

### Exercise Sheet 3. Machine learning

1. (a) find  $\theta$  that minimizes  $J(\theta) = \sum_{k=1}^n \|\theta - x_k\|^2$ , subject to  $\theta^T b = 0$

by Lagrange multiplier.

$$\begin{aligned} \text{let } L &= J(\theta) + \lambda \theta^T b \\ &= \sum_{k=1}^n \|\theta - x_k\|^2 + \lambda \theta^T b \\ &= \sum_{k=1}^n (\theta \theta^T - 2\theta^T x_k + x_k^T x_k) + \lambda \theta^T b \end{aligned}$$

$$\frac{\partial L}{\partial \lambda} = \theta^T b = 0 \quad \text{--- ①}$$

$$\frac{\partial L}{\partial \theta} = \sum_{k=1}^n (2\theta - 2x_k) + \lambda b = 0 \quad \text{--- ②}$$

$$\text{multiply ② by } \theta^T: \sum_{k=1}^n (2\theta^T \theta - 2\theta^T x_k) + \lambda \theta^T b = 0$$

$$\text{by ①: } \theta^T b = 0.$$

$$\text{so: } \sum_{k=1}^n (2\theta^T \theta - 2\theta^T x_k) = 0 \Rightarrow \theta^T \theta = \theta^T \cdot \frac{1}{n} \sum_{k=1}^n x_k$$

In order to make sure  $J(\theta)$  gets its minimization at  $\theta$ , we need to make sure the second derivation  $\frac{\partial^2 L}{\partial \theta^2} > 0$ .

$$\begin{aligned} \frac{\partial^2 L}{\partial \theta^2} &= \frac{\partial}{\partial \theta} \sum_{k=1}^n (2\theta - 2x_k) + \lambda b \\ &= \sum_{k=1}^n 2 = 2n > 0. \end{aligned}$$

so: the parameter  $\theta$  that minimizes  $J(\theta)$  should be any parameter  $\in \mathbb{R}^d$  that satisfies:  $\begin{cases} \theta^T b = 0 \\ \theta^T \theta = \theta^T \cdot \frac{1}{n} \sum_{k=1}^n x_k \end{cases} \quad //$

(b) find  $\theta$  that minimizes  $J(\theta) = \sum_{k=1}^n \|\theta - x_k\|^2$  subject to  $\|\theta - c\|^2 = 1$ .

$$\begin{aligned} \text{let } L &= J(\theta) + \lambda (\|\theta - c\|^2 - 1) \\ &= \sum_{k=1}^n (\theta \theta^T - 2\theta^T x_k + x_k^T x_k) + \lambda (\theta^T \theta - 2\theta^T c + c^T c) - \lambda \end{aligned}$$

$$\frac{\partial L}{\partial \lambda} = \|\theta - c\|^2 - 1 = (\theta - c)^T (\theta - c) - 1 = 0 \quad \text{--- ①}$$

$$\frac{\partial L}{\partial \theta} = \sum_{k=1}^n (2\theta - 2x_k) + \lambda (2\theta - 2c) = 0 \quad \text{--- ②}$$

$$\text{multiply ② by } (\theta - c)^T: \sum_{k=1}^n (2\theta - 2x_k)(\theta - c)^T + \lambda (2\theta - 2c)(\theta - c)^T = 0$$

$$\text{by ①: } (\theta - c)^T (\theta - c) = 1.$$

$$\text{so: } \sum_{k=1}^n (2\theta - 2x_k)(\theta - c)^T + 2\lambda = 0 \Rightarrow (n\theta - \sum_{k=1}^n x_k)(\theta - c)^T + \lambda = 0$$

Similar, we need to make sure  $\frac{\partial^2 L}{\partial \theta^2} > 0$ :

$$\frac{\partial^2 L}{\partial \theta^2} = \sum_{k=1}^n 2 + \lambda (2\theta - 2c) = 2n + 2\lambda > 0 \quad n > \lambda$$

So: the parameter  $\theta$  that minimizes  $J(\theta)$  subject to  $\|\theta - c\|^2 = 1$  should be any parameter that satisfies: 
$$\begin{cases} \|\theta - c\|^2 = 1 \\ (\nabla_{\theta} - \frac{\partial}{\partial \theta} \lambda R) (\theta - c)^T + \lambda = 0 \\ \lambda < n \end{cases} //$$

2. (a) by the definition of Trace:  $\text{Tr}(S) = \sum_{i=1}^d S_{ii}$

and:  $\text{Tr}(S) = \sum_{i=1}^m \lambda_i$ , where  $\lambda_i = \text{eig}(S)$

as for all  $i$  ( $i \leq m$ ):  $\lambda_i \geq 0$

$$\text{so: } \sum_{i=1}^d S_{ii} = \text{Tr}(S) = \sum_{i=1}^m \lambda_i = \lambda_1 + \lambda_2 + \dots + \lambda_m.$$

$$\text{so: } \lambda_1 \leq \sum_{i=1}^d S_{ii}. //$$

(b) when "=" holds, ie:  $\lambda_1 = \sum_{i=1}^d S_{ii}$ .

then:  $m=1$ . which means  $S$  has only one eigenvalue.

The whole database  $x_i \in \mathbb{R}^d$ ,  $x_i$  can be considered as lying in one-dimension, ie: even though  $x_i \in \mathbb{R}^d$ , they can be transformed into another space which has dimension 1. ( $x_i \in \mathbb{R}^1$ ). //

(c). The reason that  $\max_{i=1}^d S_{ii}$  is a lower bound of  $\lambda_1$  is:

$S_{11}, S_{22}, \dots, S_{dd}$  can be viewed as the variance of original data  $x$ , which measure the component of each data  $x_i$  in  $x$ .

$\lambda_1, \lambda_2, \dots, \lambda_m$  is the eigenvalue of  $S$ , by the idea of PCA,  $\lambda$  is also the component of the reshaped data in  $x$ .

PCA extracts the most principle component of a database, so  $\lambda_1$ , which is the largest eigenvalue, should be larger than any  $S_{ii}$ . ( $i \leq d$ )

$$\text{ie: } \lambda_1 \geq \max_{i=1}^d S_{ii}. //$$

(d) when "=" holds: ie:  $\lambda_1 = \max_{i=1}^d S_{ii}$ .

The reshaped data by PCA is of no difference to the original data. Or the data can not be transformed to another space  $\mathbb{R}^t$ , where  $t < d$ . //

3. (a). by  $V = S^{0.5} W$ :  $V^T V = (S^{0.5} W)^T (S^{0.5} W) = W^T S W$ .

since  $S$  is invertible:  $W = S^{-0.5} V = S^{-0.5} S^{0.5} W \Rightarrow W = S^{-0.5} V$ .

$$\text{so: } J(W) = \|S W\|^2 - \frac{1}{2} W^T S W$$

$$= \|S \cdot S^{-0.5} V\|^2 - \frac{1}{2} W^T S W$$

$$= \|S^{0.5} V\|^2 - \frac{1}{2} V^T V$$

$$\frac{\partial J}{\partial V} = \|S^{0.5}\|^2 - \frac{1}{2} V.$$

~~100-205~~

$$\text{so: } V \leftarrow V + y \frac{\partial J}{\partial V}.$$

$$V \leftarrow V + y(\|S^{0.5}\| - \frac{1}{2}V). \quad V \leftarrow (1 + \frac{y\|S^{0.5}\|}{\phi} V^{-1} - \frac{1}{2}y) V.$$

$$\text{by: } W \leftarrow \frac{S W}{\|S W\|} \text{ and } W = S^{-0.5} V:$$

$$S^{-0.5} V \leftarrow \frac{S \cdot S^{-0.5} V}{\|S \cdot S^{-0.5} V\|}$$

$$S^{0.5} S^{-0.5} V \leftarrow \frac{S^{0.5} S \cdot S^{-0.5} V}{\|S^{0.5} V\|}$$

$$\text{so: } V \leftarrow \frac{S V}{\|S^{0.5} V\|}$$

$$\text{i.e.: } 1 + y\|S^{0.5}\| V^{-1} - \frac{1}{2}y \rightarrow \frac{S}{\|S^{0.5} V\|} \quad \text{--- (3)}$$

But here we cannot find a proper way to prove (3).



(b)