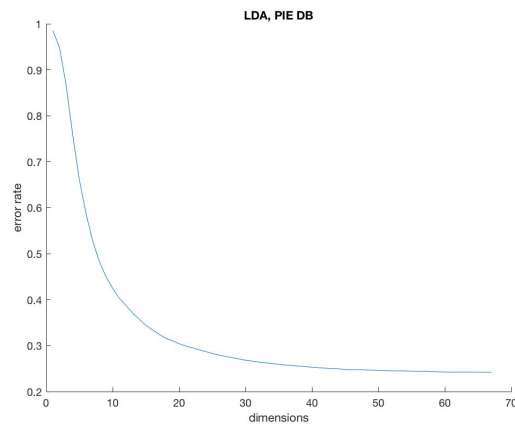


Figure 3: LDA

LDA has best performance among these three analysis algorithms. Compared to PCA, LDA is explicitly defined for classification problems, and is designed to define a latent space that helps classify data. Since PCA is also performed on projected data $\tilde{\mathbf{X}}_b = \mathbf{U}^\top \mathbf{X} \mathbf{M}$ in LDA, the columns at lower dimensions contain more information. And the performance curve shows similar pattern to that of PCA: rapidly decreases at start, and then becomes flat.

Q2

2.1

$$\min_{R, \alpha, \xi_i} R^2 + C \sum_{i=1}^n \xi_i \quad (1)$$

subject to $(\mathbf{x}_i - \alpha)^\top (\mathbf{x}_i - \alpha) \leq R^2 + \xi_i, \forall \xi_i \geq 0$.

Set $f(R, \alpha, \xi_i) = R^2 + C \sum_{i=1}^n \xi_i$, $g_i(R, \alpha, \xi_i) = (\mathbf{x}_i - \alpha)^\top (\mathbf{x}_i - \alpha) - R^2 - \xi_i$, and $h_i(R, \alpha, \xi_i) = -\xi_i$, then rewrite problem (1) to:

$$\begin{aligned} & \min_{R, \alpha, \xi_i} f(R, \alpha, \xi_i) \\ & \text{subject to } g_i(R, \alpha, \xi_i) \leq 0, \\ & h_i(R, \alpha, \xi_i) \leq 0, \forall i = 1 \dots n \end{aligned}$$

The Lagrangian is:

$$L(R, \alpha, \xi_i, \mu_i, \lambda_i) = f(R, \alpha, \xi_i) + \sum_{i=1}^n \mu_i g_i(R, \alpha, \xi_i) + \sum_{i=1}^n \lambda_i h_i(R, \alpha, \xi_i) \quad (2)$$

The primal is:

$$\min_{R, \alpha, \xi_i} \max_{\mu_i, \lambda_i \geq 0} L(R, \alpha, \xi_i, \mu_i, \lambda_i)$$

And the dual is:

$$\max_{\mu_i, \lambda_i \geq 0} \min_{R, \alpha, \xi_i} L(R, \alpha, \xi_i, \mu_i, \lambda_i)$$

Then find the partial derivatives:

$$\begin{aligned} \frac{\partial L}{\partial R} &= 2R - \sum_{i=1}^n \mu_i 2R = 0 \Rightarrow \sum_{i=1}^n \mu_i = 0 \\ \nabla_{\alpha} L &= -2 \sum_{i=1}^n \mu_i (\mathbf{x}_i - \alpha) = 0 \Rightarrow \sum_{i=1}^n \mu_i (\mathbf{x}_i - \alpha) = \sum_{i=1}^n \mu_i \mathbf{x}_i - \sum_{i=1}^n \mu_i \alpha = \sum_{i=1}^n \mu_i \mathbf{x}_i - \alpha = 0 \\ \frac{\partial L}{\partial \xi_i} &= C - \mu_i - \lambda_i = 0 \end{aligned}$$

Set $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, and then note that

$$\mathbf{x}_i - \alpha = \mathbf{x}_i - \sum_{j=1}^n \mu_j \mathbf{x}_j = \mathbf{x}_i - \mathbf{X} \boldsymbol{\mu}$$

Plug these values back to the dual problem:

$$\max_{\mu_i, \lambda_i \geq 0} \sum_{i=1}^n \mu_i (\mathbf{x}_i - \mathbf{X} \boldsymbol{\mu})^\top (\mathbf{x}_i - \mathbf{X} \boldsymbol{\mu})$$

which is equivalent to:

$$\begin{aligned}
& \min_{\mu_i \geq 0} \sum_{i=1}^n -\mu_i (\mathbf{x}_i - \mathbf{X}\boldsymbol{\mu})^\top (\mathbf{x}_i - \mathbf{X}\boldsymbol{\mu}) \\
& \text{s.t } \mu_i \leq C, \sum_{i=1}^n \mu_i = 1 \\
\\
& \sum_{i=1}^n \mu_i (\mathbf{x}_i - \mathbf{X}\boldsymbol{\mu})^\top (\mathbf{x}_i - \mathbf{X}\boldsymbol{\mu}) = \sum_{i=1}^n \mu_i (\mathbf{x}_i^\top \mathbf{x}_i - 2\mathbf{x}_i^\top \mathbf{X}^\top \boldsymbol{\mu} + \boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu}) \\
& = \sum_{i=1}^n \mu_i \mathbf{x}_i^\top \mathbf{x}_i - 2 \sum_{i=1}^n \mu_i \mathbf{x}_i^\top \mathbf{X}^\top \boldsymbol{\mu} + \sum_{i=1}^n \mu_i \boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} \\
& = \sum_{i=1}^n \mu_i \mathbf{x}_i^\top \mathbf{x}_i - 2\boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} \sum_{i=1}^n \mu_i \\
& = \sum_{i=1}^n \mu_i \mathbf{x}_i^\top \mathbf{x}_i - \boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu}
\end{aligned}$$

The problem is now

$$\begin{aligned}
& \min_{\boldsymbol{\mu}} \sum_{i=1}^n -\mu_i \mathbf{x}_i^\top \mathbf{x}_i + \boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} \\
& \text{s.t } 0 \leq \mu_i \leq C, \sum_{i=1}^n \mu_i = 1
\end{aligned} \tag{3}$$

Transform to the formula for quadprog():

$$\begin{aligned}
& \min_{\mathbf{g}} \mathbf{f}^\top \mathbf{g} + \frac{1}{2} \mathbf{g}^\top \mathbf{H} \mathbf{g} \\
& \text{s.t } \mathbf{A} \mathbf{g} \leq \mathbf{c}, \mathbf{A}_e \mathbf{g} = \mathbf{c}_e, \mathbf{g}_l \leq \mathbf{g} \leq \mathbf{g}_u
\end{aligned}$$

where $\mathbf{g} = \boldsymbol{\mu}$, $\mathbf{f} = -\text{diag}(\mathbf{X}^\top \mathbf{X})$, $\mathbf{H} = 2\mathbf{X}^\top \mathbf{X}$, $\mathbf{A} = \mathbf{0}$, $\mathbf{c} = \mathbf{0}$, $\mathbf{A}_e = \mathbf{1}^\top$, $\mathbf{g}_l = \mathbf{0}$, $\mathbf{g}_u = \mathbf{1} \cdot C$

Plug back to Lagrangian (2), and get

$$\boldsymbol{\alpha} = \sum_{i=1}^n \mu_i \mathbf{x}_i$$

Now we need to find optimal R and ξ_i 's. Note that when R^2 is larger than the square of distance from the $\lceil \frac{1}{C} \rceil$ -th farthest data point to $\boldsymbol{\alpha}$, the decrease of R^2 is more significant than the increase of $C \sum_{i=1}^n \xi_i$ (one unit decrease in R^2 causes $C \cdot \lfloor \frac{1}{C} \rfloor$ unit increase in $C \sum_{i=1}^n \xi_i$), and thus

$$\begin{aligned}
& SV = \{\mathbf{x}_i' \text{ s with largest } \lceil 1/C \rceil \mu_i' \text{ s}\}, \text{ the set of support vectors} \\
& R = \sqrt{\min_{\mathbf{x}_i \in SV} (\mathbf{x}_i - \boldsymbol{\alpha})^\top (\mathbf{x}_i - \boldsymbol{\alpha})}
\end{aligned}$$

And the centre and radius are found.

2.2

Similar to **2.1**. Replace \mathbf{x}_i with $\phi(\mathbf{x}_i)$, and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ with $\mathbf{\Phi} = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]$ in (3), then set $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ (note that \mathbf{K} is positive definite since kernel k is positive definite), and get

$$\min_{\boldsymbol{\mu}} \sum_{i=1}^n -\mu_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_i) + \boldsymbol{\mu}^\top \mathbf{\Phi}^\top \mathbf{\Phi} \boldsymbol{\mu} \quad (4)$$

$$= \min_{\boldsymbol{\mu}} \sum_{i=1}^n -\mu_i k(\mathbf{x}_i, \mathbf{x}_i) + \boldsymbol{\mu}^\top \mathbf{K} \boldsymbol{\mu} \quad (5)$$

$$\text{s.t } 0 \leq \mu_i \leq C, \sum_{i=1}^n \mu_i = 1$$

Transform to the formula for quadprog():

$$\begin{aligned} & \min_{\mathbf{g}} \mathbf{f}^\top \mathbf{g} + \frac{1}{2} \mathbf{g}^\top \mathbf{H} \mathbf{g} \\ & \text{s.t } \mathbf{A} \mathbf{g} \leq \mathbf{c}, \mathbf{A}_e \mathbf{g} = \mathbf{c}_e, \mathbf{g}_l \leq \mathbf{g} \leq \mathbf{g}_u \\ & \text{where } \mathbf{g} = \boldsymbol{\mu}, \mathbf{f} = -\text{diag}(\mathbf{K}), \mathbf{H} = 2\mathbf{K}, \mathbf{A} = \mathbf{0}, \mathbf{c} = \mathbf{0}, \mathbf{A}_e = \mathbf{1}^\top, \mathbf{g}_l = \mathbf{0}, \mathbf{g}_u = \mathbf{1} \cdot C \end{aligned}$$

The rest is similar to **2.1**.

$$\begin{aligned} \boldsymbol{\alpha} &= \sum_{i=1}^n \mu_i \phi(\mathbf{x}_i) \\ \text{SV} &= \{\mathbf{x}'_i \text{ with largest } \lceil 1/C \rceil \mu'_i \text{'s}, \text{ the set of support vectors} \} \\ R &= \sqrt{\min_{\mathbf{x}_i \in \text{SV}} (\phi(\mathbf{x}_i) - \boldsymbol{\alpha})^\top (\phi(\mathbf{x}_i) - \boldsymbol{\alpha})} \end{aligned}$$

2.3

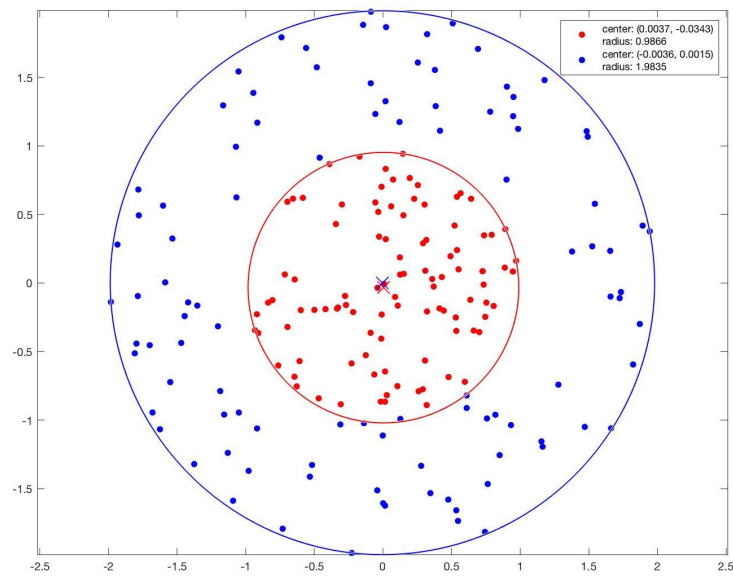


Figure 4: SVM, $C=1$

With $C = 1$ (no allowance of error):

Class 1: Centred at (0.0037, -0.0343) with radius 0.9866

Class 2: Centred at (-0.0036, 0.0015) with radius 1.9835

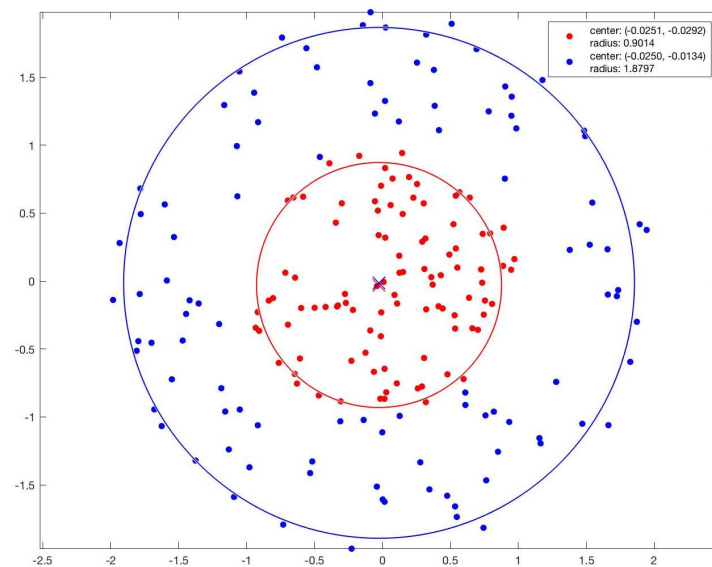


Figure 5: SVM, $C=0.05$

With $C = 0.05$:

Class 1: Centred at (-0.0251, -0.0292) with radius 0.9014

Class 2: Centred at (-0.0250, -0.0134) with radius 1.8797